

Catastrophic loss of precision

Never compute a small number as
the difference between two large numbers

x: US population (thousands)

Feb 311, 189

Mar 311, 356

→ difference $311,356 - 311,189 = 167$ (reasonable)

in our floating point scheme:

$$fl(311,189) = 3.11189 \cdot 10^5$$

$$fl(311,356) = 3.11356 \cdot 10^5$$

$$\sim \begin{array}{r} 3.11356 \cdot 10^5 \\ - 3.11189 \cdot 10^5 \\ \hline \end{array}$$

$$0.00167 \cdot 10^5 = 1.67000 \cdot 10^2$$

→ relative error of $\approx 2 \cdot 10^{-2}$ but $EPS_{REL} \approx 10^{-5}$???

→ almost all digits used up to store the
"large part" of the numbers, which cancel out.

generally. If for N digits of precision,
when the first M digits of two numbers agree,
expect only $N - M$ digits of precision in
the difference.

$M \approx N \rightarrow$ catastrophic precision loss

Ex: Finite differences

$$f'_{FD}(x) = \frac{f(x+h) - f(x)}{h}$$

$$f(x) = x, \text{ want } f'(1) = 1$$

$$a) \quad h = \frac{2}{3}, \quad fl(h) = 0.66667$$

$$f(x+h) = 1.6667$$

$$f(x) = 1.0000$$

$$\rightarrow \frac{f(x+h) - f(x)}{h} = \frac{0.66670}{0.66667}$$

differ in 4th digit.

$$\rightarrow f'_{FD}(1) \approx 1 + O(10^{-4})$$

$$b) \quad h = \frac{2}{30}, \quad fl(h) = 0.066667$$

still 5 sign. digits!

but $f(x+h) \approx 1.0667$

$$f(x) \approx 1.0000$$

$$\frac{f(x+h) - f(x)}{h} = \frac{0.066700}{0.066667}$$

differs in 3rd digit!

$$\rightarrow f'_{\text{RD}}(1) \approx 1 + \mathcal{O}(10^{-3})$$

→ decreasing h by factor 10 made error worse by factor 10!

→ lost precision because much of mantissa was used to store 1.0000 (the large number)

how to avoid floating point precision loss?

→ Rearrange things creatively to avoid the nasty differences!

Ex: $f(x, \Delta) = \sqrt{x+\Delta} - \sqrt{x}$, $\Delta \ll x$.

for instance: $x = 900$, $\Delta = 4 \cdot 10^{-3}$

→ $\underbrace{30.0000}_{6 \text{ wasted precision digits}}6667 - 30.00000000$

in our scheme: $fl(f(900, 4 \cdot 10^{-3})) = 0$

here we can use:

$$(\sqrt{x+\Delta} - \sqrt{x})(\sqrt{x+\Delta} + \sqrt{x}) = x + \Delta - x = \Delta$$

$$\Rightarrow \sqrt{x+\Delta} - \sqrt{x} = \frac{\Delta}{\sqrt{x+\Delta} + \sqrt{x}}$$

$$\leadsto \frac{4 \cdot 10^{-3}}{30.00006667 - 30.00000000} \approx \frac{4 \cdot 10^{-3}}{60000 \cdot 10^2} \approx 6.6667 \cdot 10^{-5}$$

here we got all the precision we want!

$$\Rightarrow \boxed{\vec{f}'' = \frac{1}{h^2} A \vec{f}} \quad (\text{matrix multiplication})$$

can invert : $\vec{f} = h^2 A^{-1} \vec{f}''$ (finite difference method for PDEs)

nontrivial b.c.s

$$f(x_0) = f_0, \quad f(x_{N+1}) = f_{N+1}$$

$$\rightarrow f_1'' \approx \frac{1}{h^2} (f_0 - 2f_1 + f_2)$$

$$f_N'' \approx \frac{1}{h^2} (f_{N-1} - 2f_N + f_{N+1})$$

$$\begin{pmatrix} f_1'' - \frac{1}{h^2} f_0 \\ f_2'' \\ \vdots \\ f_{N-1}'' \\ f_N'' - \frac{1}{h^2} f_{N+1} \end{pmatrix} = h^2 A \vec{f}$$

$$= \vec{f}'' - \frac{1}{h^2} \vec{\Delta} \quad \text{— sparse vector of b.c.s}$$

$$\Rightarrow \vec{f} = h^2 A^{-1} \left(\vec{f}'' - \frac{1}{h^2} \vec{\Delta} \right)$$