

AMSC460: Homework 1

JAMES ZHANG^{*}

September 9, 2024

3. (a) How many distinct positive numbers can be represented in a floating point system using base $\beta = 10$, precision $t = 2$ and exponent range $L = -9, U = 10$?
(Assume normalized fractions and don't worry about underflow.)

Solution. In general, a floating point number can be expressed in the representation

$$fl(x) = \pm \left(\frac{\tilde{d}_0}{\beta^0} + \frac{\tilde{d}_1}{\beta^1} + \cdots + \frac{\tilde{d}_{t-1}}{\beta^{t-1}} \right) \times \beta^e$$

The problem statement specifies that we are looking for positive integers, $t = 2$, $\beta = 10$, and e is bounded by -9 and 10 . Applying this information, we now have the more specific representation

$$fl(x) = + \left(\tilde{d}_0 + \frac{\tilde{d}_1}{10} \right) \times 10^e$$

Since we assume normalized fractions, $\tilde{d}_0 \neq 0$, so it can attain the digits $1 - 9$, or 9 possibilities. \tilde{d}_1 can be any digit, so 10 possibilities. Finally, e can be any number from -9 to 10 , so 20 possibilities. Multiplying these together yields

$$9 \times 10 \times 20 = 1800 \text{ distinct positive integers}$$

□

^{*}Email: jzhang72@terpmail.umd.edu

13. Consider the linear system

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

with $a, b > 0; a \neq b$.

(a) If $a \approx b$, what is the numerical difficulty in solving this linear system?

Solution. To solve this system, if the square matrix is invertible, we would invert the matrix and solve for the vector $(x \ y)^T$. By the Invertible Matrix Theorem, one of the conditionings for checking if a matrix is invertible is if its determinant is nonzero. Note that the determinant of the square matrix is

$$\det \begin{pmatrix} a & b \\ b & a \end{pmatrix} = a^2 - b^2 = (a - b)(a + b)$$

If $a \approx b$, and specifically if we take the limit $a - b \rightarrow 0$,

$$\lim_{a-b \rightarrow 0} \det \begin{pmatrix} a & b \\ b & a \end{pmatrix} = 0$$

because the $a - b$ approaches 0. Therefore, the matrix is almost singular and so the system is very ill-conditioned, meaning the system output is sensitive to small changes in coefficients and that small errors in arithmetic will get quickly propagated throughout the calculations. \square

1. Apply the bisection routine `bisect` to find the root of the function

$$f(x) = \sqrt{x} - 1.1$$

starting from the interval $[0, 2]$ (that is, $a = 0$ and $b = 2$), with `atol` = 1.e-8.

- How many iterations are required? Does the iteration count match the expectations, based on our convergence analysis?
- What is the resulting absolute error? Could this absolute error be predicted by our convergence analysis?

Solution.

- Let us apply the Bisection Method

```

1      function [root, iter] = bisection_sqrt()
2          % Define the function
3          f = @(x) sqrt(x) - 1.1;
4
5          % Set the tolerance and initial interval [a, b]
6          atol = 1e-8;
7          a = 0;
8          b = 2; % Initial guess for the root search range
9
10         % Check if the interval is valid
11         if f(a) * f(b) > 0
12             error('f(a) and f(b) must have opposite signs');
13         end
14
15         iter = 0; % Counter for number of iterations
16
17         % Bisection method loop
18         while (b - a) / 2 > atol
19             iter = iter + 1;
20             c = (a + b) / 2; % Midpoint of interval
21             if f(c) == 0
22                 break; % We've found the exact root
23             elseif f(a) * f(c) < 0
24                 b = c; % Root lies in the left subinterval
25             else
26                 a = c; % Root lies in the right subinterval
27             end
28         end
29
30         root = (a + b) / 2; % Approximate root
31
32         % Display the result
33         fprintf('Root found: %.10f\n', root);
34         fprintf('Number of iterations: %d\n', iter);
35         fprintf('Error: %.10f\n', sqrt(root) - 1.1)
36     end

```

```
>> amsc460_2
Root found: 1.210000009
Number of iterations: 27
Error: 0.000000004|
```

27 iterations were required, and this does not match the expectations based on our convergence analysis

$$\text{ceil}(\log_2(b - a)/\text{atol}) = 1 \implies 28 \text{ expected iterations}$$

- (b) The resulting absolute error is $0.000000004 < \text{atol}$, and this could have been predicted using our convergence analysis.

□

2. Consider the polynomial function⁸

$$\begin{aligned} f(x) &= (x-2)^9 \\ &= x^9 - 18x^8 + 144x^7 - 672x^6 + 2016x^5 - 4032x^4 + 5376x^3 - 4608x^2 \\ &\quad + 2304x - 512. \end{aligned}$$

- (a) Write a MATLAB script which evaluates this function at 161 equidistant points in the interval $[1.92, 2.08]$ using two methods:
- Apply nested evaluation (cf. Example 1.4) for evaluating the polynomial in the expanded form $x^9 - 18x^8 + \dots$.
 - Calculate $(x-2)^9$ directly.

Plot the results in two separate figures.

- (b) Explain the difference between the two graphs.
- (c) Suppose you were to apply the bisection routine from Section 3.2 to find a root of this function, starting from the interval $[1.92, 2.08]$ and using the nested evaluation method, to an absolute tolerance 10^{-6} . *Without computing anything*, select the correct outcome:
- The routine will terminate with a root p satisfying $|p-2| \leq 10^{-6}$.
 - The routine will terminate with a root p *not* satisfying $|p-2| \leq 10^{-6}$.
 - The routine will not find a root.

Justify your choice in one short sentence.

Solution.

