

Numerical algorithms in control

Zlatko Drmač

University of Zagreb
Department of Mathematics

Trogir 2011

Outline

- 1 Introduction
- 2 Scaling: examples
- 3 Numerical rank revealing
- 4 Eigenvalues and singular values
- 5 Jacobi method
- 6 Accurate PSVD and applications
- 7 Concluding remarks

- Problems
- Machine numbers
- Examples
- Consequences
- Goals

6 Accurate PSVD and applications

LTI systems, control, tasks

Introduction

Problems

Machine numbers

Examples

Consequences

Goals

Scaling:
examplesNumerical
rank revealingEigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

Space station, CD player, vehicle suspension system, ...

$$\begin{aligned} E\dot{x}(t) &= Ax(t) + Bu(t), \quad E, A \in \mathbb{R}^{n \times n}, \quad B \in \mathbb{R}^{n \times m} \\ y(t) &= Cx(t) + Du(t), \quad C \in \mathbb{R}^{p \times n}, \quad D \in \mathbb{R}^{p \times m}. \end{aligned} \quad (1)$$

Example: $\dot{x} = Ax + Bu + Gr$ with $x \in \mathbb{R}^6$,

$$A = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ -\frac{k_p}{m_p} & -\frac{c_p}{m_p} & \frac{k_p}{m_p} & \frac{c_p}{m_p} & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ \frac{k_p}{m_s} & \frac{c_p}{m_s} & -\frac{k_s+k_p}{m_s} & -\frac{c_s+c_p}{m_s} & \frac{k_s}{m_s} & \frac{c_s}{m_s} \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & \frac{k_s}{m_{us}} & \frac{c_s}{m_{us}} & -\frac{k_s+k_t}{m_{us}} & -\frac{c_s}{m_{us}} \end{pmatrix}, \quad B = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ -1 \end{pmatrix}$$

$G = (0, 0, 0, 0, 0, k_t/m_{us})$; $r(t)$ = road; $u(t)$ = actuator force; $x_1(t)$ = passenger's vertical displacement. Determine $u(t)$ (e.g. $u(t) = -Kx(t)$) to ensure smooth riding on a rough road.

Control, introduction, tasks

$$\begin{aligned}\dot{x}(t) &= Ax(t) + Bu(t), \quad x(0) = 0; \\ y(t) &= Cx(t) + Du(t).\end{aligned}\tag{2}$$

Apply Laplace transform to get

$$\hat{y}(s) = \underbrace{(C(sI - A)^{-1}B + D)}_{G(s) \equiv \text{transfer function}} \hat{u}(s)$$

Of interest is the input \rightarrow output behavior $\hat{y}(s) = G(s)\hat{u}(s)$.
In large scale/real time applications: try to reproduce nearly the same behavior with a system of smaller dimension $r \ll n$. Take $D = 0$.

$$\left. \begin{aligned}\dot{x}_r(t) &= A_r x_r(t) + B_r u(t) \\ y_r(t) &= C_r x_r(t)\end{aligned} \right\} \quad G_r(s) = C_r(sI - A_r)^{-1} B_r$$

$\hat{y}(s) - \hat{y}_r(s) = (G(s) - G_r(s))\hat{u}(s)$ should be small in some norm for a class of inputs $u(\cdot)$.

NLA tasks in control

Consider n dimensional LTI SISO (more general, XIXO)

$$\begin{aligned}\dot{x}(t) &= A x(t) + b u(t) & \bowtie G(s) &= c(sI - A)^{-1} b. \\ y(t) &= c x(t)\end{aligned}$$

For $r < n$ and r -dimensional $\mathcal{V}_r = \mathcal{R}(V_r)$, $\mathcal{W}_r = \mathcal{R}(W_r)$ with $\mathcal{V}_r \cap \mathcal{W}^\perp = \{0\}$ ($\Leftrightarrow \det(W_r^T V_r) \neq 0$) look for

$$\mathcal{V}_r \ni v(t) = V_r x_r(t) \text{ such that } \dot{v}(t) - A v(t) - b u(t) \perp \mathcal{W}_r.$$

The reduced output is $y_r(t) = c v(t)$. In the bases V_r , W_r ,

$$W_r^T (V_r \dot{x}_r(t) - A V_r x_r(t) - b u(t)) = 0, \text{ i.e.}$$

$$\begin{aligned}\dot{x}_r(t) &= A_r x_r(t) + b_r u(t) & \bowtie G_r(s) &= c_r(sI - A_r)^{-1} b_r \\ y_r(t) &= c_r x_r(t)\end{aligned}$$

$$A_r = (W_r^T V_r)^{-1} W_r^T A V_r, \quad b_r = (W_r^T V_r)^{-1} W_r^T b, \quad c_r = c V_r.$$

Introduction

Problems

Machine numbers

Examples

Consequences

Goals

Scaling:
examplesNumerical
rank revealingEigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

Numerical tasks

Introduction

Problems

Machine numbers

Examples

Consequences

Goals

Scaling:
examplesNumerical
rank revealingEigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

- $\|G\|_{\mathcal{H}_2} = \sqrt{\frac{1}{2\pi} \int_{-\infty}^{\infty} |G(i\omega)|^2 d\omega}$;
- $\min \|G - G_r\|_{\mathcal{H}_2}$, G_r stable of order r
- Let G_r be a local minimizer with simple poles at $\tilde{\lambda}_i$, $i = 1, \dots, r$. Then at $\sigma_i = -\tilde{\lambda}_i$: $G_r(\sigma_i) = G(\sigma_i)$, $G'_r(\sigma_i) = G'(\sigma_i)$, $i = 1, \dots, r$.
- Hermite interpolation by $\mathcal{V}_r = \text{Span}((\sigma_i I - A)^{-1} b)_{i=1}^r$, $\mathcal{W}_r = \text{Span}((\sigma_i I - A^T)^{-1} c^T)_{i=1}^r$.
- Solving linear systems for V_r and W_r . Reduce to generalized upper Hessenberg form

$$Q^T E Z = \begin{pmatrix} \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ & \blacksquare & \blacksquare & \blacksquare \\ & & \blacksquare & \blacksquare \\ & & & \blacksquare \end{pmatrix}, \quad Q^T A Z = \begin{pmatrix} \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ & \blacksquare & \blacksquare & \blacksquare \\ & & \blacksquare & \blacksquare \\ & & & \blacksquare \end{pmatrix}, \quad Q^T b = \begin{pmatrix} \blacksquare \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

and work on $(E, A, b, c) \equiv (Q^T E Z, Q^T A Z, Q^T b, c Z)$ is efficient. Simpler if $E = I$, $Z = Q$.

Numerical tasks

Introduction

Problems

Machine numbers

Examples

Consequences

Goals

Scaling:
examplesNumerical
rank revealingEigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

Generate many interesting and challenging problems.

- Simple questions, difficult answers: Compute the transfer function $G(\zeta) = C(\zeta E - A)^{-1}B$ for many complex values of ζ . Here n can be large.
- By changing the state space coordinates, $x(t) = T\hat{x}(t)$, the new representation is, e.g. for $E = I$, given with $(\hat{A}, \hat{B}, \hat{C}, \hat{D}) = (T^{-1}AT, T^{-1}B, CT, D)$. Find T such that the new representation reveals structural properties of the system. Various canonical forms.
- Solve Lyapunov equation $AH + HA^T + BB^T = 0$. Solve Riccati eqn: $XA + A^TX + Q - XSX = 0$. Many other types of matrix equations.
- Find invariant subspace that corresponds to specified eigenvalues.

... algorithms, software

Introduction

Problems

Machine numbers

Examples

Consequences

Goals

Scaling: examples

Numerical rank revealing

Eigenvalues and singular values

Jacobi method

Accurate PSVD and applications

Concluding remarks

- Solve eigenvalue and singular value problems.
- Given A with eigenvalues $\lambda_1, \dots, \lambda_n$ and B , find K such that $A - BK$ has prescribed eigenvalues $\alpha_1, \dots, \alpha_n$.
- Pressure from applications to deliver accurate solutions quickly. Computing environments changing rapidly.
- Users from applied sciences and engineering – usually not interested in math details, just solutions, software.
- Pure mathematicians not interested because the problems are "trivial", non–fundamental or just too messy.
- And we have high performance computers. So, why is this difficult?

Yes, have computer, but ..

Machine (floating-point) numbers $\mathbb{F} \subset \mathbb{Q}$.

$$f = \pm m \cdot 2^e, \quad e = -126 : 127, \quad m = 1.z_1 \dots z_{23}.$$

$$\overline{\mathbb{F}} = \mathbb{F} \cup \{+\text{Infinity}, -\text{Infinity}, \text{NaN}\}$$

Machine arithmetic $\oplus, \ominus, \odot, \oslash$.

- $\overline{\mathbb{F}}$ finite, 2^{32} (single), 2^{64} (double); $0.1 \notin \mathbb{F}$;
- $a \oplus b \equiv \mathbf{FL}(a + b) = (a + b)(1 + \epsilon_{a,b})$,

$$|\epsilon_{a,b}| \leq \mathbf{u} \equiv \text{eps} = \mathbf{round-off} \approx 10^{-8}.$$

- In general, $(a \oplus b) \oplus c \neq a \oplus (b \oplus c)$,
 $(a \odot b) \odot c \neq a \odot (b \odot c)$; $x \oplus y \oplus z = ??$
- $1 \oplus 10^{-9} = 1$; $x = y \not\Rightarrow x - y = 0$; $10^{-30} \odot 10^{-30} = 0$;
- Finite speed, finite memory.
- Faster \implies more mess per second.

Example using MATLAB,

$$\text{eps} \approx 2.2 \cdot 10^{-16}$$

$$X = (x \ y) \in \mathbf{R}^{m \times 2}, \quad \begin{pmatrix} a & c \\ c & b \end{pmatrix} = \text{computed}(X^T X),$$

Let $x = 5 \cdot 10^{153} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$, $y = 10 \begin{pmatrix} 1 \\ 2 \end{pmatrix}$. ($\cos \angle(x, y) = \frac{3}{\sqrt{10}}$).

Test the orthogonality of x and y ,

$$\cos \angle(x, y) \equiv \frac{c}{\sqrt{ab}} \leq \epsilon$$

$$(c / \text{sqrt}(a*b) <= \text{eps}) = 1,$$

$$((c / \text{sqrt}(a)) / \text{sqrt}(b) <= \text{eps}) = 0,$$

$$(c <= \text{sqrt}(a*b) * \text{eps}) = 1,$$

$$(c <= \text{sqrt}(a)*\text{sqrt}(b) * \text{eps}) = 0.$$

Example using MATLAB,

$$\text{eps} \approx 2.2 \cdot 10^{-16}$$

$$\text{Let } \mathbf{x} = 5 \cdot 10^{-153} \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad \mathbf{y} = 10^{-16} \begin{pmatrix} 1 \\ -1 \end{pmatrix}$$

Then

$$(\text{c} / \text{sqrt}(\mathbf{a} * \mathbf{b}) \leq \text{eps}) = 0,$$

$$((\text{c} / \text{sqrt}(\mathbf{a})) / \text{sqrt}(\mathbf{b}) \leq \text{eps}) = 1,$$

$$(\text{c} \leq \text{sqrt}(\mathbf{a} * \mathbf{b}) * \text{eps}) = 0,$$

$$(\text{c} \leq \text{sqrt}(\mathbf{a}) * \text{sqrt}(\mathbf{b}) * \text{eps}) = 1.$$

Another example

$$A = \begin{pmatrix} 1 & 1 & \boxed{1} \\ 0 & 1 & \xi \\ 0 & -1 & \xi \end{pmatrix}, \text{ where } \xi = 10/\epsilon_{\text{ps}}. \xi \approx 4.5\text{e}+016$$

$$\text{Givens rotation kills } A_{13}: \tilde{A}^{(1)} = \begin{pmatrix} 1 & \alpha & 0 \\ 0 & \beta & \beta \\ 0 & \boxed{\beta} & \beta \end{pmatrix};$$

$$\alpha \approx \sqrt{2}, \beta = 3.184525836262886\text{e}+016.$$

	$\text{svd}(A)$	$\text{svd}(A^T)$
σ_1	6.369051672525773e+16	6.369051672525772e+16
σ_2	5.747279316501105e+00	3.004066501831585e+00
σ_3	9.842664568695829e-01	4.220776043599739e-01

$$\tilde{A}^{(1)} = \begin{pmatrix} 1 & \alpha & 0 \\ 0 & \beta & \beta \\ 0 & \boxed{\beta} & \beta \end{pmatrix}, \quad \tilde{A}^{(2)} = \begin{pmatrix} 1 & \alpha & 0 \\ 0 & \gamma & \gamma \\ 0 & 0 & 0 \end{pmatrix},$$

$$A = BD, \kappa_2(B) < 2.$$

A 2×2 example

Introduction

Problems

Machine numbers

Examples

Consequences

Goals

Scaling:
examplesNumerical
rank revealingEigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

Take in MATLAB

$$A = \begin{pmatrix} 1.0e250 & 0 \\ 0 & 1.0e-201 \end{pmatrix},$$

$d = \text{diag}(A)$, $\sigma = \text{svd}(A)$. A is (bi)diagonal, and its singular values are on the diagonal. However,

$$d = \text{diag}(A) = \begin{pmatrix} 9.999999999999999e+249 \\ 1.0000000000000000e-201 \end{pmatrix},$$

$$\sigma = \text{svd}(A) = \begin{pmatrix} 9.999999999999999e+249 \\ 1.000000000000\textcolor{red}{16167}e-201 \end{pmatrix}.$$

$$\lambda = \text{eig}(A) = \begin{pmatrix} 9.999999999999999e+249 \\ 1.0000000000000000e-201 \end{pmatrix}$$

The 2×2 example

Introduction

Problems

Machine numbers

Examples

Consequences

Goals

Scaling:
examplesNumerical
rank revealingEigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

LAPACK's driver routine xGESVD computes $\alpha = \max_{i,j} |A_{ij}|$ and scales the input matrix A with $(1/\alpha)\sqrt{\nu}/\varepsilon$ (if $\alpha < \sqrt{\nu}/\varepsilon$) or with $(1/\alpha)\varepsilon\sqrt{\omega}$ (if $\alpha > \varepsilon\sqrt{\omega}$). Here ε , ν and ω denote the round-off unit, underflow and overflow thresholds, respectively.

Let $\alpha = \max_{i,j} |A_{ij}|$, $\varepsilon = \text{eps}/2$, $\omega = \text{realmax}$, $\nu = \text{realmin}$, $\mathbf{s} = \varepsilon\sqrt{\omega}/\alpha$, and scale A with \mathbf{s} . The singular values of $\mathbf{s}A$ are on its diagonal; scaling the diagonal of $\mathbf{s}A$ with $1/\mathbf{s}$ changes the (2, 2) entry precisely to $1.0000000000\mathbf{16167e-201}$. Five digits in the second singular value of a 2×2 diagonal matrix are lost due to scaling $\sigma = (\mathbf{1}/\mathbf{s}) * (\mathbf{s} * d)$. (In MATLAB, $\omega \approx 1.79 \cdot 10^{308}$, $\nu \approx 2.22 \cdot 10^{-308}$.) The problem is not removed if \mathbf{s} is changed to the closest integer power of two.

Note that this scaling is designed to avoid overflow in the implicit use of $A^T A$.

Introduction

Problems

Machine numbers

Examples

Consequences

Goals

Scaling: examples

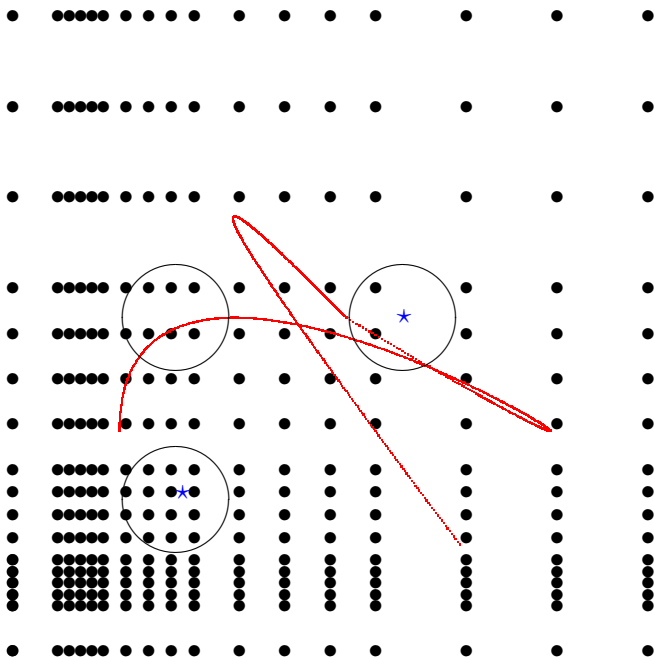
Numerical rank revealing

Eigenvalues and singular values

Jacobi method

Accurate PSVD and applications

Concluding remarks



$(0, 0)$

Early loss of definiteness

Introduction

Problems
Machine numbers

Examples

Consequences
GoalsScaling:
examplesNumerical
rank revealingEigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

The stiffness matrix of a mass spring system with 3 masses

 with spring constants $k_1 = k_3 = 1$, $k_2 = \varepsilon/2$

$$K = \begin{pmatrix} k_1 + k_2 & -k_2 & 0 \\ -k_2 & k_2 + k_3 & -k_3 \\ 0 & -k_3 & k_3 \end{pmatrix}, \quad \lambda_{\min}(K) \approx \varepsilon/4.$$

The true and the computed assembled matrix are

$$K = \begin{pmatrix} 1 + \frac{\varepsilon}{2} & -\frac{\varepsilon}{2} & 0 \\ -\frac{\varepsilon}{2} & 1 + \frac{\varepsilon}{2} & -1 \\ 0 & -1 & 1 \end{pmatrix}, \quad \tilde{K} = \begin{pmatrix} 1 & -\frac{\varepsilon}{2} & 0 \\ -\frac{\varepsilon}{2} & 1 & -1 \\ 0 & -1 & 1 \end{pmatrix}.$$

 \tilde{K} is component-wise relative perturbation of K with

$$\max_{i,j} \frac{|\tilde{K}_{ij} - K_{ij}|}{|K_{ij}|} = \frac{\varepsilon}{(2 + \varepsilon)} < \varepsilon/2.$$

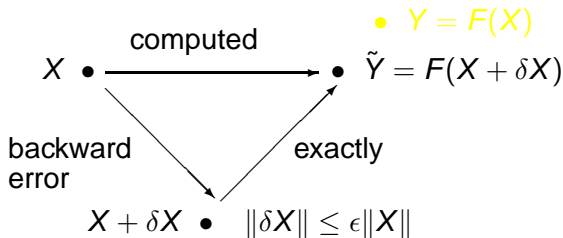
 \tilde{K} is indefinite with $\lambda_{\min}(\tilde{K}) \approx -\varepsilon^2/8$. Too late for $\lambda_{\min}(K)$. \therefore

Consequences

Almost never have exactly given data. Have $A \in \mathbb{F}^{m \times n}$ as approximation of an ideal, not accessible A_0 , $A = A_0 + E$. Do not have E , but know that $\|E\|/\|A\| \approx f(m, n)\mathbf{u}$ is small. $A \in \mathcal{B} = \{A_0 + E, \|E\| \leq \varepsilon\}$ and any $X \in \mathcal{B}$ is just as good as A .

- Full rank matrices dense in $\mathbf{M}_{m \times n}$. What is then the rank of $A_0 = A - E$? Rank of A ? Any technique will fail over \mathbb{F} .
- Chance to compute zero exactly is exactly zero.
- Matrices with simple eigenvalues dense in $\mathbf{M}_{n \times n}$. Jordan form? Diagonalizability?
- Is A +definite? Invertible? Orthogonal? Stable ($\operatorname{Re}(\lambda(A)) < 0$)? $A^{-1} = ?$ $A^\dagger = ?$
- In 1950 Goldstine and von Neuman concluded that solving linear systems with $n > 15$ with guaranteed accuracy would be nearly impossible!

Error? Distance to what?!?



Backward stability: solve exactly a problem close to X
 Not preserved under composition of mappings

$$\rightarrow \|\delta X\| \leq \epsilon \|X\|, \quad \|\delta X(:, i)\| \leq \epsilon \|X(:, i)\|$$

$$\rightarrow |\delta X_{ij}| \leq \epsilon |X_{ij}|, \quad |\delta X_{ij}| \leq \epsilon \sqrt{|X_{ii} X_{jj}|}$$

$\rightarrow X + \delta X$ same structure as X

Perturbation theory: $\|\tilde{Y} - Y\| \leq K \cdot \|\delta X\|$

Von Neumann, Turing, Givens, Wilkinson

Ill-conditioned = close to ill-posed

Relative condition number

$$\kappa(\mathcal{F}, X) = \limsup_{\Delta X \rightarrow 0} \frac{\frac{\|\mathcal{F}(X + \Delta X) - \mathcal{F}(X)\|}{\|\mathcal{F}(X)\|}}{\|\Delta X\|/\|X\|} = \frac{\|D\mathcal{F}(X)\| \|X\|}{\|\mathcal{F}(X)\|}$$

For $A \mapsto A^{-1}$, $\kappa(A) = \|A\| \cdot \|A^{-1}\|$, and the bad set is the variety of singular matrices.

$$\frac{\text{distance}(A, \text{bad})}{\|A\|} = \frac{1}{\kappa(A)}, \quad \text{bad} = \det^{-1}(\{0\}).$$

$$A_{ij} \mapsto A_{ij} + \epsilon |A_{ij}|$$

$$\inf\{|\epsilon| : \det(A + \epsilon E) = 0\} = \frac{1}{\rho_0(A^{-1}E)}$$

Probability of being too close to bad set. Algebraic and geometric properties of bad sets.

Introduction

- Problems
- Machine numbers
- Examples
- Consequences
- Goals

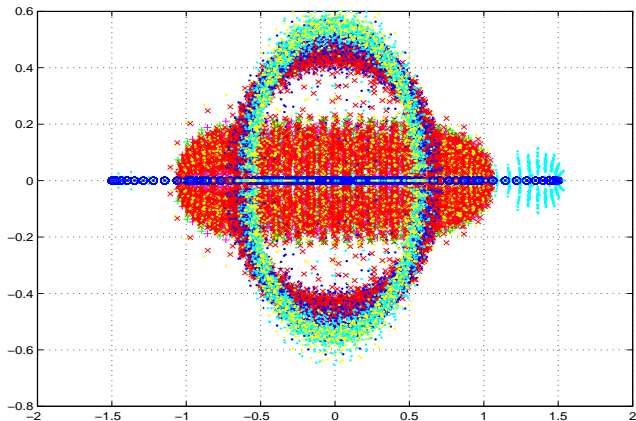
Scaling:
examplesNumerical
rank revealingEigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

Eigenvalue assignment

$\alpha_1, \dots, \alpha_n$ given. Find K such that the spectrum of $A + BK^T$ is $\{\alpha_1, \dots, \alpha_n\}$. Try many B 's and methods to hit \odot :



Placing plenty of poles is pretty preposterous (Chunyang, Laub, Mehrmann)

Goals of this lecture

We develop sharp high precision numerical tools for linear algebra problems in control theory. In this lecture, we illustrate some aspects of the development of such tools.

- We show how things can go wrong, even in computing some elementary matrix factorization in order to determine the matrix numerical rank. We stress the necessity of strict mathematical approach to numerical software development.
- The symmetric eigenvalue and the singular value problems are known to be well-conditioned. We show that numerical algorithms do not always deliver optimal accuracy. Using only orthogonal transformations in the diagonalization process does not guarantee accurate results. Perturbation theory important ingredient.
- Higher standard solutions. Accurate NLA methods make other computational tasks numerically feasible.

$$1\text{ m} = 100\text{ cm} = 1000\text{ mm}$$

$$\dot{x}(t) = Ax(t) + Bu(t):$$

$$\begin{pmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \\ \dot{x}_4 \\ \dot{x}_5 \\ \dot{x}_6 \end{pmatrix} = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ -\frac{k_p}{m_p} & -\frac{c_p}{m_p} & \frac{k_p}{m_p} & \frac{c_p}{m_p} & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ \frac{k_p}{m_s} & \frac{c_p}{m_s} & -\frac{k_s+k_p}{m_s} & -\frac{c_s+c_p}{m_s} & \frac{k_s}{m_s} & \frac{c_s}{m_s} \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & \frac{k_s}{m_{us}} & \frac{c_s}{m_{us}} & -\frac{k_s+k_t}{m_{us}} & -\frac{c_s}{m_{us}} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \end{pmatrix} + Bu$$

- x_1 displacement in meters (m); x_2 speed (m/s)
- m_p, m_s, m_{us} mass (kg)
- k_p, k_s, k_t spring stiffness (N/m)
- c_s, c_p damping coefficient (N s/m)
- $A_{23} = \frac{11000}{79} [N/m/kg]$, $A_{24} = \frac{800}{79} [Ns/m/kg]$, ...
- In different units, $x(t) = D\hat{x}(t)$, D diagonal scaling.
- How to interpret $\|x\|_2$, $\|A\|_F$, $\|\delta A\|_F$? Big? Small?

Diagonalizing the Grammians

$$\dot{x}(t) = Ax(t) + Bu(t), \quad x(0) = x_0$$

$$y(t) = Cx(t)$$

Grammians $H = L_H L_H^T$, $M = L_M L_M^T$ via Lyapunov equations:

$$AH + HA^T = -BB^T, \quad A^T M + MA = -C^T C.$$

Hankel SV, $\sigma_i = \sqrt{\lambda_i(HM)}$, $HM \mapsto T^{-1}HMT = \Sigma^2$.

Different scaling (change of units, x may contain quantities of different physical nature) $x(t) = D\hat{x}(t)$; $A \mapsto D^{-1}AD$, $B \mapsto D^{-1}B$, $C \mapsto CD$;

$$H \mapsto \hat{H} = D^{-1}HD^{-T}, \quad M \mapsto \hat{M} = D^T MD$$

Change of units (scaling) changes classical condition numbers $\kappa_2(H)$, $\kappa_2(M)$ thus making an algorithm numerically inaccurate/unstable, while the underlying problem is the same. Is this acceptable?!?

Integral equation

Consider numerical solution of the integral equation

$$y(\xi) = \int_a^b K(\xi, \zeta) x(\zeta) d\zeta$$

Here y denotes measured unknown function x distorted by the instrument with known kernel $K(\cdot, \cdot)$. If the equation is discretized at $\xi_1 < \dots < \xi_m$, and the integral is computed using quadrature rule with the nodes $\zeta_1 < \dots < \zeta_n$ and weights d_1, \dots, d_n , then

$$y(\xi_i) = \sum_{j=1}^n d_j K(\xi_i, \zeta_j) x(\zeta_j) + e_i, \quad e_i = \text{error}, \quad i = 1, \dots, m.$$

Set $y = (y(\xi_i))_{i=1}^m$, $K = (K(\xi_i, \zeta_j)) \in \mathbf{R}^{m \times n}$, $D = \text{diag}(d_i)_{i=1}^n$. An approximation $x = (x_j)_{j=1}^n$ of $(x(\zeta_j))_{j=1}^n$ is obtained by solving the linear regression problem

$$y = K D x + e, \quad x \in \mathbf{R}^n, \quad e = (e_i)_{i=1}^m.$$

$$\dots y(\xi) = \int_a^b K(\xi, \zeta) x(\zeta) d\zeta$$

Introduction

Scaling:
examplesPhysical unit
changesIntegral equations
and least squaresA symmetric 3×3
exampleNumerical
rank revealingEigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

$$y = K Dx + e, \quad x \in \mathbf{R}^n, \quad e = (e_i)_{i=1}^m.$$

with vector e dominated by statistically independent measurement errors from $\mathcal{N}(0, S^2)$, where positive definite $S = \text{diag}(s_i)_{i=1}^n$ carries standard deviations of the e_i 's. A good estimate of S is usually available.

Wanted is an estimate \tilde{x} of x . To normalize the error variances, the model is scaled with S^{-1} to get

$$b = Ax + e', \quad b = S^{-1}y, \quad A = S^{-1}KD, \quad e' = S^{-1}e.$$

Hence, we solve $\|b - Ax\|_2 \rightarrow \min$

So, what does it mean if we have $A + \delta A$ with backward error (or initial uncertainty) δA with $\|\delta A\|_F \ll \|A\|_F$?

Compare with $A + \delta A = S^{-1}(K + \delta K)D$ with $\|\delta K\|_F \ll \|K\|_F$

What is the spectrum of H ?

$$H = \begin{pmatrix} 10^{40} & 10^{29} & 10^{19} \\ 10^{29} & 10^{20} & 10^9 \\ 10^{19} & 10^9 & 1 \end{pmatrix};$$

use MATLAB, $\text{eps} \approx 2.22 \cdot 10^{-16}$

$$\text{eig}(H) = \begin{matrix} 1.0000000000000000e+040 \\ -8.100009764062724e+019 \\ -3.966787845610502e+023 \end{matrix}$$

$$L = \text{chol}(H)' \quad (H = LL^T)$$

$$L = \begin{pmatrix} 1.00000000e+20 & 0 & 0 \\ 9.99999999e+8 & 9.9498743e+9 & 0 \\ 9.99999999e-2 & 9.0453403e-2 & 9.9086738e-1 \end{pmatrix}$$

Is H **positive** definite?

Introduction

Scaling:
examplesPhysical unit
changesIntegral equations
and least squaresA symmetric 3×3
exampleNumerical
rank revealingEigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

What is the spectrum of H now?

Introduction

Scaling:
examplesPhysical unit
changesIntegral equations
and least squaresA symmetric 3×3
exampleNumerical
rank revealingEigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

	$\text{eig}(H)$	$\text{eig}(P^T H P), P \simeq (2, 1, 3)$
λ_1	1.0000000000000000e+40	1.0000000000000000e+40
λ_2	-8.100009764062724e+19	9.9000000000000000e+19
λ_3	-3.966787845610502e+23	9.818181818181818e-01
	$1./\text{eig}(\text{inv}(H))$	$\text{eig}(\text{inv}(\text{inv}(H)))$
λ_1	1.0000000000000000e+40	1.0000000000000000e+40
λ_2	9.9000000000000000e+19	9.9000000000000000e+19
λ_3	9.818181818181817e-01	9.818181818181817e-01
	$\text{eig}(H + E_1)$	$\text{eig}(H + E_2)$
λ_1	1.0000000000000000e+40	1.0000000000000000e+40
λ_2	-8.100009764062724e+19	1.208844819952007e+24
λ_3	-3.966787845610502e+23	9.899993299416013e-01

 $E_1: H_{22} = 10^{20} \rightarrow -10^{20}, E_2: H_{13}, H_{31} \rightarrow H_{13} * (1 + \text{eps}),$
 $\text{eps} \approx 2.22 \cdot 10^{-16};$ All numbers correct! ?

1 Introduction

2 Scaling: examples

3 Numerical rank revealing

Introduction

Examples

Consequences

SLICOT

Analysis

4 Eigenvalues and singular values

5 Jacobi method

6 Accurate PSVD and applications

Case study

Given $A \in \mathbb{C}^{m \times n}$, determine whether for some small δA , the matrix $A + \delta A$ is of rank $\rho < \text{rank}(A)$.

- needed and useful if A is close to matrices of lower rank (i.e. ill-conditioned)
- in the case of ill-conditioning, one does not expect much and any bad result is attributed to ill-conditioning;
- condition number can be ill-conditioned
- numerical instability in a software implementation of a basic numerical linear algebra decomposition (QR factorization with column pivoting) for almost 40 years hidden in all major numerical software packages

Eckart–Young–Mirsky–Schmidt

A $m \times n$ real of rank r

$$A = U\Sigma V^T = \sum_{k=1}^r \sigma_k u_k v_k^T$$

$$\sigma_1 \geq \cdots \geq \sigma_r > 0 = \sigma_{r+1} = \cdots = \sigma_{\min(m,n)}$$

Let $\ell < r$ and $A_\ell = \sum_{i=1}^{\ell} \sigma_i u_i v_i^T$ Then

$$\min_{\text{rank}(X) \leq \ell} \|A - X\|_F = \|A - A_\ell\|_F$$

$$\min_{\text{rank}(X) \leq \ell} \|A - X\|_2 = \|A - A_\ell\|_2$$

$$\|A - A_\ell\|_F = \sqrt{\sum_{i=\ell+1}^r \sigma_i^2}$$

$$\|A - A_\ell\|_2 = \sigma_{\ell+1}$$

QRCP with Businger–Golub pivoting

Introduction

Scaling:
examplesNumerical
rank revealing

Introduction

Examples

Consequences

SLICOT

Analysis

Eigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

permutation

$$\underbrace{A}_{m \times n} \overbrace{P}^{\text{permutation}} = Q \begin{pmatrix} R \\ 0 \end{pmatrix}, \quad R = \begin{pmatrix} \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ 0 & \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \hline 0 & 0 & \color{red}{\blacksquare} & \color{blue}{\bullet} & \color{blue}{\blacksquare} & \color{blue}{\blacklozenge} \\ 0 & 0 & 0 & \color{blue}{\bullet} & \color{blue}{\blacksquare} & \color{blue}{\blacklozenge} \\ 0 & 0 & 0 & 0 & \color{blue}{\blacksquare} & \color{blue}{\blacklozenge} \\ 0 & 0 & 0 & 0 & 0 & \color{blue}{\blacklozenge} \end{pmatrix}$$

$$Q^* Q = I_m.$$

$$|R_{ii}| \geq \sqrt{\sum_{k=i}^j |R_{kj}|^2}, \quad \text{for all } 1 \leq i \leq j \leq n. \quad (3)$$

$$|R_{11}| \geq |R_{22}| \geq \cdots \geq |R_{\rho\rho}| \gg |R_{\rho+1,\rho+1}| \geq \cdots \geq |R_{nn}| \quad (4)$$

The structure (3), (4) may not be rank revealing but it must be guaranteed by the software (e.g. LAPACK, Matlab)

QRCP as preconditioner

Introduction

Scaling:
examplesNumerical
rank revealing

Introduction

Examples

Consequences

SLICOT

Analysis

Eigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

Let $AP = Q \begin{pmatrix} R \\ 0 \end{pmatrix}$; $A_c = A \cdot \text{diag}(\frac{1}{\|A(:,1)\|_2}, \dots, \frac{1}{\|A(:,n)\|_2})$;

$$R_c = R \cdot \text{diag}(\frac{1}{\|R(:,1)\|_2}, \dots, \frac{1}{\|R(:,n)\|_2}) = \begin{pmatrix} \downarrow & \downarrow & \downarrow \\ 0 & \downarrow & \downarrow \\ 0 & 0 & \downarrow \end{pmatrix};$$

$$R_r = \text{diag}(\frac{1}{\|R(1,:)\|_2}, \dots, \frac{1}{\|R(n,:)\|_2}) \cdot R = \begin{pmatrix} \rightarrow & \rightarrow & \rightarrow \\ 0 & \rightarrow & \rightarrow \\ 0 & 0 & \rightarrow \end{pmatrix}.$$

Proposition:

Let $AP = QR$, where $|R_{ij}| \geq \sqrt{\sum_{k=i}^j |R_{kj}|^2}$, $1 \leq i \leq j \leq n$.

Then $\| |R_r^{-1}| \|_2 \leq \sqrt{n} \| |R_c^{-1}| \|_2$, $\kappa_2(R_r) \leq n^{3/2} \kappa_2(A_c)$.

Moreover, $\|R_r^{-1}\|_2$ is bounded by $O(2^n)$, independent of A .
With exception of rare pathological cases, $\|R_r^{-1}\|_2$ is below $O(n)$ for any A . **RR^* is more diagonal than R^*R .**

Example:

Let $A = \text{Hilbert}(100)$. $\kappa_2(A) > 10^{150} \gg \text{cond}(A) \approx 3.6e19$

$\kappa_2(A_c) = \kappa_2(R_c) > 10^{19}$, $\kappa_2(R_r) \approx 48.31$. Repeat with

$A \leftarrow R^T$, $P = I$, to get new $\kappa_2(R_r) \approx 3.22$.

Examples of failure (Matlab)

Zlatko Drmač

Introduction

Scaling:
examples

Numerical
rank revealing

Introduction

Examples

Consequences

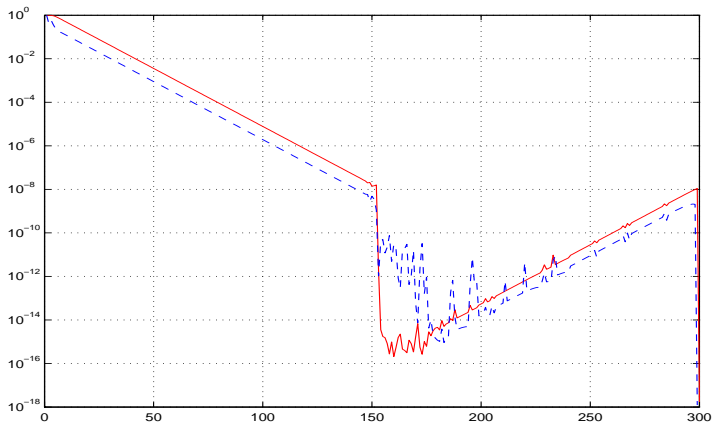
SLICOT

Analysis

Eigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applications

Concluding
remarks


$$|R_{ii}|, \max_{j \geq i} \sqrt{\sum_{k=i}^j |R_{kj}|^2}, R = \begin{pmatrix} \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ 0 & \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ 0 & 0 & \color{red}{\blacksquare} & \bullet & \color{blue}{\blacksquare} & \color{blue}{\blacklozenge} \\ 0 & 0 & 0 & \bullet & \color{blue}{\blacksquare} & \color{blue}{\blacklozenge} \\ 0 & 0 & 0 & 0 & \color{blue}{\blacksquare} & \color{blue}{\blacklozenge} \\ 0 & 0 & 0 & 0 & 0 & \color{blue}{\blacklozenge} \end{pmatrix}$$

Examples of failure (Matlab)

Zlatko Drmač

Introduction

Scaling:
examplesNumerical
rank revealing

Introduction

Examples

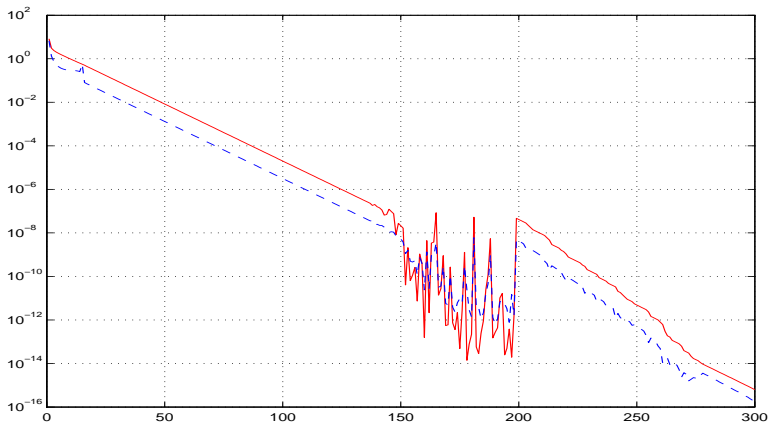
Consequences

SLICOT

Analysis

Eigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

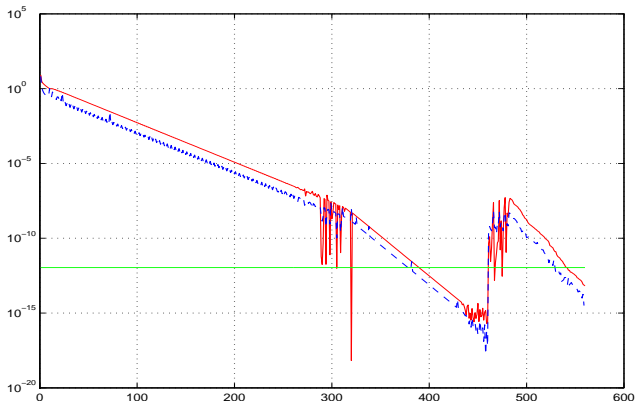
$$|R_{ii}|, \max_{j \geq i} \sqrt{\sum_{k=i}^j |R_{kj}|^2}, R = \begin{pmatrix} \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ 0 & \blacksquare & \blacksquare & \bullet & \blacksquare & \blacklozenge \\ 0 & 0 & \blacksquare & \bullet & \blacksquare & \blacklozenge \\ 0 & 0 & 0 & \bullet & \blacksquare & \blacklozenge \\ 0 & 0 & 0 & 0 & \blacksquare & \blacklozenge \\ 0 & 0 & 0 & 0 & 0 & \blacklozenge \end{pmatrix}$$

Consequences (Matlab)

$$\|Ax - d\|_2 \rightarrow \min; x = A \backslash d \text{ (Matlab LS solution)}$$

Warning: Rank deficient, rank = 304 tol =

$$1.0994e-012. R = \begin{pmatrix} \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ 0 & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ 0 & 0 & \color{red}{\blacksquare} & \bullet & \color{blue}{\blacksquare} \\ 0 & 0 & 0 & \bullet & \color{blue}{\blacksquare} \\ 0 & 0 & 0 & 0 & \color{blue}{\blacksquare} \end{pmatrix}$$



$\text{rank}(A, 1.0994e-12)$ returns 466

Consequences

Any routine based on xQRDC (LINPACK) or xGEQPF, xGEQP3 (LAPACK) can catastrophically fail.

- xGGEQPX (TOMS # 782, rank revealing QRF)
- xGELSX and xGELSY in LAPACK ($\|Ax - b\|_2 \rightarrow \min$)
- xGGSPV in LAPACK (GSVD of (A, B))

$$U^T A Q = \begin{pmatrix} 0 & A_{12} & A_{13} \\ 0 & 0 & A_{23} \\ 0 & 0 & 0 \end{pmatrix}, \quad V^T B Q = \begin{pmatrix} 0 & 0 & B_{13} \\ 0 & 0 & 0 \end{pmatrix}.$$

- ... and many others ... long list. **Need a new xGEQP3.**

Resolved by Drmač and Bujanović (ACM TOMS, 2008) and included in LAPACK.

In control, included in SLICOT in 2010.

The SLICOT (Subroutine Library In Control Theory)

- is used as computational layer in sophisticated CACSD packages such as EASY5 (since 2002. MSC.Software, initially developed in the Boeing Company), Matlab (The MathWorks) and Scilab (INRIA).
- Since its initial release, SLICOT has been growing at an impressive rate, from 90 user-callable subroutines in 1997., 200 subroutines in 2004., 470 subroutines in 2009., ...
- Efficiency and reliability based on BLAS, LAPACK and state of the art numerical linear algebra

The problem illustrated in the previous examples of QRCP failure affects SLICOT and thus many other control theory libraries.

$$\begin{aligned} E\dot{x} &= Ax + Bu \\ y &= Cx + Du. \end{aligned} \tag{5}$$

- Strategically placed "WRITE(*,*) variable" statements in the affected subroutines can completely change the computed properties of (5).
- Substantial variations of the output can also be caused by changing the compiler and optimizer options.
- This is undesired behavior, even if the computation is backward stable, and even if it is doomed to fail, due to ill-conditioning.

The problem occurs only at certain distance to singularity, and the rank revealing task itself is usually performed if the matrix is close to singularity. Since many things can happen close to singularity, any ill-behavior is usually attributed to ill-conditioning and the true cause remains inconspicuous.

SLICOT Example: MB03OY

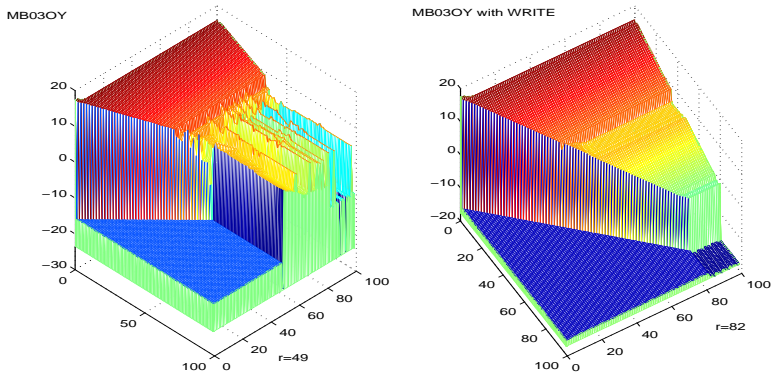


Figure: Left: The matrix R computed by MB03OY, shown by `meshz(log10(abs(R)))`. The computed rank is 49. Right: **The matrix R computed with MB03OY, with "WRITE(*,*) TEMP2" statement added after the line 339 in MB03OY.f. The computed rank is 82.**

Affected SLICOT routines

1. MB03OY \leftarrow :1.1 AB01ND \leftarrow 1.1.1 AB01OD1.2 AB08NX \leftarrow : 1.2.1 AB08ND
1.2.2 AB08MD \leftarrow 1.2.2.1 AB09JD1.3 AG08BY \leftarrow 1.3.1 AG08BD1.4 MB02QD \leftarrow 1.4.1 SB01DD1.5 TB01UD \leftarrow :1.5.1 TB01PD \leftarrow :1.5.1.1.1 SB10ZP \leftarrow 1.5.1.1 TD04AD \leftarrow 1.5.1.1.1 SB10YD \leftarrow

1.5.1.1.1.1 SB10MD

1.5.1.2 AB09ID

1.5.2 TB03AD \leftarrow 1.5.2.1 TD03AD1.5.3 TB04AY \leftarrow 1.5.3.1 TB04AD1.6 TG01FD \leftarrow 1.6.1 AG08BD

Introduction

Scaling:
examplesNumerical
rank revealing

Introduction

Examples

Consequences

SLICOT

Analysis

Eigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

... affected SLICOT routines

Introduction

Scaling:
examplesNumerical
rank revealing

Introduction

Examples

Consequences

SLICOT

Analysis

Eigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks2. MB04GD \leftarrow 2.1 MB03PD3. MB03OD \leftarrow :3.1 IB01ND \leftarrow 3.1.1 IB01AD \leftarrow : 3.1.1.1 IB03AD
3.1.1.2 IB03BD3.2 IB01PD \leftarrow 3.2.1 IB01BD \leftarrow : 3.2.1.1 IB03AD
3.2.1.2 IB03BD3.3 IB01PY \leftarrow 3.3.1 IB01PD \leftarrow 3.3.1.1 IB01BD \leftarrow : 3.3.1.1.1 IB03AD
3.3.1.1.2 IB03BD
3.4 MB02QD \leftarrow 3.4.1 SB01DD3.5 * MB02YD \leftarrow : 3.5.1 NF01BQ \leftarrow 3.5.1.1 NF01BP
3.5.2 MD03BY \leftarrow : 3.5.2.1 MD03BB
3.5.2.2 NF01BP3.6 * MD03BY \leftarrow : 3.6.1 MD03BB
3.6.2 NF01BP3.7 * NF01BR \leftarrow : 3.7.1 NF01BP
3.7.2 NF01BQ \leftarrow 3.7.2.1 NF01BP

60 out of 470 affected!

Introduction

Scaling:
examplesNumerical
rank revealing

Introduction

Examples

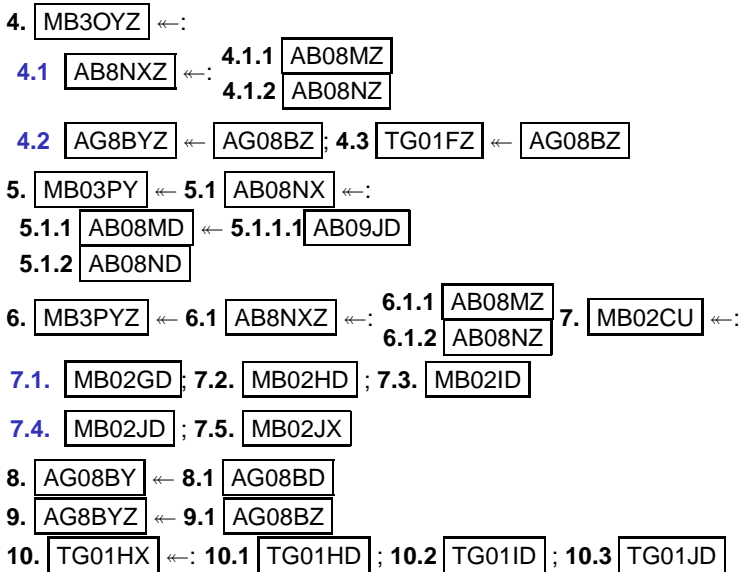
Consequences

SLICOT

Analysis

Eigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

Implementation: L(IN+A)PACK, 1970s, 1990s, ...

Introduction

Scaling:
examplesNumerical
rank revealing

Introduction

Examples

Consequences

SLICOT

Analysis

Eigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

$$A^{(k)} \Pi_k = \begin{pmatrix} \cdot & \cdot & \odot & \cdot & \oplus & \cdot \\ & \cdot & \odot & \cdot & \oplus & \cdot \\ & & \blacksquare & \circledast & \circledast & \circledast \\ & & \odot & * & * & * \\ & & \odot & * & * & * \\ & & \odot & * & * & * \end{pmatrix}, \quad \mathbf{a}_j^{(k)} = \begin{pmatrix} \oplus \\ \oplus \\ \circledast \\ * \\ * \\ * \end{pmatrix} \equiv \begin{pmatrix} \mathbf{x}_j^{(k)} \\ \mathbf{z}_j^{(k)} \end{pmatrix}$$

$$\mathbf{H}_k \mathbf{z}_k^{(k)} = \begin{pmatrix} R_{kk} \\ 0 \end{pmatrix}, \quad \mathbf{H}_k \mathbf{z}_j^{(k)} = \begin{pmatrix} \beta_j^{(k+1)} \\ \mathbf{z}_j^{(k+1)} \end{pmatrix}, \quad \omega_j^{(k)} = \|\mathbf{z}_j^{(k)}\|$$

$$\|\mathbf{z}_j^{(k+1)}\| \equiv \omega_j^{(k+1)} = \sqrt{(\omega_j^{(k)})^2 - (\beta_j^{(k+1)})^2}, \quad \text{provided that}$$

$$\text{computed} \left(\left(1 - \left(\frac{\tilde{\beta}_j^{(k+1)}}{\tilde{\omega}_j^{(k)}} \right)^2 \right) \cdot \left(\frac{\tilde{\omega}_j^{(k)}}{\tilde{\nu}_j} \right)^2 \right) > \text{tol}, \quad \text{tol} \approx 20 \cdot \text{eps},$$

L(IN+A)PACK update

Introduction

Scaling:
examplesNumerical
rank revealing

Introduction

Examples

Consequences

SLICOT

Analysis

Eigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

```

DO 30 J = I+1, N
  IF ( WORK( J ).NE.ZERO ) THEN
    TEMP = ONE - ( ABS( A( I, J ) ) / WORK( J ) )**2
    TEMP = MAX( TEMP, ZERO )
    TEMP2 = ONE + 0.05*TEMP*( WORK( J ) / WORK( N+J ) )**2
    WRITE(*,*) TEMP2
    IF( TEMP2.EQ.ONE ) THEN
      IF( M-I.GT.0 ) THEN
        WORK( J ) = SNRM2( M-I, A( I+1, J ), 1 )
        WORK( N+J ) = WORK( J )
      ELSE
        WORK( J ) = ZERO
        WORK( N+J ) = ZERO
      END IF
    ELSE
      WORK( J ) = WORK( J )*SQRT( TEMP )
    END IF
  END IF
30 CONTINUE

```

g77 -c -O -ffloat-store

Critical part in the column norm update. (For the full source see <http://www.netlib.org/lapack/single/sgeqpf.f>)

NEW update

Introduction

Scaling:
examplesNumerical
rank revealing

Introduction

Examples

Consequences

SLICOT

Analysis

Eigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

$$A^{(k)} \Pi_k = \begin{pmatrix} \cdot & \cdot & \odot & \cdot & \oplus & \cdot \\ & \cdot & \odot & \cdot & \oplus & \cdot \\ & & \blacksquare & \circledast & \circledast & \circledast \\ & & \odot & * & * & * \\ & & \odot & * & * & * \\ & & \odot & * & * & * \end{pmatrix}, \quad \mathbf{a}_j^{(k)} = \begin{pmatrix} \oplus \\ \oplus \\ \circledast \\ * \\ * \\ * \end{pmatrix} \equiv \begin{pmatrix} \mathbf{x}_j^{(k)} \\ \mathbf{z}_j^{(k)} \end{pmatrix} \quad (6)$$

$$\mathbf{H}_k \mathbf{z}_k^{(k)} = \begin{pmatrix} R_{kk} \\ \mathbf{0} \end{pmatrix}, \quad \mathbf{H}_k \mathbf{z}_j^{(k)} = \begin{pmatrix} \beta_j^{(k+1)} \\ \mathbf{z}_j^{(k+1)} \end{pmatrix}, \quad \omega_j^{(k)} = \|\mathbf{z}_j^{(k)}\| \quad (7)$$

$$\|\mathbf{a}_j^{(k)}\| = \alpha_j^{(k)} = \alpha_j^{(0)}; \quad \xi_j^{(k+1)} = \sqrt{(\xi_j^{(k)})^2 + (\beta_j^{(k+1)})^2}$$

$$\|\mathbf{z}_j^{(k+1)}\| \equiv \omega_j^{(k+1)} = \sqrt{(\alpha_j^{(0)})^2 - (\xi_j^{(k+1)})^2}$$

New update – conclusion

Introduction

Scaling:
examplesNumerical
rank revealing

Introduction

Examples

Consequences

SLICOT

Analysis

Eigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

- Provably delivers Businger–Golub structured R (up to roundoff)
- For the computed $\tilde{R} = R + \delta R$, not only $\|\delta R\|/\|R\|$, but also $\|\delta R(:, i)\|/\|R(:, i)\|$ and $\|\delta R(i, :)\|/\|R(i, :)\|$ are small.
- Row scaled \tilde{R}_r well conditioned. $\begin{pmatrix} \overrightarrow{0} & \overrightarrow{0} & \overrightarrow{0} \\ 0 & \overrightarrow{0} & \overrightarrow{0} \\ 0 & 0 & \overrightarrow{0} \end{pmatrix}$
- Same efficiency as original routines
- Makes many other solvers more robust and can prevent catastrophes in mission critical applications
- **Included in LAPACK, SLICOT**

$$H - \lambda I, HM - \lambda I$$

$$\dot{x}(t) = Ax(t) + Bu(t), \quad x(0) = x_0$$

$$y(t) = Cx(t)$$

Grammians $H = L_H L_H^T$, $M = L_M L_M^T$ via Lyapunov equations:

$$H = \int_0^\infty e^{tA} B B^T e^{tA^T} dt, \quad M = \int_0^\infty e^{tA^T} C^T C e^{tA} dt$$

$$AH + HA^T = -BB^T, \quad A^T M + MA = -C^T C.$$

Hankel singular values, $\sigma_i = \sqrt{\lambda_i(HM)}$. Need spectral decomposition of the product HM of positive definite matrices, $HM \mapsto T^{-1} HMT = \Sigma^2$. New state coo's $x(t) = T\hat{x}(t)$; $A \mapsto T^{-1}AT$, $B \mapsto T^{-1}B$, $C \mapsto CT$;

$$H \mapsto \hat{H} = T^{-1}HT^{-T} = \Sigma, \quad M \mapsto \hat{M} = T^T MT = \Sigma$$

Backward stability: eig()

$H = H^T$, $n \times n$ symmetric.

$$Hu_i = \lambda_i u_i, \quad H = U \Lambda U^T, \quad \Lambda = \text{diag}(\lambda_i)_{i=1}^n$$

Symm. EigenValue Problem perfect ★ ★ ★:

★ eigenvalues real, eigenvectors orthogonal

★ algorithms use orthogonal transformations

★ Weyl: If $H \rightsquigarrow H + \delta H$, then $\lambda \rightsquigarrow \lambda + \delta \lambda$, with

$$\max_{\lambda} |\delta \lambda| \leq \|\delta H\|$$

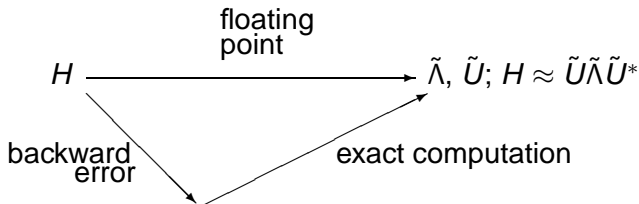
$$\cdots U_k^T \cdots (U_2^T (U_1^T H U_1) U_2) \cdots U_k \cdots \xrightarrow{U} \Lambda$$

Computed (finite prec., $O(n^3)$) $\tilde{U} \approx U$, $\tilde{\Lambda} \approx \Lambda$.

Backward stability:

$$\tilde{U}^T (H + \delta H) \tilde{U} \approx \tilde{\Lambda}, \quad \frac{\|\delta H\|}{\|H\|} \leq \epsilon \approx 10^{-16} \text{ small.}$$

Forward error



$$H + \delta H, \|\delta H\| \leq \epsilon \|H\|, \epsilon \text{ small}$$

Weyl: $|\delta \lambda_i| \leq \|\delta H\|$, $i = 1, \dots, n$. **Bad news for small λ_i 's:**

$$\frac{|\delta \lambda_i|}{|\lambda_i|} \leq \epsilon \frac{\|H\|}{|\lambda_i|}$$

Let $\kappa_2(H) = \|H\| \|H^{-1}\|$. Then

$$\max_i \left| \frac{\delta \lambda_i}{\lambda_i} \right| \leq \kappa_2(H) \frac{\|\delta H\|}{\|H\|}.$$

Want better accuracy for better inputs.

Error in the eigenvalues

Let $H = LL^T \succ 0$ and $\tilde{L}\tilde{L}^T = H + \delta H \succ 0$, $|\delta H_{ij}| \leq \eta_C \sqrt{H_{ii}H_{jj}}$.

Compare the eigenvalues of H and $\tilde{H} = H + \delta H = \tilde{L}\tilde{L}^T$:

- $H = LL^T$ is similar to L^TL , $H \sim L^TL$.
- Let $Y = \sqrt{I + L^{-1}\delta HL^{-T}}$. Then

$$H + \delta H = L(I + L^{-1}\delta HL^{-T})L^T = LY Y^T L^T \sim Y^T L^T LY.$$

Compare $\lambda_i(L^TL) = \lambda_i(H)$ and $\lambda_i(Y^T L^T LY) = \lambda_i(H + \delta H)$.

- Ostrowski: $\tilde{M} = Y^T M Y$, then, for all i , $\lambda_i(\tilde{M}) = \lambda_i(M) \xi_i$, $\lambda_{\min}(Y^T Y) \leq \xi_i \leq \lambda_{\max}(Y^T Y)$. Here $Y^T Y = I + L^{-1}\delta HL^{-T}$.
- Hence $|\lambda_i(H) - \lambda_i(\tilde{H})| \leq \lambda_i(H) \|L^{-1}\delta HL^{-T}\|_2$,

$$\begin{aligned} \|L^{-1}\delta HL^{-T}\|_2 &= \|L^{-1}DD^{-1}\delta HD^{-1}DL^{-T}\|_2 = \|L^{-1}D(\delta H_s)DL^{-T}\|_2 \\ &\leq \|L^{-1}D\|_2^2 \|\delta H_s\|_2 = \|DL^{-T}L^{-1}D\|_2 \|\delta H_s\|_2 \\ &= \|(D^{-1}HD^{-1})^{-1}\|_2 \|\delta H_s\|_2 = \|H_s^{-1}\|_2 \|\delta H_s\|_2 \end{aligned}$$

Error in the eigenvalues

Since $\delta H_s = (\delta H_{ij} / \sqrt{H_{ii} H_{jj}})$

$$\max_i \left| \frac{\delta \lambda_i}{\lambda_i} \right| \leq \|H_s^{-1}\|_2 \underbrace{\left\| \left[\frac{\delta H_{ij}}{\sqrt{H_{ii} H_{jj}}} \right] \right\|_2}_{\leq n \eta_C}$$

Compare with $\max_i \left| \frac{\delta \lambda_i}{\lambda_i} \right| \leq \kappa_2(H) \frac{\|\delta H\|_2}{\|H\|_2}$

Van der Sluis: $\|H_s^{-1}\|_2 \leq \kappa_2(H_s) \leq n \min_{D=\text{diag}} \kappa_2(DHD)$.

Our 3×3 example: $H = DH_s D$, $D = \text{diag}(10^{20}, 10^{10}, 1)$,

$$\begin{pmatrix} 10^{40} & 10^{29} & 10^{19} \\ 10^{29} & 10^{20} & 10^9 \\ 10^{19} & 10^9 & 1 \end{pmatrix} = DH_s D = D \begin{pmatrix} 1 & 0.1 & 0.1 \\ 0.1 & 1 & 0.1 \\ 0.1 & 0.1 & 1 \end{pmatrix} D,$$

$$\kappa_2(H) > 10^{40}, \kappa_2(H_s) < 1.4, \|H_s^{-1}\|_2 < 1.2.$$

Positive definiteness in floating-point

Demmel and Veselić

Let $H = DH_sD$, where $D = \text{diag}(\sqrt{H_{ii}})_{i=1}^n$, and let $\lambda_{\min}(H_s)$ be the minimal eigenvalue of H_s .

If δH is symmetric perturbation such that $H + \delta H$ is not positive definite, then

$$\max_{1 \leq i, j \leq n} \frac{|\delta H_{ij}|}{\sqrt{H_{ii}H_{jj}}} \geq \frac{\lambda_{\min}(H_s)}{n} = \frac{1}{n \|H_s^{-1}\|_2}.$$

If $\delta H = -\lambda_{\min}(H_s)D^2$, then $\max_{i,j} \frac{|\delta H_{ij}|}{\sqrt{H_{ii}H_{jj}}} = \lambda_{\min}(H_s)$ and

$H + \delta H$ is singular.

If $\|H_s^{-1}\|_2$ is too big ($\gtrsim 1/\varepsilon$) then H is entry-wise close to a non-definite matrix. Can say: H is numerically definite iff $\|H_s^{-1}\|_2 < 1/\varepsilon$.

Implicit definiteness

Introduction

Scaling:
examplesNumerical
rank revealingEigenvalues
and singular
values

Introduction

Backward stability

Perturbations of the
spectrumPositive definiteness
in floating point

Jacobi method

Accurate
PSVD and
applicationsConcluding
remarks

$$K = \begin{pmatrix} k_1 + k_2 & -k_2 & 0 \\ -k_2 & k_2 + k_3 & -k_3 \\ 0 & -k_3 & k_3 \end{pmatrix}.$$

Note that $K = LL^T$, where

$$L = \begin{pmatrix} \sqrt{k_1} & 0 & 0 \\ -\sqrt{k_2} & \sqrt{k_2} & 0 \\ 0 & -\sqrt{k_3} & \sqrt{k_3} \end{pmatrix} = \begin{pmatrix} \sqrt{k_1} & 0 & 0 \\ 0 & \sqrt{k_2} & 0 \\ 0 & 0 & \sqrt{k_3} \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 0 & -1 & 1 \end{pmatrix}$$

6 Accurate PSVD and applications

Jacobi, 1844, 1846

Ein leichtes Verfahren ...

$$H = H^T, H^{(k+1)} = U_k^T H^{(k)} U_k \longrightarrow \Lambda = \text{diag}(\lambda_i) \quad (k \longrightarrow \infty)$$

Each U_k annihilates (p_k, q_k) , (q_k, p_k) positions in $H^{(k)}$.

$$\cdots U_3^T U_2^T U_1^T \begin{pmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{pmatrix} U_1 U_2 U_3 \cdots = \begin{pmatrix} \bullet & \circledast & \otimes & 0 \\ \circledast & \bullet & \star & \bullet \\ \otimes & \star & \bullet & \bullet \\ 0 & \bullet & \bullet & \bullet \end{pmatrix}$$

$$U_1 = \begin{pmatrix} \cos \psi_1 & \sin \psi_1 \\ -\sin \psi_1 & \cos \psi_1 \end{pmatrix} \oplus I_{n-2}, \quad U_2 = \cdots$$

$$\boxed{\text{Jacobi rotation}} \quad \cot 2\psi_k = \frac{H_{q_k q_k}^{(k)} - H_{p_k p_k}^{(k)}}{2H_{p_k q_k}^{(k)}},$$

$$\tan \psi_k = \frac{\text{sign}(\cot 2\psi_k)}{|\cot 2\psi_k| + \sqrt{1 + \cot^2 2\psi_k}} \in \left(-\frac{\pi}{4}, \frac{\pi}{4}\right],$$

 $(p, q) = \mathcal{P}(k)$ pivot strategy, $\mathcal{P} : \mathbb{N} \rightarrow \{(i, j) : i < j\}$

Convergent strategies

Jacobi: $|h_{pq}^{(k)}| = \max_{i \neq j} |h_{ij}^{(k)}|$, $\mathcal{P}(k) = (p, q)$.

Reading Jacobi's 1846. paper recommended.

Cyclic: \mathcal{P} periodic, one full period called sweep.

Row-cyclic and column-cyclic:

$$\begin{pmatrix} \bullet & 1 \rightarrow & 2 \rightarrow & 3 \\ \bullet & \bullet & 4 \rightarrow & 5 \\ \bullet & \bullet & \bullet & 6 \\ \bullet & \bullet & \bullet & \bullet \end{pmatrix}, \quad \begin{pmatrix} \bullet & 1 \downarrow & 2 \downarrow & 4 \downarrow \\ \bullet & \bullet & 3 \downarrow & 5 \downarrow \\ \bullet & \bullet & \bullet & 6 \downarrow \\ \bullet & \bullet & \bullet & \bullet \end{pmatrix}$$

$$\text{Off}(H^{(k)}) = \sqrt{\sum_{i \neq j} (H^{(k)})_{ij}^2} \rightarrow 0 \quad (k \rightarrow \infty)$$

$H^{(k)} \rightarrow \Lambda$, $U_1 \cdots U_k \cdots \rightarrow U$ as $(k \rightarrow \infty)$; $U^T H U = \Lambda$
Asymptotically quadratic reduction of $\text{Off}(H^{(k)})$.

Forsythe, Henrici, Wilkinson, Rutishauser, Hari, Veselić
Asymptotically cubic strategies exist.

One-sided Jacobi SVD

Hestenes used implicit Jacobi for SVD of $A \in \mathbb{R}^{m \times n}$:

Diagonalize $H = H^{(0)} = A^T A$; $A \equiv A_0$.

$$H^{(1)} = V_0^T H^{(0)} V_0 = V_0^T A^T (AV_0) = A_1^T A_1$$

$$H^{(k+1)} = V_k^T H^{(k)} V_k = A_{k+1}^T A_{k+1} \longrightarrow \Lambda = \text{diag}(\lambda_i)$$

$$\Leftrightarrow A_{k+1} = A_k V_k, \text{ where } H^{(k)} = A_k^T A_k$$

V_k uses Jacobi rotation to diagonalize

$$\begin{pmatrix} h_{pp}^{(k)} & h_{pq}^{(k)} \\ h_{qp}^{(k)} & h_{qq}^{(k)} \end{pmatrix} \quad \begin{aligned} h_{pp}^{(k)} &= \|A_k(1:m, p)\|^2 \\ h_{qq}^{(k)} &= \|A_k(1:m, q)\|^2 \\ h_{pq}^{(k)} &= A_k(1:m, p)^T A_k(1:m, q) \end{aligned}$$

$h_{pp}^{(k)}, h_{qq}^{(k)}$ scalar update; $h_{pq}^{(k)}$ BLAS1 SDOT

$$A_k \longrightarrow U\Sigma, \Sigma = \text{diag}(\sqrt{\lambda_i}), U^T U = I$$

$$V_1 \cdots V_k \cdots \longrightarrow V, V^T V = I, AV = U\Sigma$$

$$A = U\Sigma V^T \text{ the SVD of } A.$$

One-sided rotation

$$d_p = \|A_k(1:m, p)\|^2, \quad d_q = \|A_k(1:m, q)\|^2, \\ \xi = A_k(1:m, p)^T A_k(1:m, q);$$

$$\text{ROTATE}(A_{1:m,p}, A_{1:m,q}, d_p, d_q, \xi, [V_{1:m,p}, V_{1:m,q}])$$

$$1: \quad \vartheta = \frac{d_q - d_p}{2 \cdot \xi}; \quad t = \frac{\text{sign}(\vartheta)}{|\vartheta| + \sqrt{1 + \vartheta^2}};$$

$$c = \frac{1}{\sqrt{1 + t^2}}; \quad s = t \cdot c;$$

$$2: \quad (A_{1:m,p} \ A_{1:m,q}) = (A_{1:m,p} \ A_{1:m,q}) \begin{pmatrix} c & s \\ -s & c \end{pmatrix};$$

$$3: \quad d_p = d_p - t \cdot \xi; \quad d_q = d_q + t \cdot \xi;$$

$$4: \quad \text{if } V \text{ is wanted then}$$

$$5: \quad (V_{1:n,p} \ V_{1:n,q}) = (V_{1:n,p} \ V_{1:n,q}) \begin{pmatrix} c & s \\ -s & c \end{pmatrix}$$

$$6: \quad \text{end if}$$

Can avoid squared norms. Can use fast rotations. Unit stride memory access. Vectorizable. Parallelizable.

Introduction

Scaling:
examplesNumerical
rank revealingEigenvalues
and singular
values

Jacobi method

Ein leichtes
VerfahrenOne-sided Jacobi
SVD

Floating point Jacobi

Provable accuracy

SVD computation in
floating pointAccurate
PSVD and
applicationsConcluding
remarks

Jacobi SVD

$\hat{p} = n(n-1)/2$; $s = 0$; **convergence = false** ;

if V is wanted **then initialize** $V = I_n$ **end if**

for $i = 1$ **to** n **do** $d_i = A_{1:m,i}^T A_{1:m,i}$ **end for**;

repeat

$s = s + 1$; $p = 0$;

for $i = 1$ **to** $n - 1$ **do**

for $j = i + 1$ **to** n **do**

$\xi = A_{1:m,i}^T A_{1:m,j}$;

if $|\xi| > m\varepsilon\sqrt{d_i d_j}$ **then**

call **ROTATE**($A_{1:m,i}$, $A_{1:m,j}$, d_i , d_j , ξ , [$V_{1:m,i}$, $V_{1:m,j}$]) ;

else $p = p + 1$ **end if**

end for

end for

if $p = \hat{p}$ **then** **convergence=true**; **go to** ► **end if**

until $s > 30$

► **if** **convergence** **then** $\Sigma_{ii} = \sqrt{d_i}$, $U_{1:m,i} = A_{1:m,i}\Sigma_{ii}^{-1}$, $i = 1 : n$;
else Error: *Did not converge in 30 sweeps.* **end if**

Jacobi in floating-point

Breakthrough: Jacobi method is more accurate than QR!

Demmel and Veselić: Let $\tilde{H}^{(k)}$, denote the computed matrices. Then, in the positive definite case, one step of Jacobi in floating-point arithmetic reads

$$\tilde{H}^{(k+1)} = \hat{U}_k^T (\tilde{H}^{(k)} + \delta \tilde{H}^{(k)}) \hat{U}_k$$

where \hat{U}_k is exactly orthogonal and ε -close to the actually used Jacobi rotation \tilde{U}_k , and $\delta \tilde{H}^{(k)}$ is sparse with

$$\mathbf{e}_k = \max_{i,j} \frac{|(\delta \tilde{H}^{(k)})_{ij}|}{\sqrt{(\tilde{H}^{(k)})_{ii}(\tilde{H}^{(k)})_{jj}}} \leq \epsilon$$

Relative perturbation of eigenvalues in the k -step bounded by $n \mathbf{e}_k \|(\tilde{H}_s^{(k)})^{-1}\|_2$, $\tilde{H}_s^{(k)}$ scaled to have unit diagonal.

IMPORTANT: Stop when $\max_{i \neq j} |(\tilde{H}_s^{(k)})_{ij}| \leq \epsilon$

The accuracy depends on $\max_k \|(\tilde{H}_s^{(k)})^{-1}\|_2$

Jacobi in floating-point

If the entries of the initial H are given with relative uncertainty ε , then:

- The spectrum is determined up to relative error of order of $n\varepsilon \|H_s^{-1}\|$ (H_s diagonally scaled H to have unit diagonal)
- The symmetric Jacobi method introduces perturbation of the order of $n\varepsilon \max_k \|(\tilde{H}_s^{(k)})^{-1}\|_2$

Numerical evidence: $\max_k \|(\tilde{H}_s^{(k)})^{-1}\|_2$ behaves well.

Theoretical (still open) problem: Bound

$$\max_{k \geq 1} \frac{\| (H_s^{(k)})^{-1} \|_2}{\| H_s^{-1} \|_2} \quad \text{or} \quad \max_{k \geq 1} \frac{\kappa_2(H_s^{(k)})}{\kappa_2(H_s)}$$

Demmel, Veselić, Slapničar, Mascarenhas, Drmač

Provable accuracy

Let $H = LL^T \succ 0$, L Cholesky factor.

Use Veselić–Hari trick:

- If we apply Jacobi SVD to L , $LV = U\Sigma$, where V is the product of Jacobi rotations, then $H = U\Sigma^2U^T$.
- So, can apply Jacobi and get eigenvectors without accumulation of Jacobi rotations! This reduces flop count, memory requirements and memory traffic!
- This implicitly diagonalizes L^TL , which is similar to $H = LL^T$, and it is actually one step of the Rutishauser's LR method. If L is computed with pivoting, then L^TL is 'more diagonal' than H .
- The cost of Cholesky ($n^3/3$) much less than one sweep of Jacobi ($2n^3$ with fast rotations).
- In floating point

$$\tilde{L}\tilde{L}^T = H + \delta H, \quad \max_{i,j} \frac{|\delta H_{ij}|}{\sqrt{H_{ii}H_{jj}}} \leq \eta_C \lesssim n\varepsilon$$

Provable accuracy

Introduction

Scaling:
examplesNumerical
rank revealingEigenvalues
and singular
values

Jacobi method

Ein leichtes
VerfahrenOne-sided Jacobi
SVD

Floating point Jacobi

Provable accuracy

SVD computation in
floating pointAccurate
PSVD and
applicationsConcluding
remarksNow to the SVD of \tilde{L} :One sided Jacobi SVD $\tilde{L} V_1 V_2 \cdots V_k \cdots V_\ell \rightarrow \tilde{U} \tilde{\Sigma}$

In floating point

$$\bullet \quad \tilde{L} \leftarrow (((\tilde{L}_1 + \delta \tilde{L}_1) \hat{V}_1 + \delta \tilde{L}_2) \hat{V}_2 + \delta \tilde{L}_3) \hat{V}_3 + \cdots$$

- If $y = xV$, V rotation, x, y row vectors, then $\tilde{y} = (x + \delta x) \hat{V}$, \hat{V} orthogonal, $\|\delta x\| \leq 6\varepsilon \|x\|$.
- Hence, each row of $\delta \tilde{L}_j$ is ε small relative to the corresponding row of \tilde{L}_i . The \hat{V}_j with $j \neq i$ do not change the row norms of $\delta \tilde{L}_i$.
- At convergence, $\tilde{U} \tilde{\Sigma} = (\tilde{L} + \delta \tilde{L}) \hat{V}$, with $\tilde{\Sigma} = \text{diag}(\tilde{\sigma}_i)$, $\|\delta \tilde{L}(i, :)\| \leq O(n)\varepsilon \|\tilde{L}(i, :)\|$ for all i .
- $\tilde{\lambda}_i = \tilde{\sigma}_i^2$ are the eigenvalues of $(\tilde{L} + \delta \tilde{L})(\tilde{L} + \delta \tilde{L})^T$

Provable accuracy

$$(\tilde{L} + \delta\tilde{L})(\tilde{L} + \delta\tilde{L})^T = \tilde{L}\tilde{L}^T + \underbrace{\tilde{L}\delta\tilde{L}^T + \delta\tilde{L}\tilde{L}^T + \delta\tilde{L} + \delta\tilde{L}^T}_E$$

By Cauchy–Schwarz,

$$\begin{aligned} |E_{ij}| &\leq 2O(n\varepsilon)\|\tilde{L}(i, :)\|\|\tilde{L}(j, :)\| + O(\varepsilon^2)\|\tilde{L}(i, :)\|\|\tilde{L}(j, :)\| \\ &\approx (O(n\varepsilon) + O(\varepsilon^2))\sqrt{(\tilde{L}\tilde{L}^T)_{ii}(\tilde{L}\tilde{L}^T)_{jj}} \\ &\approx (O(n\varepsilon) + O(\varepsilon^2))\sqrt{H_{ii}H_{jj}}, \end{aligned}$$

$$\text{since } \tilde{L}\tilde{L}^T = H + \delta H, \quad \max_{i,j} \frac{|\delta H_{ij}|}{\sqrt{H_{ii}H_{jj}}} \leq \eta_C \lesssim n\varepsilon$$

So, we have the eigenvalues of

$$\tilde{L}\tilde{L}^T + E = H + \delta H + E = H + \Delta H, \quad \max_{i,j} \frac{|\Delta H_{ij}|}{\sqrt{H_{ii}H_{jj}}} \leq O(n\varepsilon)$$

Provable accuracy–conclusion

If $H \succ 0$ then

- The algorithm:
 1. Compute the Cholesky factorization $H = LL^T$;
 2. Compute $L = U\Sigma V^T$ using one–sided Jacobi SVD;
 3. Output: Set $\Lambda = \Sigma^2$; $H = U\Lambda U^T$

computes the eigenvalues and eigenvectors of H with entry–wise small backward error $\max_{i,j} \frac{|\Delta H_{ij}|}{\sqrt{H_{ii}H_{jj}}} \leq O(n\varepsilon)$.
- The forward error is $\max_i |\delta \lambda_i| / \lambda_i \leq O(n^2 \varepsilon) \|H_s^{-1}\|_2$.
- Most of the forward error comes from Step 1. Step 2. in floating point is as good as exact SVD.
- If Cholesky in Step 1 fails to compute L , then the matrix is entry–wise close to a non–definite matrix, and smallest eigenvalue can be lost due to symmetric tiny entry–wise perturbations.
- All computations in one $n \times n$ array.

Introduction

Scaling:
examplesNumerical
rank revealingEigenvalues
and singular
values

Jacobi method

Ein leichtes
VerfahrenOne–sided Jacobi
SVD

Floating point Jacobi

Provable accuracy

SVD computation in
floating pointAccurate
PSVD and
applicationsConcluding
remarks

SVD perturbation theory

Let $\text{rank}(A) = n \leq m$, $D = \text{diag}(\|A(:, i)\|)$, and

$$A \mapsto A + \delta A \implies \sigma_j \mapsto \sigma_j + \delta \sigma_j.$$

$$A + \delta A = (I + \delta A A^\dagger) A \implies \max_j \frac{|\tilde{\sigma}_j - \sigma_j|}{\sigma_j} \leq \|\delta A A^\dagger\|,$$

$$\|\delta A A^\dagger\| \leq \begin{cases} \frac{\|\delta A\|}{\|A\|} (\|A^\dagger\| \|A\|) = \epsilon \cdot \kappa(A), \\ \|\delta A D^{-1}\| \|(A D^{-1})^\dagger\|. \end{cases}$$

$$\|\delta A D^{-1}\| \leq \sqrt{n} \max_j \frac{\|A(:, j)\|}{\|A(:, j)\|} \leq \sqrt{n} \epsilon;$$

$$\|(A D^{-1})^\dagger\| \equiv \|A_s^\dagger\| \leq \sqrt{n} \min_{\Delta = \text{diag}} \kappa(A \Delta)$$

Possible: $\|A_s^\dagger\| \ll \kappa(A)$; always $\|A_s^\dagger\| \leq \sqrt{n} \kappa(A)$.

Jacobi SVD: $\|A_s^\dagger\| \longrightarrow$ more accurate .

bidagonal SVD: $\kappa(A) \longrightarrow$ less accurate ,

bidagonalization provokes $\kappa(A)$.

Jacobi++ SVD: $A = D_1 C D_2 \rightarrow D_1 (C + \delta C) D_2$.

Introduction

Scaling:
examplesNumerical
rank revealingEigenvalues
and singular
values

Jacobi method

Ein leichtes
VerfahrenOne-sided Jacobi
SVD

Floating point Jacobi

Provable accuracy

SVD computation in
floating pointAccurate
PSVD and
applicationsConcluding
remarks

QRF preprocessor for Jacobi

$A = QR$; $[\tilde{Q}, \tilde{R}] = \text{qr}(A)$, \tilde{Q}, \tilde{R} computed.

Backward error analysis:

$$(\exists \delta A) \quad (\exists \hat{Q}, \quad \hat{Q}^T \hat{Q} = I) \quad A + \delta A = \hat{Q} \tilde{R},$$

$$\|\delta A(:, i)\| \leq \epsilon_1 \|A(:, i)\|, \quad i = 1, \dots, n.$$

Perturbation analysis: $\sigma_i(\tilde{R}) = \sigma_i((I + \delta A A^\dagger)A)$

$$1 - \|\delta A A^\dagger\| \leq \frac{\sigma_i(\tilde{R})}{\sigma_i(A)} \leq 1 + \|\delta A A^\dagger\|, \quad \text{for all } i.$$

Let $A = A_S D$, $D = \text{diag}(\|A(:, i)\|)$.

$$\|\delta A A^\dagger\| = \|\delta A D^{-1} (A D^{-1})^\dagger\| \leq \sqrt{n} \max_i \frac{\|\delta A(:, i)\|}{\|A(:, i)\|} \|A_S^\dagger\|$$

$$\|A_S^\dagger\| \leq \kappa(A_S) \leq \sqrt{n} \min_{\Delta=\text{diag}} \kappa(A\Delta)$$

If $\kappa(A_S)$ is moderate, then $SVD(\tilde{R})$ is OK for the $SVD(A)$.

Introduction

Scaling:
examplesNumerical
rank revealingEigenvalues
and singular
values

Jacobi method

Ein leichtes
VerfahrenOne-sided Jacobi
SVD

Floating point Jacobi

Provable accuracy

SVD computation in
floating pointAccurate
PSVD and
applicationsConcluding
remarks

Strong backward stability

Jacobi SVD(\tilde{R}): $\tilde{R}^T V = U \Sigma$. Computed:

$[\tilde{U}, \tilde{V}, \tilde{\Sigma}] = \text{JacobiSVD}(\tilde{R}^T)$. Jacobi rotations \tilde{V} such that

$$\max_{i \neq j} \left| \cos \angle((\tilde{U} \tilde{\Sigma}) e_i, (\tilde{U} \tilde{\Sigma}) e_j) \right| \leq O(n) \mathbf{u}$$

Error analysis:

$$(\exists \delta \tilde{R}) \quad (\exists \hat{V}, \quad \hat{V}^T \hat{V} = I) \quad (\tilde{R} + \delta \tilde{R})^T \hat{V} = (\tilde{U} \tilde{\Sigma})$$

$$\|\delta \tilde{R}(:, i)\| \leq \epsilon_2 \|\tilde{R}(:, i)\|, \quad i = 1, \dots, n.$$

Finally,

$$\begin{aligned} \tilde{R} + \delta \tilde{R} &= \hat{Q}^T (A + \delta A) + \hat{Q}^T \hat{Q} \delta \tilde{R} \\ &= \hat{Q}^T (A + \underbrace{\delta A + \hat{Q} \delta \tilde{R}}_{\Delta A}) \end{aligned}$$

where $\|\Delta A(:, i)\| \leq (\epsilon_1 + \epsilon_2(1 + \epsilon_1)) \|A(:, i)\|$ for all i , and the SVD is $(A + \Delta A)^T \hat{Q} \hat{V} = \tilde{U} \tilde{\Sigma}$. Very nice and simple.

Accurate.

+definite $HMx = \lambda x$

Introduction

Scaling:
examplesNumerical
rank revealingEigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsAccurate PSVD
RRD of structured
matrices
Rational
approximationConcluding
remarks

Entry-wise backward stability also possible for $HMx = \lambda x$.

- Use contragredient scaling $H := DHD$, $M := D^{-1}MD^{-1}$ to get all $M_{ii} = 1$. Here $D = \text{diag}(\sqrt{M_{ii}})_{i=1}^n$.
- Cholesky f. $H := P^T H P = L_H L_H^T$, $M = L_M L_M^T$
- $HM = PL_H L_H^T P^T L_M L_M^T = L_M^{-T} (L_M^T P L_H L_H^T P^T L_M) L_M^T$
- $HM = L_M^{-T} (A A^T) L_M$, $A = L_M^T P L_H$, $L_H = L_{H,s} D_H$
- Compute $A = (L_M^T P) L_H$. (No fast matrix-multiply allowed. Must pay $O(n^3)$.)
- Compute the SVD $A = U \Sigma V^T$ using the Jacobi SVD ($AV = U \Sigma$, $AA^T = U \Sigma^2 U^T$).
- Assemble: $T = D L_M^{-T} U \Sigma^{1/2}$.
- It holds $T^{-1} H T^{-T} = T^T M T = \Sigma$

Accurate solution

The algorithm solves

$$(H + \delta H)(M + \delta M)x = \tilde{\lambda}x$$

exactly, with symmetric $\delta H, \delta M$,

$$\frac{|\delta H_{ij}|}{\sqrt{H_{ii}H_{jj}}} \leq f(n) \cdot \varepsilon, \quad \frac{|\delta M_{ij}|}{\sqrt{M_{ii}M_{jj}}} \leq g(n, L_{H,s}) \cdot \varepsilon, \quad 1 \leq i, j \leq n$$

$$\frac{|\delta \lambda|}{\lambda} \leq h(n)(\|H_s^{-1}\| + \|M_s^{-1}\|) \cdot \varepsilon, \quad \varepsilon = \text{eps}.$$

$$H_s = \text{diag}(H)^{-1/2} H \text{diag}(H)^{-1/2}, \quad \kappa_2(H_s) \leq n \min_{D=\text{diag}} \kappa_2(DHD)$$

All λ 's stable IFF $\|H_s^{-1}\|$ and $\|M_s^{-1}\|$ moderate.

We have optimal accuracy.

Implicit diagonalization of HM is actually computing the SVD of a product of two matrices, $BC^T = U\Sigma V^T$.

$A = BC^T = U\Sigma V^T$, B , C full column rank

$$BC^T = \begin{pmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{pmatrix} \begin{pmatrix} \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare \end{pmatrix}$$

- $A = \text{yxGEMM}(B, C^T)$ fastest matrix multiply
- $\text{CALL yxGESDD}(A)$ fastest SVD
- $\begin{pmatrix} 1 & \epsilon \\ -1 & \epsilon \end{pmatrix} \begin{pmatrix} 2 & 2 \\ 2 & 1 \end{pmatrix} = \begin{pmatrix} 2+2\epsilon & 2+\epsilon \\ -2+2\epsilon & -2+\epsilon \end{pmatrix} \approx \begin{pmatrix} 2 & 2 \\ -2 & -2 \end{pmatrix}$
- $U_2 B U_1^T U_1 C^T U_3 \rightsquigarrow \Sigma$, U_i orthogonal

Introduction

Scaling:
examplesNumerical
rank revealingEigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsAccurate PSVD
RRD of structured
matrices
Rational
approximationConcluding
remarks

- $\begin{pmatrix} 1 & \epsilon \\ -1 & \epsilon \end{pmatrix} \begin{pmatrix} 2 & 2 \\ 2 & 1 \end{pmatrix} \approx \begin{pmatrix} 2 & 2 \\ -2 & -2 \end{pmatrix}$, ϵ not happy
- $\begin{pmatrix} 1 & \epsilon \\ -1 & \epsilon \end{pmatrix} U_1^T U_1 \begin{pmatrix} 2 & 2 \\ 2 & 1 \end{pmatrix}$, $U_1 = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}$
- $U_1 C^T = \begin{pmatrix} \sqrt{8} & \frac{\sqrt{18}}{2} \\ 0 & -\frac{\sqrt{2}}{2} \end{pmatrix}$
- $BU_1^T = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 + \epsilon & -1 + \epsilon \\ -1 + \epsilon & 1 + \epsilon \end{pmatrix} \approx \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$
- ϵ happy because U_1 is orthogonal ?!
- backward errors: $\|\delta B\| \lesssim \text{eps} \|B\|$, $\|\delta C\| \lesssim \text{eps} \|C\|$
- Is that the best we can do?

PSVD: $SVD(BC^T)$

How to compute the SVD of a product of two matrices,
 $BC^T = U\Sigma V^T$, accurately?

$$B\Delta_B^{-1}\Delta_B C^T = \underbrace{\begin{pmatrix} \text{green squares} \end{pmatrix}}_{\text{unit columns}} \begin{pmatrix} \text{blue squares} \end{pmatrix}$$

- $C\Delta_B P = Q \begin{pmatrix} R \\ 0 \end{pmatrix}$; $BC^T = (B\Delta_B^{-1}P) \begin{pmatrix} R^T & 0 \end{pmatrix} Q^T$;
- $A = (B\Delta_B^{-1}P)R^T$; $R^T = \begin{pmatrix} \text{blue square} & & \\ \text{yellow square} & \text{blue square} & \\ \text{yellow square} & \text{yellow square} & \text{cyan square} \end{pmatrix} = \text{well.cond} \times \text{diag.}$
- $[U, \Sigma, V_1] = \text{SVD}(A)_{\text{Jacobi}}$; $V = Q \begin{pmatrix} V_1 & 0 \\ 0 & I_{n-p} \end{pmatrix}$

And what about BPR^T ?

$$BPT^T \equiv \begin{pmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{pmatrix} \begin{pmatrix} \blacksquare & & \\ \blacksquare & \blacksquare & \\ \blacksquare & \blacksquare & \blacksquare \end{pmatrix}$$

$$\text{Consider } (a_1 \ a_2 \ a_3) = (b_1 \ b_2 \ b_3) \begin{pmatrix} \ell_{11} & 0 & 0 \\ \ell_{21} & \ell_{22} & 0 \\ \ell_{31} & \ell_{32} & \ell_{33} \end{pmatrix}$$

$$\tilde{a}_3 = (b_3 + \delta b_3) \ell_{33}$$

$$\tilde{a}_2 = (b_2 + \delta_2 b_2) \ell_{22} + (b_3 + \delta_2 b_3) \ell_{32}$$

$$= (b_2 + \delta_2 b_2) \ell_{22} + (b_3 + \delta b_3 - \delta b_3 + \delta_2 b_3)$$

$$= (b_2 + \delta_2 b_2 + (\delta_2 b_3 - \delta b_3) \frac{\ell_{32}}{\ell_{22}}) \ell_{22} + (b_3 + \delta b_3) \ell_{33}$$

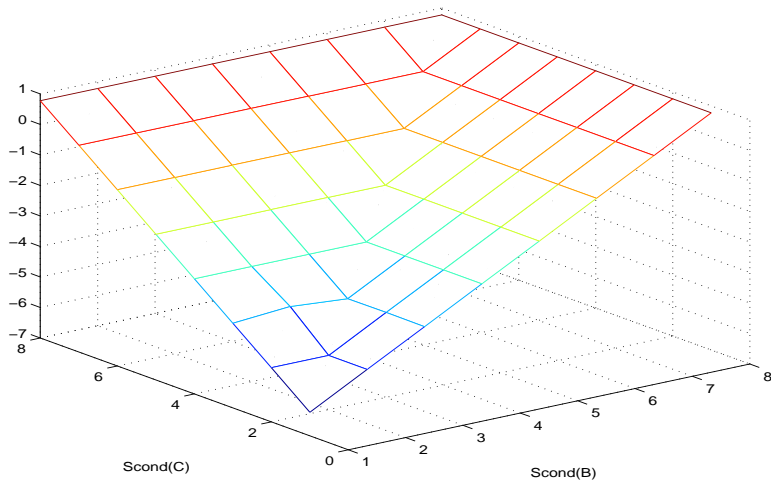
$$= (b_2 + \delta b_2) \ell_{22} + (b_3 + \delta b_3) \ell_{33}$$

$$(\tilde{a}_2 \ \tilde{a}_3) = (b_2 + \delta b_2 \ b_3 + \delta b_3) L, \quad L = \tilde{R}^T$$

Backward stability

- $C = Q \begin{pmatrix} R \\ 0 \end{pmatrix};$
 - $C + \delta C = \tilde{Q} \begin{pmatrix} \tilde{R} \\ 0 \end{pmatrix};$
 - $\|\delta C(:, i)\| \leq \epsilon \|C(:, i)\|$, for all columns i
- $A = BR^T;$
 - $\tilde{A} = (B + \delta B)\tilde{R}^T,$
 - $\|\delta B(:, i)\| \leq \epsilon \|B(:, i)\|$, for all columns i
- $(B + \delta B)(C + \delta C)^T = (I + \delta BB^\dagger)BC^T(I + \delta CC^\dagger)^T$
- $B = B_{scaled}D$, $scond(B) = cond(B_{scaled})$

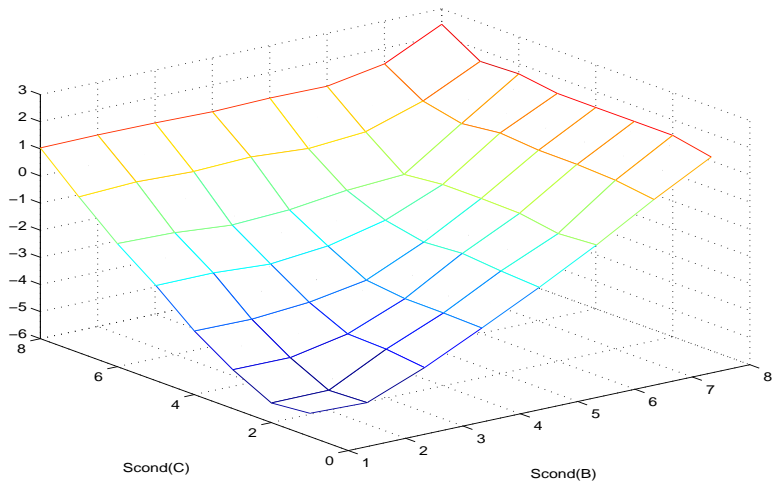
Theoretical accuracy



theory:

$$\max_i \frac{|\tilde{\sigma}_i - \sigma_i|}{\sigma_i} \leq f(m, p, n) \cdot \epsilon \cdot \max\{\text{scnd}(B), \text{scnd}(C)\}$$

Measured accuracy



theory: $\max_i \frac{|\tilde{\sigma}_i - \sigma_i|}{\sigma_i}$; in $(0.3, 46) \times$ theory

Rank Revealing Decomposition

Introduction

Scaling:
examplesNumerical
rank revealingEigenvalues
and singular
values

Jacobi method

Accurate
PSVD and
applicationsAccurate PSVD
RRD of structured
matricesRational
approximationConcluding
remarks

In a +60 pages LAA paper Demmel, Drmač, Gu, Eisenstat, Slapničar, Veselić (DGESVD paper) noted that some classes of matrices allow so called Rank Revealing Decomposition (RRD),

$$P_1 A P_2 = L D U, \quad P_1, P_2 \text{ permutations,}$$

where D is diagonal, and L and U are well conditioned. Moreover, L, D, U can be computed in a forward stable way. An example of a RRD of A is obtained by non-standard Gaussian eliminations using certain structural properties of A . More examples by Demmel and Koev. Then, we can use the accurate PSVD algorithm and get the SVD of LDU .

Example: Cauchy matrix $C_{ij} = 1/(x_i + y_j)$
(displacement rank one, $XC + CY = d_1 d_2^T$)

Cauchy matrices

$$\det(C) = \frac{\prod_{i < j} (x_j - x_i)(y_j - y_i)}{\prod_{i,j} (x_i + y_j)}$$

Can get accurate LDU at high cost, $O(n^5)$. Then Demmel reduced it to the usual $O(n^3)$ using the recursive structure of the Schur complement.

$$\begin{pmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{pmatrix} = \begin{pmatrix} I & 0 \\ C_{21}C_{11}^{-1} & I \end{pmatrix} \begin{pmatrix} C_{11} & C_{12} \\ 0 & S^{(k)} \end{pmatrix}$$

$$S_{ij}^{(k)} = S_{ij}^{(k-1)} \frac{(x_i - x_k)(y_j - y_k)}{(x_k + y_j)(x_i + y_k)}$$

Straightforward extension to Cauchy-like matrices $D_1 C D_2$, D_i diagonal. Simplified for symmetric positive definite cases.

An illustration

After computing $C = LDU$, one applies accurate Jacobi PSVD to the product $(LD)U$. **All forward stable, but the spectrum is ill-conditioned!**

An illustration of the power of this algorithm is the example of 100×100 Hilbert matrix H_{100} . Computation done by Demmel:

- The singular values of H_{100} range over 150 orders of magnitude and are computed using the package Mathematica with 200–decimal digit software floating point arithmetic. The computed singular values are rounded to 16 digits and used as reference values.
- The singular values computed in IEEE double precision floating–point ($\varepsilon \approx 10^{-16}$) by the Jacobi PSVD agree with the reference values with relative error less than $34 \cdot \varepsilon$.

Rational approximation

New highly accurate NLA algorithms open new possibilities in other computational tasks.

For instance, Haut and Beylkin (2011) used Adamyan–Arov–Krein theory to show that nearly L^∞ –optimal rational approximation on $|z| = 1$ of

$$f(z) = \sum_{i=1}^n \frac{\alpha_i}{z - \gamma_i} + \sum_{i=1}^n \frac{\overline{\alpha_i} z}{1 - \overline{\gamma_i} z} + \alpha_0$$

with $\max_{|z|=1} |f(z) - r(z)| \rightsquigarrow \min$,

$$r(z) = \sum_{i=1}^m \frac{\beta_i}{z - \eta_i} + \sum_{i=1}^m \frac{\overline{\beta_i} z}{1 - \overline{\eta_i} z} + \alpha_0$$

is numerically feasible if one can compute the con–eigenvalues and con–eigenvectors

$$Cu = \lambda \overline{u}, \quad C_{ij} = \frac{\sqrt{\alpha_i} \sqrt{\overline{\alpha_j}}}{\gamma_i^{-1} - \overline{\gamma_j}} \rightsquigarrow \frac{\alpha_i \overline{\alpha_j}}{1 - \gamma_i \overline{\gamma_j}}$$

Con-eigenvalues

Here $C = \left(\frac{\sqrt{\alpha_i} \sqrt{\alpha_j}}{\gamma_i^{-1} - \gamma_j} \right)$ is positive definite Cauchy matrix C .

The con-eigenvalue problem $Cu = \lambda \bar{u}$ is equivalent to solving

$$\overline{C}Cu = |\lambda|^2 u,$$

where C is factored as $C = XD^2X^*$. The problem reduces to computing the SVD of the product $G = DX^T XD$. Accurate SVD via the PSVD based on the Jacobi SVD. Haut and Beylkin tested the accuracy with $\kappa_2(C) > 10^{200}$ and using **Mathematica with 300 hundred digits for reference values**. Over 500 test examples of size 120, the maximal error in **IEEE 16 digit arithmetic** ($\varepsilon \approx 2.2 \cdot 10^{-16}$) was

$$\frac{|\tilde{\lambda}_i - \lambda_i|}{|\lambda_i|} < 5.2 \cdot 10^{-12}, \quad \frac{\|\tilde{u}_i - u_i\|_2}{\|u_i\|_2} < 5.4 \cdot 10^{-12}.$$

- 1 Introduction
- 2 Scaling: examples
- 3 Numerical rank revealing
- 4 Eigenvalues and singular values
- 5 Jacobi method
- 6 Accurate PSVD and applications
- 7 Concluding remarks**

Concluding remarks

- Ill-conditioning can be artificial, an artifact of a particular algorithm, and not the underlying problem. In many cases accurate computation is possible, despite high classical condition numbers.
- Backward stability is often used to justify the result. Structured backward error can yield better results.
- Using only orthogonal transformations does not automatically guarantee good results.
- Users from applied sciences and engineering – often not interested in math details, just solutions, software. Need robust reliable and efficient numerical software. Is trading accuracy for speed avoidable?
- Challenging problems for numerical linear algebra. Higher standards for new algorithms.