

Data Error and Computational Error.

3/

compute $f(x)$ for a given problem.

- x true value of input
- $f(x)$ true value of output.
- \hat{x} approximate input.
- \hat{f} computed value of $f(x)$ (approximate of the function).

$$\begin{aligned}\text{Total error : } & \hat{f}(\hat{x}) - f(x) \\ &= \underbrace{\hat{f}(\hat{x}) - f(\hat{x})}_{\text{computational error}} + \underbrace{f(\hat{x}) - f(x)}_{\text{data error}}\end{aligned}$$

- The algorithm that is chosen to compute $f(x)$ has no effect on the data error.
- We can only know the computational error and never the total error unless $f(x)$ is known.

Truncation and Rounding Errors

- truncation error : difference between true result and result produced by given algorithm using exact arithmetic

$$\begin{array}{ccc}\frac{1}{1-h} = 1 + h + h^2 + h^3 + \dots & \approx & \underbrace{1 + h + h^2}_{\text{algorithm}} + \underbrace{\mathcal{O}(h^3)}_{\text{truncation error}} \\ \uparrow & & \\ \text{true result} & & \text{big-O of } h^3\end{array}$$

4/

- Rounding error: difference between result produced by algorithm using exact arithmetic and result produced by the same algorithm using limited precision arithmetic

example: $|1.085679 - 1.0857| = 0.000029$

- computational error = truncation error + rounding error.

$$= \underbrace{\hat{f}(\hat{x}) - f(\hat{x})}$$

truncation error
will come from the order of
accuracy that which was used
to compute $f(\hat{x})$ and $\hat{f}(\hat{x})$

rounding error comes
from the precision
used to represent
 $\hat{f}(\hat{x})$ or \hat{x} .

Absolute and Relative errors

- absolute error: approximate value - true value

- relative error: $\frac{\text{absolute value}}{\text{true value}}$

Example: Finite Difference Approximation

$$f'(x) \approx \frac{f(x + \epsilon e) - f(x)}{\epsilon} \text{ is a first order}$$

approximation. , where ϵ is very small, $\epsilon \ll 1$ and
 e is a vector, $e = [1, \dots, 1]^T$ $\|e\| = 1$

How do we derive this formula?

Finite-Difference Approach

5/

From Taylor's Theorem, if f is twice differentiable, then f can be expanded at point x , such as

$$f(x+p) = f(x) + \nabla f(x)^T p + \frac{1}{2} p^T \nabla^2 f(x+tp) p, \quad t \in (0,1)$$

If L is a bound on the magnitude of the second derivative, $\nabla^2 f(x+tp)$, then

$$\|\nabla^2 f(x+tp)\| \leq L.$$

Thus, $\|f(x+p) - f(x) - \nabla f(x)^T p\| \leq \left(\frac{L}{2}\right) \|p\|^2 \leq \frac{1}{2} \|\nabla^2 f\| \cdot \|p\|^2$

choose $p = \epsilon e_i$, then

$$\|f(x + \epsilon e_i) - f(x) - \nabla f(x)^T \epsilon e_i\| \leq \left(\frac{L}{2}\right) \|\epsilon e_i\|^2$$

Here ϵ is a very small scalar value, $\epsilon \ll 1$ and e is a vector, $e = [1, \dots, 1]^T$.

Then

$$\|f(x + \epsilon e_i) - f(x) - \nabla f(x)^T \epsilon e_i\| \leq \left(\frac{L}{2}\right) \epsilon^2.$$

since $\|e\|^2 = 1$.

Next divide by ϵ ,

$$\left\| \frac{f(x + \epsilon e_i) - f(x)}{\epsilon} - \nabla f(x)^T e_i \right\| \leq \left(\frac{L}{2}\right) \epsilon$$

$\nabla f(x)^T e_i$ is then the first derivative of $f(x)$, or $\frac{\partial f}{\partial x_i}$.

So,

$$\frac{\partial f}{\partial x_i} = \frac{f(x + \epsilon e_i) - f(x)}{\epsilon} + \delta_e, \quad |\delta_e| \leq \left(\frac{L}{2}\right) \epsilon.$$

This is simply called a forward difference formula. As $\epsilon \rightarrow 0$, then

b/

$$\frac{\partial f}{\partial x_i} = \frac{f(x + \epsilon e_i) - f(x)}{\epsilon} + O(\epsilon), \text{ and}$$

$$O(\epsilon) = |\delta_e| \leq \left(\frac{L}{2}\right) \epsilon, \text{ so } O(\epsilon) = \|\delta_e\| \rightarrow 0 \text{ as } \epsilon \rightarrow 0.$$

And $\frac{f(x + \epsilon e_i) - f(x)}{\epsilon}$ approximates the first derivative.

To improve the accuracy, a second order approximation can be derived, from the Taylor series expansion, if p is set to $+ \epsilon e_i$ and $- \epsilon e_i$, then

$$f(x + \epsilon e_i) = f(x) + \epsilon \nabla f(x)^T + \frac{1}{2} \epsilon^2 \nabla^2 f(x + t \epsilon e_i) + O(\epsilon^3)$$

$$f(x - \epsilon e_i) = f(x) - \epsilon \nabla f(x)^T + \frac{1}{2} \epsilon^2 \nabla^2 f(x - t \epsilon e_i) + O(\epsilon^3)$$

Subtract the second from the first,

$$f(x + \epsilon e_i) - f(x - \epsilon e_i) = 2\epsilon \nabla f(x)^T + \frac{1}{3} \epsilon^2 (\nabla^2 f(x + t \epsilon e_i) - \nabla^2 f(x - t \epsilon e_i)) + O(\epsilon^3)$$

$$\frac{f(x + \epsilon e_i) - f(x - \epsilon e_i)}{2\epsilon} = \nabla f(x) + \frac{1}{4} \epsilon^2 (\nabla^2 f(x + t \epsilon e_i) - \nabla^2 f(x - t \epsilon e_i)) + O(\epsilon^2)$$

as $\epsilon \rightarrow 0$,

$$\frac{f(x + \epsilon e_i) - f(x - \epsilon e_i)}{2\epsilon} = \frac{\partial f}{\partial x_i} + \frac{1}{4} \epsilon \left(\frac{\partial^2 f}{\partial x_i^2} - \frac{\partial^2 f}{\partial x_i^2} \right) + O(\epsilon)$$

$$\frac{\partial f}{\partial x_i} = \frac{f(x + \epsilon e_i) - f(x - \epsilon e_i)}{2\epsilon} + O(\epsilon^2)$$

when using a computer whose processor is based on IEEE 754 floating-point arithmetic standard, the numerical value of $f(x)$ is bounded by u , where $u = 1.1 \times 10^{-16}$.

Thus

$$|f(x)_{\text{numerical}} - f(x)| \leq u L_f$$

$$|f(x + \epsilon e_i)_{\text{numerical}} - f(x + \epsilon e_i)| \leq u L_f$$

where L_f is a bound on $|f'(x)|$.

Therefore $\frac{\partial f}{\partial x_i}(x) = \frac{f(x + \epsilon e_i) - f(x)}{\epsilon} + (1/2)\epsilon + 2u L_f / \epsilon.$

Error is bounded by $(\frac{L}{2})\epsilon + \frac{2u L_f}{\epsilon} = \delta_e$

$$\text{so } \delta_e \epsilon^2 = \frac{4 L_f u}{L}$$

$$\delta_e \epsilon = \sqrt{u}$$

assuming $\frac{L_f}{L}$ is scaled constantly.

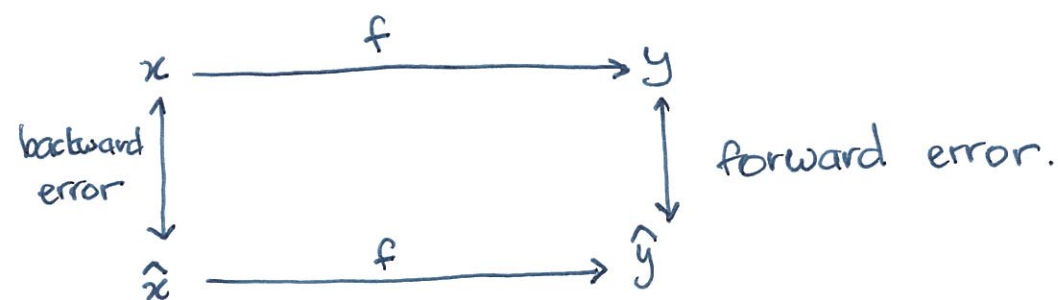
This gives us an estimate for ϵ .

Forward and Backward Error

8

If we want to compute $y = f(x)$, but then obtain an approximate value \hat{y} .

- The forward error : $\Delta y = \hat{y} - y$, where $\hat{y} = f(\hat{x})$.
- The backward error : $\Delta x = \hat{x} - x$



Example 1

If $y = \sqrt{2}$, then $\hat{y} = 1.4$ is an approximation

$$\begin{aligned}\text{Its absolute forward error, } |\Delta y| &= |\hat{y} - y| \\ &= |1.4 - 1.4142\dots| \\ &\approx 0.0142\end{aligned}$$

$$\text{relative forward error, } \frac{|\Delta y|}{y} \approx 1\%$$

Since $\sqrt{1.96} = 1.4$, then

$$\begin{aligned}\text{the absolute backward error} &= |\Delta x| = |\hat{x} - x| \\ &= |1.96 - 2|\end{aligned}$$

$$\begin{aligned}\text{relative backward error, } \frac{|\Delta x|}{x} &\approx 2\% \\ &= 0.04\end{aligned}$$

Example 2

Approximate $f(x) = e^x$ with a simpler function, and evaluate its accuracy at $x=1$.

$$\text{We know that } f(x) = e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots$$

If we truncate the series as,

$$\hat{f}(x) = 1 + x + \frac{x^2}{2} + \frac{x^3}{6} + O(x^4)$$

The forward error is $|\hat{f}(x) - f(x)|$

The backward error is $|\hat{x} = \log(\hat{f}(x)) - x|$

$$|\hat{f}(x) - f(x)| = |2.666667 - 2.718282| \approx 0.051615$$

$$|\hat{x} - x| = |0.980829 - 1| = 0.019171$$

Sensitivity and Conditioning

- Problem is sensitive or ill-conditioned, if relative change in solution can be much larger than that in input data or condition number $\gg 1$.
- Condition number :

$$\begin{aligned} \text{cond} &= \frac{|\text{relative change in solution}|}{|\text{relative change in input data}|} \\ &= \left| \frac{f(x) - f(\hat{x})}{f(x)} \right| \bigg/ \left| \frac{\hat{x} - x}{x} \right| \end{aligned}$$

$$\text{cond} := \frac{\text{relative forward error}}{\text{relative backward error}}$$

Example 1: Given a function $f(x)$, evaluate the condition number if $\hat{x} = x + \Delta x$.

$$f(\hat{x}) = f(x + \Delta x)$$

$$\frac{f(\hat{x}) - f(x)}{f(x)} = \frac{f(x + \Delta x) - f(x)}{f(x)} \approx \frac{f'(x) \Delta x}{f(x)}$$

$$\text{condition number} \approx \frac{\left| \frac{f'(x) \Delta x}{f(x)} \right|}{\left| \frac{\Delta x}{x} \right|} = \left| \frac{x f'(x)}{f(x)} \right|$$

Example 2 : Evaluate the sensitivity of the tangent function, near $\pi/2$

$$\frac{\pi}{2} = 1.570796 \dots$$

$$f(x) = \tan(1.57079) \approx 1.58 \times 10^5$$

$$f(\hat{x}) = \tan(1.57078) \approx 6.124 \times 10^4$$

$$\text{condition number} = \frac{\left| \frac{6.124 \times 10^4 - 1.58 \times 10^5}{1.58 \times 10^5} \right|}{\left| \frac{1.57078 - 1.57079}{1.57079} \right|} = 9.6 \times 10^4$$

Stability

- Algorithm is stable if result produced is relatively insensitive to perturbations to the input, which also means that the condition number is very small.

Accuracy

- Difference between computed and true solution.
- Stability alone does not guarantee accurate results.

Floating Point Numbers

$$x = \pm .b_1 b_2 \dots b_t \times \beta^e$$

- $b_1 b_2 \dots b_t$ is the mantissa or $t = \text{precision}$
- e is the exponent, $L \leq e \leq U$
- b_i are the base digits, $0 \leq b_i \leq \beta - 1$.
- β is the base.

example : if $x = 14.7$ $(\beta, t, L, U) = (10, 3, -9, 9)$

$$x = .147 \times 10^2$$

In computers we use IEEE Standards,

	β	t	L	U
IEEE (SP - single precision)	2	24	-126	127
IEEE (DP - double precision)	2	53	-1022	1023

IEEE SP translates to 8 significant digits in the decimal system and 16 in DP. 12/

Theorem : Suppose we are given the set of floating point numbers $F(\beta, t, L, U)$ and that $x \in \mathbb{R}$ satisfies $\beta^L < |x| < \beta^U$. It follows that

$$\frac{|f(x) - x|}{|x|} \leq \frac{1}{2} \beta^{1-t}$$

→ $f(x)$ is the floating-point approximation of a given real number x .

so $\frac{|f(x) - x|}{|x|}$ is the relative error in representing real number x within range of floating-point system.

→ For IEEE SP $\frac{1}{2} \beta^{1-t} = \frac{1}{2} 2^{1-24} = 5.96 \text{e-}8$

IEEE DP $\frac{1}{2} \beta^{1-t} = \frac{1}{2} 2^{1-53} = 1.11 \text{e-}16$

→ Therefore the accuracy of the floating point system is given by $\frac{1}{2} \beta^{1-t}$ or also called machine precision.

Two Issues with Finite - precision Arithmetic

13

1) Rounding

- chop : round toward zero.

- round to nearest :

	2.442	2.4	2.4
example :	2.481	2.4	2.5
	2.591	2.5	2.6
		chop	round to nearest.

2) Cancellation

When two numbers that are almost identical except for the last few digits are used in subtraction, then the resulting number will have only a few digits of accuracy. If we use this new number in a calculation, then the results will only have a few digits of accuracy.

example : $a = 654738.2919$

$b = 654737.7921$

If we represent a as 0.6547383×10^6 and b as 0.6547378×10^6 , then $a - b = 0.5$. However the true number is 0.4998 .