

Homework Problem Sheet 1

Introduction. Round-off error analysis and cancellation.

Problem 1.1 Avoiding Cancellation

In [NMI, Sect. 1.5] we saw that the subtraction of numbers of about the same size can lead to massive amplification of relative errors in these numbers, a scourge called *cancellation* (Auslöschung). Often cancellation can be avoided by judiciously rewriting an expression, replacing the dangerous subtraction by another operation. This problem will demonstrate this approach in the case of a few examples.

Rewrite the following functions $f(x)$ so that cancellation is avoided. Implement the original function and its more stable reformulation in MATLAB and compute the relative error of the unstable implementation using the stable version as a substitute for the exact value. Tabulate the errors obtained.

(1.1a) $f(x) := (1 - x)/(1 + x) - 1/(3x + 1)$ for $x \approx 0$.

Compute relative errors for $x = 10^{-17}, 10^{-16}, \dots, 10^{-1}, 1$.

(1.1b) $f(x) := \sin(x) - \sin(y)$ for $x \approx y$.

Compute relative errors for $x = 1.5, y - x = 10^{-17}, 10^{-16}, \dots, 10^{-1}, 1$.

HINT: Use trigonometric identities.

(1.1c) $f(x) := \sqrt{x - (1/x)}$ for $x \approx 1$.

Compute relative errors for $x = 1 \pm 10^{-17}, 10^{-16}, \dots, 10^{-10}$.

(1.1d) $f(x) := \sqrt{x + (1/x)} - \sqrt{x - (1/x)}$ for $x \gg 1$.

Compute relative errors for $x = 10, 100, 1000, \dots, 10^{16}, 10^{17}$.

HINT: Use $a - b = \frac{a^2 - b^2}{a + b}$.

Problem 1.2 Summing the Harmonic Series

In analysis you have seen that the harmonic series diverges. On a computer this will not happen, of course!

The series $\sum_{k=1}^{+\infty} k^{-1}$ is called the harmonic series. The partial sums, $S_n = \sum_{k=1}^n k^{-1}$, can be computed recursively by setting $S_1 = 1$ and using $S_n = S_{n-1} + n^{-1}$. If this computation were carried out on your computer, what is the largest S_n that would be obtained? (Do not do this experimentally on the computer; it is too expensive.)

HINT: Find n such that $|\frac{S_n - S_{n-1}}{S_n}| < u(\mathbb{F})$, where $u(\mathbb{F})$ is the unit round-off of the floating-point number system \mathbb{F} . To this end, first prove that $\sum_{k=1}^n \frac{1}{k} > \ln(n)$.

Problem 1.3 Numerical Differentiation

Numerical differentiation aims to compute point values of the derivative of a functions approximately based on point values of the function itself. A natural idea is to use a difference quotient with a small h as replacement for the derivative:

$$f'(x) \approx \frac{f(x+h) - f(x)}{h}.$$

Analysis tells us that the smaller $|h|$ the better the approximation, but, as in [NMI, Ex. 1.1], round-off errors foil this reasoning. This problem is devoted to a detailed study of this phenomenon.

Let $f(x) = e^x$ and consider the approximation of the first derivative by means of the difference quotient.

(1.3a) Write a MATLAB-function `calcdderiv(x)` that calculates $f'(x)$ numerically for $h = 10^{-n}$, $n = 1, 2, \dots, 16$. Plot the relative error of the approximation of $f'(0)$ for $h = 10^{-n}$, $n = 1, 2, \dots, 16$, $x = 0$.

A characteristic resulting plot is in Figure 1.1.

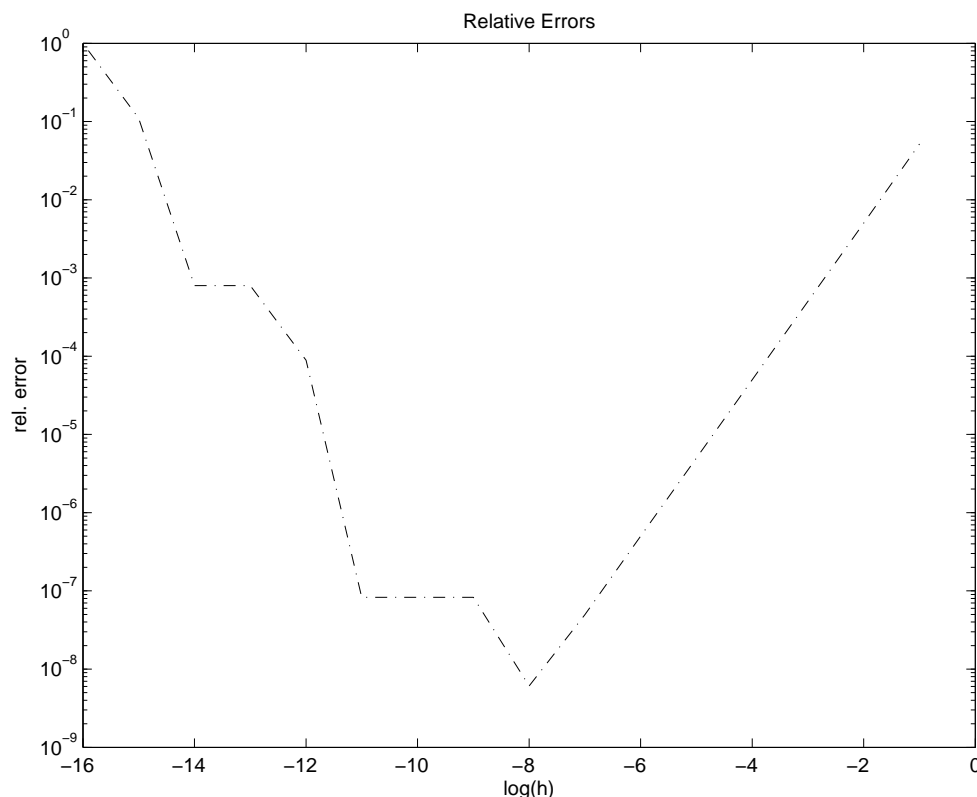


Figure 1.1: Relative error.

(1.3b) What is the value of h for which the relative error is minimized. Prove your answer analytically.

HINT: Conduct a round-off error analysis of the difference quotient based on [NMI, Def. 1.13] and [NMI, Eq. (1.8)] to derive a bound for the relative error of the approximation of the derivative by means of a computed difference quotient. Also use bounds for the remainder term of an approximation of e^x by means of a Taylor polynomial.

Problem 1.4 Round-off Error: Stable Formulation

Again, as in 1.1, we face cancellation in this problem and study ways to avoid it.

(1.4a) Consider the following functions

$$f_1(x) = 1 - \cos(x), \quad f_2(x) = 2 \sin^2\left(\frac{x}{2}\right)$$

Show analytically that f_1 and f_2 are equivalent.

HINT: Use Euler's formula, $e^{ix} = \cos(x) + i\sin(x)$ for all $x \in \mathbb{R}$.

(1.4b) Implement two MATLAB functions $y = f1(x)$ and $y = f2(x)$ to evaluate at x the functions f_1 and f_2 respectively. Plot the functions in the interval $[0, 3e-8]$ and provide a legend of the plot.

(1.4c) Explain why the graphs of the two functions differ significantly.

(1.4d) Equivalent formulations of the same function can lead to significantly different round-off errors. Consider the function

$$f(x) = \ln(\sqrt{x^2 + 1} - x), \quad x \in [10^3, 10^{18}].$$

What happens when $f(x)$ is evaluated with floating-point numbers for large values of x ? Derive analytically an equivalent formulation g of f where no subtraction occurs.

(1.4e) Implement the equivalent formulations $f(x)$ and $g(x)$ in a MATLAB routine and print the values of the two functions evaluated in $x = 10^3, 10^4, \dots, 10^{15}$. What do you observe?

Problem 1.5 Round-off Error Analysis

In this problem we delve into asymptotic round-off analysis as presented in [NMI, Sect. 1.3] and [NMI, Sect. 1.4]. The attribute asymptotic indicates that you may assume all relative errors δ introduced by elementary operation to be very small so that you can always use linearization (Taylor expansion) around zero and subsequently drop terms of size $O(\delta^2)$.

Let $|x| < 1$, the MATLAB functions `asin(x)` and `atan(x)` compute $\arcsin(x)$ and $\arctan(x)$ respectively, with relative error $\leq u(\mathbb{F})$. It holds

$$f(x) := \arctan(x) = \arcsin\left(\frac{x}{\sqrt{1+x^2}}\right) =: g(x). \quad (1.5.1)$$

(1.5a) Implement a MATLAB routine that computes and prints the values of the relative error

$$\left| \frac{g(x) - f(x)}{f(x)} \right|$$

with respect to the `atan`-function, for $x = 10^{-5}, 10^{-4}, \dots, 1$ and for $x = 10^6, 10^7, \dots, 10^{11}$. For which values of x formula (1.5.1) is unstable?

(1.5b) Gauge the propagation of round-off errors introduced by the division in $f(x)$. Compute the relative error of

$$\tilde{f}(x) = \arcsin\left(\frac{x}{\sqrt{1+x^2}}(1+\delta)\right)$$

with respect to $f(x)$. When is the error large for small values of δ ?

HINT: Use Taylor expansions.

(1.5c) Analyze the propagation of round-off errors in floating-point arithmetic by performing a complete round-off analysis of (1.5.1) as in [NMI, Sect. 1.4].

Published on February 24, 2014.

To be submitted on March 4/5, 2014.

MATLAB: Submit all files in the online system. Include the files that generate the plots. Label all your plots. Include commands to run your functions. Comment on your results.

Written exercises: Hand-in the solutions during the exercise class or in the labeled boxes in HG G 53.x.

References

[NMI] [Lecture Notes](#) for the course “Numerical Analysis I”.

Last modified on March 5, 2014