

## Precisión de Algoritmo Backward-Stable

ALN2022  
Clase 16  
26/10/21

De las discusiones anteriores vimos que la precisión de un algoritmo bad-st. depende del condicionamiento del problema. Podemos concluir este resultado con el siguiente teorema.

Teo: Sea  $S: I \rightarrow O$  un problema (input-output map)

y sea  $\tilde{S}$  un algoritmo bad-st. que aproxima  $S$ .

Luego la precisión del output satisface.

$$\|\tilde{S}(a) - S(a)\| \leq O(\mu(a) E_{\text{machine}}) \cdot \|S(a)\|$$

Demo. Como  $\tilde{S}$  es bad-st. se tiene  $\tilde{S}(a) = S(\tilde{a})$  para cierto

$$\tilde{a} / \frac{\|\tilde{a} - a\|}{\|a\|} = O(E_{\text{machine}}). \text{ Luego tenemos}$$

$$\|\tilde{S}(a) - S(a)\| \leq \|S(\tilde{a}) - S(a)\| \leq \|DS(a)\| \cdot \|a - \tilde{a}\| + o(\|a - \tilde{a}\|)$$

$$= \left( \frac{\|DS(a)\|}{\|S(a)\|} \cdot \|a\| \right) \frac{\|a - \tilde{a}\|}{\|a\|} \cdot \|S(a)\| + o(\|a - \tilde{a}\|)$$

$$= \mu(a) \cdot O(E_{\text{machine}}) \|S(a)\| + o(\|a - \tilde{a}\|).$$

Pensando que podemos asumir que  $E_{\text{machine}} \rightarrow 0$  tenemos que el término  $o(\|a - \tilde{a}\|) = o(\|a\| \cdot O(E_{\text{machine}})) \rightarrow 0$ .  $\square$

## Estabilidad de Triangularización Householder

Comencemos realizando experimentos en MATLAB para ver la precisión de este algoritmo. (ver notebook Estabilidad-Householder)

Lo que está de fondo es el siguiente resultado (que enunciaremos en pruebas)

Teo: Consideremos una máquina con la Aritmética de P.F.

(satisfaciendo axiomas  $\forall x \in \mathbb{R}, \exists \epsilon > 0, |\epsilon| \leq \epsilon_{med} \text{ t.q. } f(x) = x(1 \pm \epsilon)$ )  
 $\forall x, y \in \mathbb{R}, \exists \epsilon > 0, |\epsilon| \leq \epsilon_{med} \text{ t.q. } x \odot y = x \cdot y(1 \pm \epsilon)$

El algoritmo Householder para computar la descomp.

QR de una matriz  $A \in \mathbb{C}^{m \times n}$ , es back-st., i.e. el

output  $\hat{Q}, \hat{R}$  del algoritmo satisficen

$$\hat{Q} \cdot \hat{R} = A + \delta A \quad \text{con} \quad \frac{\|\delta A\|}{\|A\|} = O(\epsilon_{med})$$

para alguna matriz  $\delta A \in \mathbb{C}^{m \times n}$ .

Comentarios: En otras palabras, la Q computada ( $\hat{Q}$ ) y

la R computada ( $\hat{R}$ ) son ~~soluciones~~ la factorización QR de un problema (lineal), perturbado.

La matriz  $\hat{R}$  es exactamente la matriz triangular superior que resulta de la aplicación de las reflexiones Householder en punto flotante.

Sin embargo la  $\tilde{Q}$  es una verdadera matriz ortogonal que resulta del tomar el producto  $\overbrace{\text{Id} - 2\frac{\tilde{v}_k \tilde{v}_k^T}{\tilde{v}_k^T \tilde{v}_k}}^{\text{reflexión } H}$  (las reflexiones  $H$ ) pero donde  $\tilde{v}_k$  son los vectores en punto flotante ( $\tilde{Q} = \tilde{Q}_1 \dots \tilde{Q}_m$ ,  $\tilde{Q}_k = \tilde{Q}_k - \tilde{Q}_k$ )

Esta restricción en el resultado en realidad no es tal dada que si recordamos el algoritmo Householder la matriz  $Q$  no la escribimos explícita (solo actualizamos la matriz hasta llegar a una triangular superior).

- Como vimos, individualmente  $\tilde{Q}$  y  $\tilde{R}$  no son precisos pero en realidad lo que ocurre es que se supone que el problema de encontrar  $Q$  (y  $R$ ) está mal condicionado. (Sería interesante saber hasta qué punto este resultado está probado teóricamente.)

- En general, la descomposición  $QR$  no es en fin en sí mismo pero en medios para estudiar o resolver otro problema. Recordar que por ejemplo nos da una forma sencilla de resolver un sistema lineal  $Ax = b$ .

La pregunta es si los errores individuales en  $Q$  y  $R$  pueden afectar la precisión de resolver el sistema. La respuesta es que felizmente no, al menos con la precisión conjunta de  $Q, R$ .

## Analicemos el algoritmo de resolver $AX=b$ : Alg-SL

Recorden que el algoritmo consiste en

- i)  $QR=A$  # realizar factorización via Reflexions H.
- ii)  $y=Q^*b$  # se construye  $Q^*b$
- iii)  $x=R^{-1}y$  # resolvamos el sistema triangular  $Rx=y$  por sustitución "hacia atrás".

Se puede probar que i), ii) y iii), son back-st. (i) ya lo mencionamos) y esto tiene como resultado que el Alg-SL también lo es.

El paso ii), realiza el cálculo  $Q^*b$  mediante el alg. ya visto en clase (cálculo implícito). Que sea back-st. implica que el output  $\tilde{y}$  satisface  $(Q + \delta Q)^* b = \tilde{y}$

$$(\tilde{Q} + \delta Q) \tilde{y} = b \quad \text{con} \quad \|\delta Q\| = O(\epsilon_{mach}) \quad (I)$$

Análogamente el paso iii) da un output  $\tilde{x}$  que satisface

$$(\tilde{R} + \delta R) \tilde{x} = \tilde{y} \quad \text{con} \quad \frac{\|\delta R\|}{\|\tilde{R}\|} = O(\epsilon_{mach}) \quad (II)$$

Juego de (I) y (II) tenemos que

$$b = (\tilde{Q} + \delta Q) (\tilde{R} + \delta R) \tilde{x} = (\tilde{Q}\tilde{R} + \delta Q \cdot \tilde{R} + \tilde{Q} \delta R + \delta Q \cdot \delta R) \tilde{x}$$

Utilizando que el Alg. Householder es b-s. tenemos  $(A + \delta A = \tilde{Q} \tilde{R})$

$$b = \underbrace{(A + \delta A + \delta Q \cdot \tilde{R} + \tilde{Q} \cdot \delta R + \delta Q \cdot \delta R)}_{\Delta A} \tilde{x}$$

i.e.  $b = (A + \Delta A) \tilde{x}$ , por lo que resta probar que  $\|\Delta A\|$  es "relativamente chico", i.e.

$$\frac{\|\Delta A\|}{\|A\|} \leq O(\epsilon_{mach}). \quad \text{Para eso observamos que}$$

$$\tilde{R} = \tilde{Q}^*(A + \delta A) \Rightarrow \frac{\|\tilde{R}\|}{\|A\|} = \frac{\|A + \delta A\|}{\|A\|} = O(1)$$

$$\frac{\|\delta Q \tilde{R}\|}{\|A\|} \leq \|\delta Q\| \cdot \frac{\|\tilde{R}\|}{\|A\|} \leq O(\epsilon_{mach}) \cdot \frac{\|\tilde{R}\|}{\|A\|} = O(\epsilon_{mach})$$

Análogo para  $\frac{\|\tilde{Q} \delta R\|}{\|A\|} = O(\epsilon_{mach})$

Y además  $\frac{\|\delta Q \delta R\|}{\|A\|} \leq O(\epsilon_{mach}^2)$ , por lo tanto,

$$\frac{\|\Delta A\|}{\|A\|} \leq \frac{\|\delta A\|}{\|A\|} + \frac{\|\delta Q \tilde{R}\|}{\|A\|} + \frac{\|\tilde{Q} \delta R\|}{\|A\|} + \frac{\|\delta Q \delta R\|}{\|A\|}$$

$$\leq O(\epsilon_{mach}) \quad \square$$