

Estabilidad

ALN2022
Clase 14
20/10

En un mundo ideal, los algoritmos en análisis numérico podrían dar soluciones exactas a nuestros problemas, pero este es el mundo de la realidad. Los problemas que estamos mirando pertenecen al mundo continuo pero nuestras máquinas digitales viven en un mundo discreto. En este sentido la noción de estabilidad nos da una medida de la precisión de los outputs de nuestros algoritmos.

Recordando el mapa solución $S: I \rightarrow O$ del espacio de inputs en el espacio de outputs (en general definido localmente como hemos visto) podemos pensar en algoritmo como otro mapa $\tilde{S}: I \rightarrow O$ que a un input $a \in I$ devuelve una solución aproximada.

Más formalmente, tenemos $S: I \rightarrow O$, una computadora con la aritmética dada anteriormente (de punto flotante), un algoritmo ~~para~~ y su implementación. Luego dado un input $a \in I$, lo ~~se~~ introducimos en la computadora ~~para~~ resultando en $fl(a)$ (el mapa redondeo) para luego implementar el algoritmo y dar como resultado en n^o en punto flotante $\tilde{S}(a)$:

Resumiendo: $a \in I \xrightarrow{\text{computadora}} fl(a) \xrightarrow{\text{algoritmo}} \tilde{S}(a) \in \mathbb{F}$ " " " " " "

Lo anterior de hecho está idealizado, de hecho la misma secuencia de pasos nos podría dar una solución distinta por lo que quizá $\tilde{S}(a)$ puede no estar bien definido.

Pero a pesar de las idealizaciones y complicaciones podemos hacer un análisis riguroso de los algoritmos.

Precisión

Una forma natural de decir que nuestro algoritmo es preciso es que el output dado por el algoritmo esté muy cerca del output real.

En términos ya vistos, decimos que el algoritmo es preciso si:

$$\forall a \in I, \quad \frac{\|\tilde{S}(a) - S(a)\|}{\|S(a)\|} = O(\epsilon_{\text{machine}}).$$

Esta expresión parece natural pero esconde en detalle importante.
¿Qué entendemos por $O(\epsilon_{\text{machine}})$? ¿Es uniforme en todo los inputs?

Sabemos de la teoría del condicionamiento que el n.º de cond. es inevitable. Una vez que nuestro input se transpara de a a $f(a)$ tenemos una pequeña perturbación que puede ser muy significativa si el input está mal condicionado.

Mencionado esto podemos decir que el algoritmo es preciso si:

$$\| \tilde{S}(a) - S(a) \| \leq C_{\mu}(a) \cdot \epsilon_{\text{machine}} \cdot \| S(a) \| \quad \text{para alguna cte. } C \text{ que sólo puede depender de los dimensiones.}$$

Estabilidad

La estabilidad de un algoritmo \tilde{S} asociado al mapa S se describe de la siguiente manera. Decimos que \tilde{S} es estable si $\forall a \in I$

$$\frac{\|\tilde{S}(a) - S(\tilde{a})\|}{\|S(\tilde{a})\|} = O(\epsilon_{\text{machine}}) \quad \text{para } \tilde{a} \in I \text{ t.º } \frac{\|a - \tilde{a}\|}{\|a\|} = O(\epsilon_{\text{machine}})$$

En palabras de Trefler:

"Un algoritmo estable da casi la respuesta exacta a casi la pregunta correcta".

Backward Stability (Estabilidad Inversa)

Una condición más fuerte que la estabilidad ~~de~~ es la est. inversa:

\tilde{S} es backward-stable si $\forall a \in I \quad \tilde{S}(a) = S(\tilde{a})$

para $\tilde{a} \in I$ tq $\frac{\| \tilde{a} - a \|}{\| a \|} = O(\epsilon_{\text{mach}})$

i.e. "Un algoritmo es b.stable si da la respuesta exacta a casi la pregunta correcta".

Observar que si \tilde{S} es b.s. entonces

$$\frac{\| \tilde{S}(a) - S(a) \|}{\| S(a) \|} = \frac{\| \tilde{S}(a) - S(\tilde{a}) \| + \| S(\tilde{a}) - S(a) \|}{\| S(a) \|} \\ = \frac{\| S(\tilde{a}) - S(a) \|}{\| S(a) \|} \leq \frac{\mu(a) \cdot \| \tilde{a} - a \|}{\| a \|} = O(\epsilon_{\text{mach}})$$

Estabilidad a la Smale

Supongamos que estamos en las condiciones ~~ver~~ anteriormente de que tenemos que la variedad solución $V = \{ (a, x) \in I \times \mathcal{O} : F(a, x) = 0 \}$ con F regular (al menos $F(a, \cdot) \in C^2$), y las dimensiones adecuadas.

Luego x es output de a si $F_a(x) = 0$ con $F_a(\cdot) = F(a, \cdot)$.

Se puede construir un mapa de Newton asociado a la función

$F_a : \mathcal{O} \rightarrow \mathbb{R}^k$, y lo definimos por N_{F_a}

Smale define el output de ~~un~~ algoritmo a un output $\tilde{x} = \tilde{S}(a)$ tal que $\tilde{x} \in N_{F_a}(\tilde{x})$ está en la cuenca de atracción inmediata ~~del~~ ~~del~~ de convergencia cuadrática del ~~del~~ N_{F_a} al cero x .

i.e. $d(N_{F_a}^k(\tilde{x}), x) \leq \left(\frac{1}{2}\right)^{2^k-1} d(\tilde{x}, x)$.

Observa que la definición de Smale

Smale le da el nombre de "cero aproximado" a los ~~primeros~~ puntos que están en la cuenca de atracción anterior.

Por lo tanto si nuestro algoritmo nos da un cero aproximado, es muy fácil ~~para~~ diseñar un algoritmo estable ya que a veces ~~output~~ le podemos aplicar el método de Newton para obtener la precisión que deseamos.

Algunos Ejemplos de Estabilidad

Estabilidad de Aritmética de Punto Flotante

Veamos que las operaciones básicas (sus correspondientes algoritmos) son back-estables.

Resta: Sea $(a_1, a_2) \in \mathbb{I}^2$ y nuestra función $S(a_1, a_2) = a_1 - a_2$

Luego nuestra computadora (con el axioma fundamental) define el algoritmo

$$\tilde{S}(a_1, a_2) = fl(a_1) \ominus fl(a_2)$$

i.e. redondea $a_i \mapsto fl(a_i)$, para luego proceder a realizar un cálculo en \mathbb{F} (que desconocemos) para dar un elemento

en \mathbb{F} . Recordando el mapa $fl: \mathbb{R} \rightarrow \mathbb{F}$ y el axioma se tiene

que:

$$fl(a_1) = a_1(1 + \varepsilon_1) \quad fl(a_2) = a_2(1 + \varepsilon_2) \quad \text{con } |\varepsilon_1|, |\varepsilon_2| \leq \varepsilon_{\text{mach}}$$

$$\text{Y luego. } fl(a_1) \ominus fl(a_2) = (fl(a_1) - fl(a_2))(1 + \varepsilon_3) \\ |\varepsilon_3| \leq \varepsilon_{\text{mach}}$$

Por lo tanto

$$\begin{aligned} \tilde{S}(a_1, a_2) &= [a_1(1 + \varepsilon_1) - a_2(1 + \varepsilon_2)] \cdot (1 + \varepsilon_3) \\ &= a_1(1 + \varepsilon_1)(1 + \varepsilon_3) - a_2(1 + \varepsilon_2)(1 + \varepsilon_3) \\ &= a_1(1 + \varepsilon_4) - a_2(1 + \varepsilon_5) \end{aligned}$$

$$\left(\begin{array}{l} \text{con } |\varepsilon_4| \leq 2\varepsilon_{\text{mach}} + \mathcal{O}(\varepsilon_{\text{mach}}^2) \\ \text{y } |\varepsilon_5| \end{array} \right)$$

Hemos encontrado que $\tilde{S}(a_1, a_2) = S(\tilde{a}_1, \tilde{a}_2)$

con $\tilde{a}_1 := a_1(1 + \varepsilon_4)$, $\tilde{a}_2 := a_2(1 + \varepsilon_5)$ satisfaciendo

$$\frac{\|\tilde{a}_1 - a_1\|}{\|a_1\|} = \mathcal{O}(\varepsilon_{\text{mach}}) \quad , \quad \frac{\|\tilde{a}_2 - a_2\|}{\|a_2\|} = \mathcal{O}(\varepsilon_{\text{mach}})$$

$$\text{y por lo tanto } \frac{\|(\tilde{a}_1, \tilde{a}_2) - (a_1, a_2)\|}{\|(a_1, a_2)\|} = \mathcal{O}(\varepsilon_{\text{mach}})$$

concluyendo que es backward-estable.

Punto Flotante (Revisado)

Fijamos la base $\beta=2$, y nos permitiremos trabajar con t y s distintos siendo t la cantidad de bits para la mantisa y s cantidad de bits para el exponente.

Recordar que en \mathbb{R} se puede escribir en expansión binaria como

$$x = \pm (d_n 2^n + d_{n-1} 2^{n-1} + \dots + d_0 + d_{-1} 2^{-1} + \dots)$$

donde los d_k son dígitos binarios, i.e. $d_k \in \{0, 1\}$.

Una forma abreviada de lo anterior es escribir

$$x = \pm (d_n d_{n-1} \dots d_0 d_{-1} d_{-2} \dots)_2$$

Un n.º en punto flotante lo escribimos como $x = m \times 2^e$

$$x = \pm (d_{-1} d_{-2} \dots d_{-t})_2$$

mantisa

exponent

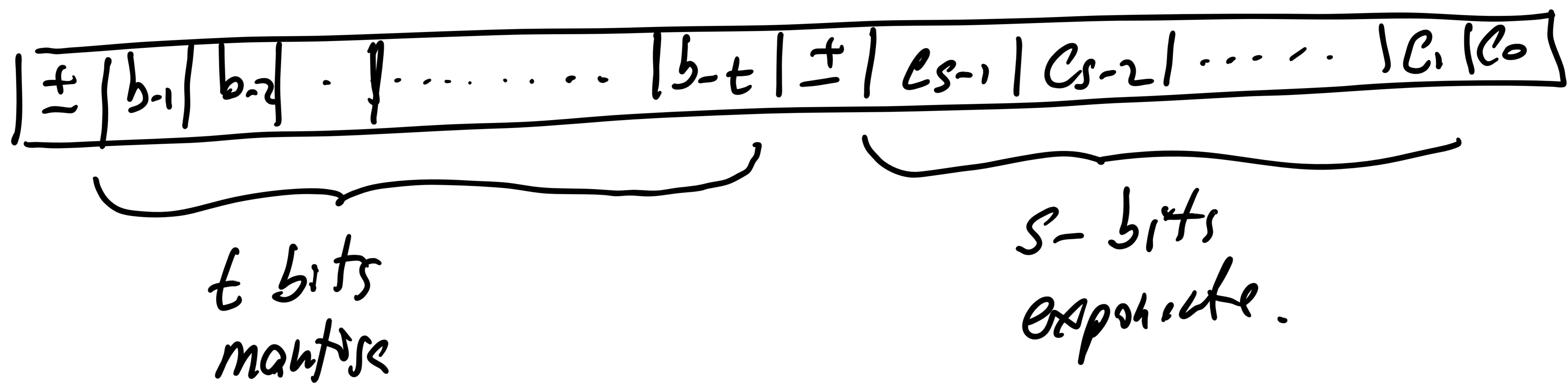
donde $m = \pm (d_{-1} d_{-2} \dots d_{-t})_2$, $e = \pm (c_{s-1} \dots c_1 c_0)_2$

A diferencia de la clase pasada al exponente lo escribimos explícitamente el exponente. Recordar que \mathbb{F} son los normalizados, es decir $b_{-1} = 1$.

~~Tomando~~ Observar que ahora \mathbb{F} queda bien definida y sólo depende de los parámetros (t, s) y podemos

notar el $\mathbb{F} = \mathbb{R}(t, s)$. (Si x no está normalizado casi siempre poder llevarlo a \mathbb{F})

Luego la representación física del punto flotante se modela con los siguientes esquemas



El conjunto $\mathbb{R}(t,s)$ es finito. Calculemos explícitamente, en función de los parámetros t, s cuál es el n.º mayor, y cuál es el más cercano al 0.

$$\max_{x \in \mathbb{R}(t,s)} |x| = \left(\frac{1}{2} + \frac{1}{2^2} + \dots + \frac{1}{2^t} \right) \times 2^{1 + \frac{1}{2} + 2 + \dots + 2^{s-1}}$$

$$= \left(1 - \frac{1}{2^t} \right) 2^{\frac{2^s - 1}{2}}$$

$$\min_{x \in \mathbb{R}(t,s)} |x| = \frac{1}{2} \times 2^{-2^{s-1}}$$

En 32-bits $\rightarrow t=23, s=7$ se tiene $\max \approx 2^7 = 2^{128} \approx 10^{38}$
 y mínimo $2^{-128} \approx 10^{-38}$

En 64-bits ($t=53, s=10$) el resultado es 10^{308} y 10^{-308} respectivamente.

