

7.6 Accuracy of a Backward Stable Algorithm

Theorem: Suppose a backward stable algorithm is used to solve $f : x \rightarrow y$ with condition number κ on a computer satisfying

$$fl(x \odot y) = (x \odot y)(1 + \delta) \quad |\delta| \leq \epsilon_{\text{machine}} \quad (7.17)$$

then

$$\frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|} = O(\kappa(x) \cdot \epsilon_{\text{machine}}) \quad (7.18)$$

Proof: By definition $\tilde{f}(x) = f(\tilde{x})$ for $\tilde{x} : \frac{\|\tilde{x} - x\|}{\|x\|} = O(\epsilon_{\text{machine}})$.

(Continued on next page.)

$$\kappa(x) = \lim_{\delta \rightarrow 0} \sup_{\|\delta x\| \leq \delta} \frac{\frac{\|\delta f\|}{\|f\|}}{\frac{\|\delta x\|}{\|x\|}} \quad (7.19)$$

$$\begin{aligned} \frac{\|\delta f\|}{\|f\|} &\leq (\kappa(x) + O(1)) \frac{\|\delta x\|}{\|x\|} \\ &= O(\kappa(x)\epsilon_{\text{machine}}) \end{aligned} \quad (7.20)$$

Let P be an exact Householder or Givens, $\tilde{P} = fl(P)$, then

$$fl(\tilde{P}A) = PA + E \quad (7.21)$$

$$\|E\|_2 = O(\epsilon) \|A\|_2 \quad (7.22)$$

Proof: Only for Givens

$$G = \begin{bmatrix} c & s \\ -s & c \end{bmatrix} \quad (7.23)$$

$$a = \begin{bmatrix} x \\ y \end{bmatrix} \quad (7.24)$$

$$\begin{aligned} fl(\tilde{G}a) &= fl\left(\begin{bmatrix} \tilde{c} & \tilde{s} \\ -\tilde{s} & \tilde{c} \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}\right) \\ &= fl\left(\begin{bmatrix} \tilde{c}x + \tilde{s}y \\ -\tilde{s}x + \tilde{c}y \end{bmatrix}\right) \\ &= \begin{pmatrix} cx(1 + \delta_1)(1 + \delta_2)(1 + \delta_3) + sy(1 + \delta_4)(1 + \delta_5)(1 + \delta_6) \\ \dots \end{pmatrix} \end{aligned} \quad (7.25)$$

$$\|fl(\tilde{G}a) - Ga\|_2 = O(\epsilon) \|Ga\|_2 \quad (7.26)$$

$$\begin{aligned} \|E\|_2 &= \|fl(\tilde{G}A) - GA\|_2 \\ &= O(\epsilon) \|GA\|_2 \\ &= O(\epsilon) \|A\|_2 \end{aligned} \quad (7.27)$$

Similarly for Householder.

7.7 Backward Substitution (L-lower Triangular)

If $Lx = b$, and \tilde{x} is the floating point solution, then

$$(L + \delta L)\tilde{x} = b \quad (7.28)$$

$$\|\delta L\| \leq \|L\| O(\epsilon_{\text{machine}}) \quad (7.29)$$

Proof:

$$l_{11}x_1 = b_1 \quad (7.30)$$

$$l_{21}x_1 + l_{22}x_2 = b_2 \quad (7.31)$$

...

$$l_{n1}x_1 + \cdots + l_{nn}x_n = b_n \quad (7.32)$$

then

$$\begin{aligned} \tilde{x}_1 &= \frac{b_1}{l_{11}}(1 + \delta_1) \\ &= \frac{b_1}{\frac{l_{11}}{(1 + \delta_1)}} \\ &= \frac{b_1}{\tilde{l}_{11}} \end{aligned} \quad (7.33)$$

$$\begin{aligned} \tilde{x}_2 &= \frac{(b_2 - l_{21}\tilde{x}_1(1 + \delta_2))(1 + \delta_3)}{l_{22}}(1 + \delta_4) \\ &= \frac{b_2 - l_{21}(1 + \delta_2)\tilde{x}_1}{\frac{l_{22}}{(1 + \delta_3)(1 + \delta_4)}} \\ &= \frac{b_2 - \tilde{l}_{21}\tilde{x}_1}{\tilde{l}_{22}} \end{aligned} \quad (7.34)$$

$$\begin{aligned} \dots \\ \tilde{x}_n &= \frac{((b_n - l_{n1}\tilde{x}_1(1 + \delta_{n1}))(1 + \delta'_{n1}) - l_{n2}\tilde{x}_2(1 + \delta_{n2}))(1 + \delta_{n2})}{l_{nn}} \\ &= \frac{b_n - l_{n1}(1 + \delta_{ij} \text{ products})\tilde{x}_1 - l_{n2}(1 + \delta_{ij} \text{ products})\tilde{x}_2}{\frac{l_{nn}}{1 + \delta_{ij} \text{ products}}} \end{aligned} \quad (7.35)$$