

BLG202E

HOMEWORK 1

Seyfölmülük Kutluk

150180073

Question 1)

1. The function $f_1(x_0, h) = \sin(x_0 + h) - \sin(x_0)$ can be transformed into another form, $f_2(x_0, h)$ using the trigonometric formula given below:

$$\sin(\phi) - \sin(\psi) = 2 \cos\left(\frac{\phi + \psi}{2}\right) \sin\left(\frac{\phi - \psi}{2}\right)$$

Thus, f_1 and f_2 have the same values, in exact arithmetic, for any given argument values x_0 and h .
a. Derive $f_2(x_0, h)$.

b. Suggest a formula that avoids cancellation errors for computing the approximation $(f(x_0 + h) - f(x_0))/h$ to the derivative of $f(x) = \sin(x)$ at $x = x_0$. Write a Python program that implements your formula and computes an approximation of $f'(1.2)$, for $h = 1e-20, 1e-19, \dots, 1$.

Solution 1)

$$f_1(x_0, h) = \sin(x_0 + h) - \sin(x_0)$$

we are going to implement the formula of $\sin(\phi) - \sin(\psi) = 2 \cos\left(\frac{\phi + \psi}{2}\right) \sin\left(\frac{\phi - \psi}{2}\right)$ to f_1

$$\phi = x_0 + h \quad \text{and} \quad \psi = x_0$$

$$\sin(x_0 + h) - \sin(x_0) = \sin(\phi) - \sin(\psi) = 2 \cos\left(\frac{\phi + \psi}{2}\right) \sin\left(\frac{\phi - \psi}{2}\right)$$

$$\sin(x_0 + h) - \sin(x_0) = 2 \cos\left(\frac{(2x_0 + h)}{2}\right) \sin\left(\frac{h}{2}\right)$$

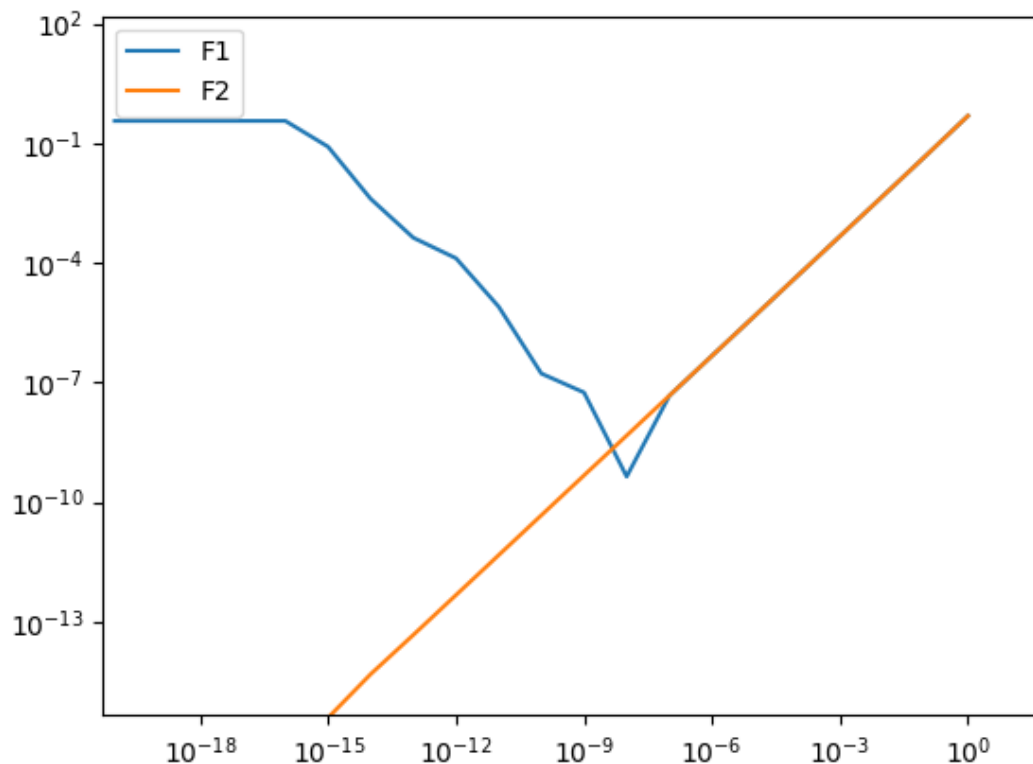
$2 \cos\left(\frac{(2x_0 + h)}{2}\right) \sin\left(\frac{h}{2}\right)$ is the f_2 and we find it from the formula above

Derivative for f_1 . using this formula $f'(x_0) \approx (f(x_0 + h) - f(x_0))/h$ is

$$(\sin(x_0 + h) - \sin(x_0))/h$$

Derivative of f_2 using this formula if $f'(x_0) \approx (f(x_0 + h) - f(x_0))/h$ is

$$(2 \cos\left(\frac{(2x_0 + h)}{2}\right) \sin\left(\frac{h}{2}\right))/h$$



We can clearly observe cancellation errors for computing the approximation $(f(x_0 + h) - f(x_0))/h$ to the derivative of $f(x) = \sin(x)$ at $x = x_0$ in f1. This occurs because of the subtracting two quantities that are close to each other.

To remove the cancellation errors, we can change the f1 to f2 with the given formula above. As can be seen in the figure there is no cancellation errors in f2.

Question 2) 2. Consider the linear system $\begin{bmatrix} a & b \\ b & a \end{bmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ with $a, b > 0$.

- 1) If $a \approx b$, what is the numerical difficulty in solving this linear system?
- 2) Suggest a numerically stable formula for computing $z = x + y$ given a and b
- 3) Determine whether the following statement is true or false, and explain why: "When $a \approx b$, the problem of solving the linear system is ill-conditioned but the problem of computing $x + y$ is not ill-conditioned."

Solution 2)

- 1) The determinant of the matrix $\begin{pmatrix} a & b \\ b & a \end{pmatrix}$ is $a^2 - b^2$ and since the $a \approx b$ the determinant is close to zero.
Since the solution of $Ax=B$ is $x=(A^{-1}) * B$
and $A^{-1} = 1/|A|$ so the result of the $\begin{pmatrix} x \\ y \end{pmatrix}$ is $(1/|A|)/B$

Since a small perturbation in the data produce a large difference in the result the problem is ill-conditioned.

$$2) \begin{pmatrix} a & b \\ b & a \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$ax+by=1 \longrightarrow ax+by-1=0$$

$$bx+ay=0 \longrightarrow bx+ay=0$$

if we add two equations we will have $ax + bx + ay + by = 1$

we will combine ax, bx in x parentheses and ay, by in y parentheses. $\longrightarrow x(a+b) + y(a+b)=1$

we will combine $x(a+b), y(a+b)$ in $a+b$ parentheses. $\longrightarrow (a+b)(x+y)=1$

$$(x+y)=1/(a+b)$$

$$z=(x+y)$$

$$z=1/(a+b) \longrightarrow \text{Numerically stable formula}$$

3) "When $a \approx b$, the problem of solving the linear system is ill-conditioned but the problem of computing $x + y$ is not ill-conditioned."

As explained in part 1 the determinant is close to zero. And this produces an ill-conditioned problem.

However, in part 2 we observe that $x+y=1/(a+b)$ it depends on the $a + b$ and it is not ill-conditioned.

The statement is true.

Question 3) With exact rounding, we know that each elementary operation has a relative error which is bounded in terms of the rounding unit η , e.g., for two floating point numbers x and y ,

$fl(x + y) = (x + y)(1 + \epsilon)$, $|\epsilon| \leq \eta$. But is this true also for elementary functions such as \sin , \ln , and exponentiation?

Consider exponentiation, which is performed according to the formula below.

$$x^y = e^{y \ln(x)} \text{ (assuming } x > 0 \text{)}$$

Estimate the relative error in calculating x^y in floating point, assuming $fl(\ln z) = (\ln z)(1 + \epsilon)$, $|\epsilon| \leq \eta$, and that everything else is exact. Show that the sort of bound we have for elementary operations and for \ln does not hold for exponentiation when x^y is very large.

Solution 3) $fl(fl(x) \pm fl(y)) = (fl(x) \pm fl(y))(1 + \epsilon_1)$ $|\epsilon| \leq \eta$.

If we take the exact rounding of $\ln x$ for example exact rounding of the $f(x)=\ln(x)$ is $f(x)=\ln(x)(1+E)$

If we implement this to our example we can compare them and find whether

$$F1(x^y) = e^{y \ln(x)}$$

$$F1(x^y) = e^{y \ln(x)(1+E)} \text{ -----} \rightarrow \text{exact rounding}$$

$$F1(x^y) = e^{y \ln(x)} * e^{y \ln(x)E} \text{ -----} \rightarrow \text{since } e^{y \ln(x)} \text{ is equal to } x^y \text{ we can put it instead}$$

$$F1(x^y) = x^y * e^{y \ln(x)E} \text{ -----} \rightarrow \text{at this point we can compare the two of them and find the relative error}$$

Relative error is $|u-v|/u$

$$|x^y - x^y * e^{y \ln(x)E}| / x^y \text{ -----} \rightarrow \text{we will take them in } x^y \text{ parenthesis.}$$

$$|x^y (1 - e^{y \ln(x)E})| / x^y \text{ -----} \rightarrow$$

$$(1 - e^{y \ln(x)E}) \text{ -----} \rightarrow$$

$$(1 - e^{y \ln(x)E}) = 1 - 1 + \epsilon y \ln x + (1/2) * (\epsilon y \ln x)^2 + (1/6) (\epsilon y \ln x)^3 + O(\epsilon^3) \text{ -----} \rightarrow \text{Taylor expansion}$$

$$(1 - e^{y \ln(x)E}) = \epsilon y \ln x + O(\epsilon^2) \text{ -----} \rightarrow$$

$$(1 - e^{y \ln(x)E}) = \epsilon y \ln x \text{ -----} \rightarrow \text{For small numbers this is true however for large numbers this is not true}$$

Lets say $x=e^{10000}$, $y=1000$, $E=10^{-10}$

$$|x^y - x^y * e^{y \ln(x)E}| / x^y$$

$$F1(x^y) = e^{y \ln(x)}$$

$$F1(x^y) = e^{1000 \ln(e^{10000})} = e^{1000000} \text{ -----} \rightarrow \text{This is the real result}$$

$$x^y * e^{y \ln(x)E} = e^{10000000 * 10^{-10}} = e^{10^{-3}} \text{ -----} \rightarrow \text{Measured result}$$

The relative error = $|e^{10000000} - e^{10^{-3}}| / e^{10000000}$ It can be observed that the error is too much for big numbers.

Question 4) You are required by a computer manufacturer to write a library function for a given floating point system to find the cube root $y^{1/3}$ of any given positive number y . Any such relevant floating-point number can be represented as $y = a * 2^e$, where a normalized fraction ($0.5 \leq a < 1$) and e is an integer exponent. This library function must be very efficient, and it should always work. For efficiency purposes it makes sense to store some useful constants ahead of computation time, e.g., the constants $2^{1/3}$, $2^{2/3}$ and $a/3$, and should these prove useful.

a. Show how $y^{1/3}$ can be obtained, once $a^{1/3}$ has been calculated for the corresponding fraction, in at most five additional flops.

b. Derive the corresponding Newton iteration. What is the flop count per iteration?

c. How would you choose an initial approximation? Roughly how many iterations are needed? (The machine rounding unit is 2^{-52} .)

Solution 4)

Let's say given floating point is y when we take the cube root of y and $a^{(1/3)}$ is already calculated.

$$y^{1/3} = (a^{1/3}) * 2^{e/3}$$

$(a^{1/3}) = (2^{e/3}) / (y^{1/3})$ this means $(2^{e/3}) / (y^{1/3})$ is already calculated

since we can calculate $(2^{e/3})$ we can calculate $(y^{1/3})$

b)

$x_{n+1} = x_n - (f(x_n)/f'(x_n))$ the Newtons iterations rule

$$y_{n+1} = y_n - (y_n^{1/3} / y_n^{-2/3}) = y_n - 3 y_n^{-2} y_n = -2 y_n \rightarrow \text{Instead of } y \text{ I will write } y = a * 2^e$$

$$y_{n+1} = -2 y_n$$

$$a_{n+1} * 2^{e_{n+1}} = - a_{n+1} * 2^{(e_{n+1}+1)} \rightarrow n \text{ flops needed}$$

c) I will choose first iteration in the middle of the array.

Later, I will observe in which side the real number stay.

Let's say it stayed at the right I will iterate the middle part of the right part.

I will continue until the number is matched or it is too close.