

LECTURE 17BACKWARD ERROR ANALYSIS

$$f: V \rightarrow W$$

V, W are vector spaces

$\tilde{f}(x)$: computed value
of $f(fl(x))$ by
the numerical algorithm.

Forward Error

Absolute error

$$\|\tilde{f}(x) - f(x)\|$$

Relative error

$$\frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|}$$

Accurate algorithm (for all x)

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

REMARK

(*) means there exists a constant c
(independent of x)

$$\frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|} \leq c \epsilon_{\text{mach}}$$

for all $\epsilon_{\text{mach}} > 0$ close to zero.

Forward error directly is difficult to analyze.

Backward Error

DEFN

Suppose there exists a $\delta x \in V$ s.t.

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

Then

(i) $\|\delta x\|$ is the absolute backward error.

(ii) $\|\delta x\| / \|x\|$ is the relative backward error.

EXAMPLE

Matrix-vector product

$$b = f(A) = Ax \quad (A \in \mathbb{C}^{n \times n})$$

\tilde{b} : computed product

Find the backward error in an IEEE floating point arithmetic.

SOLN

$$\tilde{b}_j = \bigoplus_{k=1}^n (a_{jk} \otimes x_k)$$

$$= \bigoplus_{k=1}^n (a_{jk} \times x_k) (1 + \epsilon_k) \quad \left(\begin{array}{l} \text{where} \\ |\epsilon_k| \leq \epsilon_{\text{mach}} \end{array} \right)$$

$$= (a_{j1} \times x_1) (1 + \epsilon_1) \oplus \left(\bigoplus_{k=2}^n (a_{jk} \times x_k) (1 + \epsilon_k) \right)$$

$$= (a_{j1} \times x_1) (1 + \epsilon_1) (1 + \tilde{\epsilon}_1) \quad \left(\begin{array}{l} \text{where} \\ |\tilde{\epsilon}_1| \leq \epsilon_{\text{mach}} \end{array} \right)$$

$$(**) \quad \bigoplus_{k=2}^n (a_{jk} \times x_k) (1 + \epsilon_k) (1 + \tilde{\epsilon}_1)$$

Here note that (for $k=1, \dots, n$)

$$(1 + \epsilon_k) (1 + \tilde{\epsilon}_1) = 1 + \epsilon_k + \tilde{\epsilon}_1 + \epsilon_k \tilde{\epsilon}_1$$

$$= 1 + \epsilon_k + \tilde{\epsilon}_1 + O(\epsilon_{\text{mach}}^2) \quad (3)$$

Repeat (**) by replacing

⊕ with +.

$$\begin{aligned} \tilde{b}_j &= (a_{j1} \times X_1) (1 + \epsilon_1) (1 + \tilde{\epsilon}_1) \\ &\quad + \\ &\quad (a_{j2} \times X_2) (1 + \epsilon_2) (1 + \tilde{\epsilon}_1) (1 + \tilde{\epsilon}_2) \\ &\quad + \\ &\quad \vdots \\ &\quad + \\ &\quad (a_{jn-1} \times X_{n-1}) (1 + \epsilon_1) (1 + \tilde{\epsilon}_1) \dots (1 + \tilde{\epsilon}_{n-1}) \\ &\quad + \\ &\quad (a_{jn} \times X_n) (1 + \epsilon_1) (1 + \tilde{\epsilon}_1) \dots (1 + \tilde{\epsilon}_{n-1}) \end{aligned}$$

Above note

$$\begin{aligned} (1 + \epsilon_1) (1 + \tilde{\epsilon}_1) \dots (1 + \tilde{\epsilon}_{\ell_1}) &= (1 + \epsilon_1 + \tilde{\epsilon}_1 + \dots + \tilde{\epsilon}_{\ell_1} + O(\epsilon_{mach}^2)) \\ &\leq (1 + \ell \epsilon_{mach} + O(\epsilon_{mach}^2)) \end{aligned}$$

Consequently

~~$$\tilde{b}_j = \sum_{k=1}^{n-1} X_k \cdot a_{jk} (1 + \epsilon_1) \prod_{l=1}^k (1 + \tilde{\epsilon}_l)$$~~

$$\begin{aligned} \tilde{b}_j &= \sum_{k=1}^{n-1} X_k \cdot a_{jk} (1 + \epsilon_1) \prod_{l=1}^k (1 + \tilde{\epsilon}_l) \\ &\quad + \\ &\quad X_n \cdot a_{jn} (1 + \epsilon_1) \prod_{l=1}^{n-1} (1 + \tilde{\epsilon}_l) \end{aligned}$$

$$= \sum_{k=1}^n x_k \cdot a_{jk} (1 + \delta \epsilon_k) \left(\text{where } |\delta \epsilon_k| \leq n \epsilon_{\text{mach}} + O(\epsilon_{\text{mach}}^2) \right)$$

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

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

Summary

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

where

$$* \quad |\delta a_{jk}| \leq n \epsilon_{\text{mach}} |a_{jk}| + O(\epsilon_{\text{mach}})$$

* equivalently

$$\|\delta a_k\|_1 = \sum_{j=1}^n |\delta a_{jk}|$$

$$\leq n \epsilon_{\text{mach}} \sum_{j=1}^n |a_{jk}|$$

$$\leq n \epsilon_{\text{mach}} \|a_k\|_1$$

* equivalently

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

$$\leq \max_{k=1, \dots, n} n \epsilon_{\text{mach}} \|a_k\|_1$$

$$= n \epsilon_{\text{mach}} \|A\|_1$$

Backward error

$$\frac{\|\delta A\|_1}{\|A\|_1} \leq n \epsilon_{\text{mach}}$$

EXAMPLE

Least squares problem, \hat{x} is s.t.

$$\|A\hat{x} - b\|_2 = \underset{x \in \mathbb{C}^n}{\text{minimize}} \|Ax - b\|_2$$

$$(A \in \mathbb{C}^{m \times n}, b \in \mathbb{C}^m \text{ with } m \geq n)$$

\tilde{x} : computed value of \hat{x} using the QR factorization by Householder reflectors.

Assuming A is full-rank it can be shown that

$$\|A\tilde{x} - b\|_2 = \underset{x \in \mathbb{C}^n}{\text{minimize}} \|(A + \delta A)x - b\|_2$$

where

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

DEFN (Backward Stability)

An algorithm is called backward stable if for all x there exist δx s.t.

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

with

$$\|\delta x\| / \|x\| = O(\epsilon_{\text{mach}}).$$

EXAMPLES

- * Matrix-vector product is backward stable.
- * Least squares problem using QR factor. by Householder reflectors is backward stable.

Backward Error Analysis

Relates the backward error with forward error.

THM

Suppose there exists δx s.t.

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

Then

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

where

$$\tilde{\kappa}_\delta := \sup_{\|\hat{\delta}_x\| \leq \delta} \frac{(\|f(x + \hat{\delta}_x) - f(x)\| / \|f(x)\|)}{(\|\hat{\delta}_x\| / \|x\|)}$$

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

PROOF

$$\begin{aligned} \frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|} &= \frac{\|f(x + \delta x) - f(x)\|}{\|f(x)\|} \\ &= \frac{\|f(x + \delta x) - f(x)\| / \|f(x)\|}{\|\delta x\| / \|x\|} \cdot \frac{\|\delta x\|}{\|x\|} \\ &\leq \left(\sup_{\|\hat{\delta}_x\| \leq \delta} \frac{\|f(x + \hat{\delta}_x) - f(x)\| / \|f(x)\|}{\|\hat{\delta}_x\| / \|x\|} \right) \cdot \frac{\|\delta x\|}{\|x\|} \end{aligned}$$

□

EXAMPLES

① Backward Error Analysis of Matrix-Vector Product

$$b = f(A) = Ax$$

\tilde{b} : computed product

δ -relative condition number

$$\tilde{\kappa}_\delta = \sup_{\|\delta A\| \leq \delta} \frac{\|f(A+\delta A) - f(A)\|}{\|\delta A\|} \cdot \frac{\|A\|}{\|f(A)\|}$$
$$\leq \|A\| \|A^{-1}\|$$

Backward stability

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

for some δA s.t.

$$\frac{\|\delta A\|}{\|A\|} \leq n \epsilon_{mach}$$

Forward error

$$\frac{\|\tilde{b} - b\|}{\|b\|} = \frac{\|f(A + \delta A) - f(A)\|}{\|\delta A\|}$$

$$\stackrel{\text{(THEM ABOVE)}}{\leq} \tilde{\kappa}_\delta \frac{\|\delta A\|}{\|A\|}$$

$$= \underbrace{\|A\| \|A^{-1}\|}_{\kappa_1(A)} n \epsilon_{mach}$$

② Backward Error Analysis of LSP

$$\|A\hat{x} - b\|_2 = \underset{x \in \mathbb{C}^n}{\text{minimize}} \|Ax - b\|_2$$

\tilde{x} : computed value using
QR factor by HH reflectors

δ - relative condition number

$$\tilde{\kappa}_\delta = \sup_{\|\delta A\| \leq \delta} \frac{\|\tilde{x}(A + \delta A) - \tilde{x}(A)\|}{\|\tilde{x}(A)\|} \frac{\|\delta A\| / \|A\|}{\|\delta A\| / \|A\|}$$

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

where

$$\kappa(A) = \|A\| \|A^+\|$$

$$\theta = \arccos\left(\frac{y}{b}\right) \quad \left(\begin{array}{l} \text{where } y \text{ is} \\ \text{the orthogonal} \\ \text{proj. of } b \text{ onto} \\ \text{Range}(A) \end{array} \right)$$

$$n = \frac{\|\tilde{x}\| \|A\|}{\|A\tilde{x}\|}$$

Backward stability (assuming A is full rank)

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

for some δA s.t.

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

Forward Error

$$\frac{\|\tilde{x} - \hat{x}\|}{\|\hat{x}\|} \leq \tilde{\kappa}_\delta \frac{\|\delta A\|}{\|A\|}$$

$$\leq c \left(\kappa(A) + \kappa(A) \frac{\tan^2 \theta}{n} \right) \epsilon_{\text{mach}}$$

for some constant c .