

Math 128A: Homework 2

Due: June 28

1. In problems where high precision is not needed, the IEEE standard provides a specification for single precision numbers, which occupy 32 bits of storage. Look up this specification on the web, such as the number of bits attributed to the mantissa and exponent, the machine precision ϵ , the range of representable numbers, etc. Give the sequence of bits representing the real number 11.25 in single precision.
2. Suppose that as $x \rightarrow a$, $f_1 = O(g_1)$ and $f_2 = O(g_2)$. Prove that
 - (a) $f_1 + f_2 = O(|g_1| + |g_2|)$.
 - (b) $f_1 f_2 = O(g_1 g_2)$.
 - (c) for any constant c , $c f_1 = O(g_1)$.

3. In class, we saw that k -digit chopping is accurate to at least $(k - 1)$ significant figures. In this problem, we find the analogous result for k -digit rounding.

- (a) Convince yourself that, given a decimal number of the form

$$a = d_1.d_2d_3 \dots d_k d_{k+1} \dots \times 10^n$$

with $d_1 \neq 0$, k -digit rounding is equivalent to k -digit chopping applied to $a + 5 \times 10^{n-k}$.

- (b) By combining the steps in the chopping proof with (a), show that k -digit rounding is accurate to at least k significant figures.

4. Write a MATLAB function of the form

```
function p = Horner(z,A)
```

that uses Horner's method to evaluate the polynomial

$$p(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0$$

at a point z with $\mathbf{A} = [a_0, a_1, a_2, \dots, a_{n-1}, a_n]$.

1. Define $b_n = a_n$.
2. Define $b_{i-1} = a_{i-1} + z b_i$ for $i = n, n - 1, \dots, 1$.
3. Output b_0 .

5. Use the MATLAB function `BisMethod.m` to solve the following equations and describe your results (that is, the interval you started from, the tolerance chosen, the number of iterations needed to get the desired accuracy and, finally, the answer).

$$(a) \ x^x = 50, \quad (b) \ \ln(x) = \cos(x), \quad (c) \ x + e^{-x^2} = 0.$$

6. An eclipse of the sun is said to be α -complete when $\alpha\%$ of the area of the sun's disc is hidden behind the disc of the moon. One can model this with two discs, each of radius 1. Use the bisection method to find the distance between the centers of the discs when the eclipse is 10-complete.

7. Transform the following equations algebraically into a fixed point form and compute solutions using the fixed point MATLAB solver `FixPoi.m`.

$$(a) \ x \ln(x) - 1 = 0.$$

$$(b) \ x + \sqrt{x} = 1 + x^2.$$

$$(c) \ 3x^2 + \tan(x) = 0.$$

8. Let $x = g(x)$ be a fixed point form for which the sequence $\{p_n\}$ diverges, because $|g'(x)| > 1$.

- (a) Prove that then the fixed point form

$$x = g^{-1}(x)$$

generates a convergent sequence to the fixed point p of $x = g(x)$.

- (b) Compute the first positive solution of $x - \tan x = 0$.

Hint: The fixed point form $x = \tan(x)$ does not converge; however, using (a), it can be argued that

$$x = \arctan(x) + \pi$$

generates a convergent sequence over an appropriate interval.

9. (i) Determine analytically the fixed point that is computed by the following iterations.
 (ii) Examine, numerically, the speeds of convergence for $a = 10$ (i.e., determine the number of iterations needed to converge for all three and rank them in order of least-to-most iterations).

$$(a) \ x_0 = 1, \ x_{n+1} = 0.2 \left(4x_n + \frac{a}{x_n} \right).$$

$$(b) \ x_0 = 1, \ x_{n+1} = 0.5 \left(x_n + \frac{a}{x_n} \right).$$

$$(c) \ x_0 = 1, \ x_{n+1} = \frac{x_n(x_n^2 + 3a)}{3x_n^2 + a}.$$

Submission Details: Turn in all problems on paper except problem 4. Upload your MATLAB file for that to bCourses.