

## Homework 1, due Wednesday, February 19.

1. IEEE standard requires the approximate result of an arithmetic operation (for example, addition  $a$  plus  $b$ , where plus denotes addition with machine precision is equal to the precise sum of  $a$  and  $b$  rounded to the nearest floating-point number. Find an example showing that this rule does not mean that the machine addition is associative, that is, numbers  $a, b, c$ , such that  $(a \text{ plus } b) \text{ plus } c \neq a \text{ plus } (b \text{ plus } c)$
2. Modify the program used to plot the absolute error of the forward difference  $(f(x+h) - f(x))/h$  for the derivative of  $\sin$  to plot the relative error of the central difference  $(f(x+h) - f(x-h))/(2h)$ . Plot the *relative errors* and *absolute* errors for the forward and central differences at  $x = 1.2$  (the original plot),  $x = \pi/2 - 10^{-4}$ ,  $x = \pi/2 - 10^{-10}$ . What are the best values of  $h$  in each of these cases? What is the minimal absolute and relative error?
3. Plot  $y = (2^m + 1)x - 2^m x$ ,  $x = 0 \dots 1$ , for each of  $m = 48 \dots 54$  (put the plots for all values of  $m$  on a single figure). Explain the behavior of the plots, relating it to the way floating numbers are represented (64-bit floats are used by Python). why some are closer to the line  $y = x$ ? Why do we see flat regions in the plots
4. How many distinct numbers can be represented in a floating-point system following IEEE 754 standard but with only 6 bits in mantissa and 3 bits in the exponent? Count both normalized and unnormalized numbers. What is the largest and smallest magnitude of numbers that can be represented?
5. Problems 12, 15, Chapter 2. Hint for 12: Example 2.3 from the text.
6. Problems 7,8, Chapter 4.