

LECTURE 16BACKWARD ERROR ANALYSIS (PART II)

A backward stable algorithm is accurate if and only if the problem is well-conditioned.

THM (Forward and Backward Error)

Suppose there exists a δx such that

$$\tilde{f}(x) = f(x + \delta x).$$

Then

$$\frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|} \leq \tilde{\kappa}_f \frac{\|\delta x\|}{\|x\|}$$

RELATIVE FORWARD ERROR RELATIVE BACKWARD ERROR

where $\delta := \|\delta x\|$.

PROOF

$$\frac{\|\tilde{f}(x) - f(x)\| / \|f(x)\|}{\|\delta x\| / \|x\|} = \frac{\|f(x + \delta x) - f(x)\| / \|f(x)\|}{\|\delta x\| / \|x\|} \quad \text{(since } \tilde{f}(x) = f(x + \delta x)\text{)}$$

①

$$\leq \underbrace{\sup_{\|\tilde{\delta x}\| \leq \delta} \frac{\|f(x+\tilde{\delta x}) - f(x)\|}{\|\tilde{\delta x}\| / \|x\|}}_{\tilde{K}_\delta}$$

\implies

$$\|\tilde{f}(x) - f(x)\| / \|f(x)\| \leq \tilde{K}_\delta (\|\delta x\| / \|x\|) \quad \square$$

COROLLARY

Suppose \tilde{f} is a backward stable algorithm and $\tilde{K} > 0$. Then

$$\frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|} = O(\tilde{K} \epsilon_{mach})$$

PROOF

From the previous thm

$$\frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|} \leq \tilde{K}_\delta \underbrace{c \epsilon_{mach}}_{\left(\text{since } \frac{\|\delta x\|}{\|x\|} = O(\epsilon_{mach})\right)}$$

where $\delta := \|\delta x\| \leq c \epsilon_{mach} \|x\|$.

But since $\tilde{K} = \lim_{\delta \rightarrow 0} \tilde{K}_\delta$, for small δ we have $\tilde{K}_\delta \leq d\tilde{K}$ where $d > 1$ is a constant. Therefore for ϵ_{mach} small enough (ensuring δ is small)

$$\frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|} \leq cd\tilde{K}\epsilon_{mach}$$

□

BACKWARD ERROR ANALYSIS FOR MATRIX-VECTOR PRODUCT

View the product Ax as a function of $A \in \mathbb{C}^{n \times n}$ for a fixed $x \in \mathbb{C}^n$.

$$f: \mathbb{C}^{n \times n} \rightarrow \mathbb{C}^n, \quad f(A) = Ax$$

Recall that

$$\tilde{K} \leq \|A\| \|A^{-1}\|$$

THM (Backward Error of **Matrix-Vector Product**)
 Let $\tilde{b} = \tilde{f}(A)$ be the computed value of $b = f(A) = Ax$ in a floating point arithmetic.

Then

$$\tilde{b} = \tilde{f}(A) = (A + \delta A)x$$

for some δA such that

$$\frac{\|\delta A\|_1}{\|A\|_1} \leq n \epsilon_{mach} + O(\epsilon_{mach}^2)$$

PROOF

Consider the j th entry of \tilde{b}

$$\begin{aligned} \tilde{b}_j &= (a_{j1} \otimes x_1) \oplus (a_{j2} \otimes x_2) \oplus \dots \oplus (a_{jn} \otimes x_n) \\ &= (a_{j1} x_1 (1 + \epsilon_1)) \oplus (a_{j2} x_2 (1 + \epsilon_2)) \oplus \dots \oplus (a_{jn} x_n (1 + \epsilon_n)) \\ &\checkmark = ((a_{j1} x_1 (1 + \epsilon_1) + a_{j2} x_2 (1 + \epsilon_2)) (1 + \tilde{\epsilon}_1)) \oplus \dots \oplus (a_{jn} x_n (1 + \epsilon_n)) \\ &= (a_{j1} x_1 (1 + \epsilon_1 + \tilde{\epsilon}_1 + O(\epsilon_{mach}^2)) + a_{j2} x_2 (1 + \epsilon_2 + \tilde{\epsilon}_1 + O(\epsilon_{mach}^2))) \\ &\quad \oplus \dots \oplus (a_{jn} x_n (1 + \epsilon_n)) \\ &= (a_{j1} (1 + \epsilon_{11} + \epsilon_{12} + \dots + \epsilon_{1n} + O(\epsilon_{mach}^2))) x_1 \\ &\quad + \dots + \\ &\quad (a_{jn} (1 + \epsilon_{n1} + \epsilon_{n2} + \dots + \epsilon_{nn} + O(\epsilon_{mach}^2))) x_n \end{aligned}$$

where ϵ_{ij} and $\epsilon_k, \tilde{\epsilon}_k$ are at most ϵ_{mach} in absolute value.

Consequently

$$\tilde{b}_j = \sum_{k=1}^n (a_{jk} + \delta a_{jk}) x_k$$

with

$$\delta a_{jk} = a_{jk} (\epsilon_{k1} + \dots + \epsilon_{kn} + O(\epsilon_{mach}^2))$$

$$\implies |\delta a_{jk}| \leq |a_{jk}| n \epsilon_{mach} + O(\epsilon_{mach}^2).$$

It follows that

$$\tilde{b} = (A + \delta A) x$$

where

$$\|\delta A\|_1 = \max_{1 \leq k \leq n} \|\delta a_k\|_1$$

$$= \max_{1 \leq k \leq n} \sum_{j=1}^n |\delta a_{jk}|$$

$$\leq \max_{1 \leq k \leq n} \sum_{j=1}^n (|a_{jk}| n \epsilon_{mach} + O(\epsilon_{mach}^2))$$

$$= n \epsilon_{mach} \max_{1 \leq k \leq n} \sum_{j=1}^n |a_{jk}| + O(\epsilon_{mach}^2)$$

$$= n \epsilon_{mach} \|A\|_1 + O(\epsilon_{mach}^2)$$

□ (5)

Forward error of the matrix-vector product

$$\frac{\|\tilde{b} - b\|_1}{\|b\|_1}$$

FORWARD RELATIVE
ERROR

\leq

$\tilde{\kappa}_s$

$$\frac{\|\delta A\|_1}{\|A\|_1}$$

BACKWARD RELATIVE
ERROR

$$\leq \|A\|_1 \|A^{-1}\|_1 (n \epsilon_{mach} + O(\epsilon_{mach}^2))$$

\implies

$$\frac{\|\tilde{b} - b\|_1}{\|b\|_1} = O(\kappa_1(A) \epsilon_{mach})$$

BACKWARD ERROR ANALYSIS FOR LEAST-SQUARES PROBLEM

Least squares problem

Given $A \in \mathbb{C}^{m \times n}$ (with $m \geq n$) and $b \in \mathbb{C}^m$. Find $x \in \mathbb{C}^n$ such that $\|Ax - b\|_2$ is as small as possible.

Relative condition number of x w.r.t. A

$$\tilde{\kappa} = \kappa(A) + \frac{\kappa(A)^2 \tan \theta}{n}$$

(6)

where

$$\theta = \arccos\left(\frac{\|Ax\|_2}{\|b\|_2}\right) \text{ and } \kappa = \frac{\|A\| \|x\|}{\|Ax\|}$$

THM (Backward Stability of QR based algorithm)

Suppose the LSP solved using the QR factorization of A computed by Householder reflectors. Then the computed solution \hat{x} solves

Find \hat{x} so that $\|(A + \delta A)\hat{x} - b\|_2$ is as small as possible

where

$$\frac{\|\delta A\|}{\|A\|} = O(\epsilon_{mach}).$$

Forward error of LSP

$$\begin{aligned} \frac{\|\hat{x} - x\|}{\|\hat{x}\|} &= O(\tilde{\kappa} \epsilon_{mach}) \\ &= O\left(\left(\kappa(A) + \frac{\kappa^2(A) \tan \theta}{\kappa}\right) \epsilon_{mach}\right) \end{aligned}$$