

# Math 5601 Homework 1

Jacob Hauck

September 5, 2023

## Problem 1.

---

(a) See `bisect.m` – also copied here for convenience.

```
1 function result = bisect(f, a, b, epsilon, epsilon_f, max_it)
2     for k = 0:max_it
3         x_k = (a + b) / 2;
4         fk = f(x_k);
5         max_error = (b - a) / 2;
6
7         fprintf( ...
8             "k = %d, x_k = %.5g, max error = %.5g, f(x_k) = %.5g\n", ...
9             k, x_k, max_error, fk ...
10        );
11
12        if (b - a) / 2 < epsilon || abs(fk) < epsilon_f
13            break;
14        elseif f(a) * fk < 0 % root lies in [a, x_k]
15            b = x_k;
16        else % if root is not in [a, x_k], it must be in [x_k, b]
17            a = x_k;
18        end
19    end
20
21    result = x_k;
```

(b) (1) The following is copied from MATLAB output. For  $\varepsilon = 10^{-2}$ :

```
>> bisect(@(x) atan(x), -4.9, 5.1, 1e-2, 0, 50)
k = 0, x_k = 0.1, max error = 5, f(x_k) = 0.099669
k = 1, x_k = -2.4, max error = 2.5, f(x_k) = -1.176
k = 2, x_k = -1.15, max error = 1.25, f(x_k) = -0.85505
k = 3, x_k = -0.525, max error = 0.625, f(x_k) = -0.48345
k = 4, x_k = -0.2125, max error = 0.3125, f(x_k) = -0.20939
k = 5, x_k = -0.05625, max error = 0.15625, f(x_k) = -0.056191
k = 6, x_k = 0.021875, max error = 0.078125, f(x_k) = 0.021872
k = 7, x_k = -0.017188, max error = 0.039062, f(x_k) = -0.017186
k = 8, x_k = 0.0023437, max error = 0.019531, f(x_k) = 0.0023437
k = 9, x_k = -0.0074219, max error = 0.0097656, f(x_k) = -0.0074217

ans =

-0.0074
```

For  $\varepsilon = 10^{-4}$ :

```
>> bisection(@(x) atan(x), -4.9, 5.1, 1e-4, 0, 50)
k = 0, x_k = 0.1, max error = 5, f(x_k) = 0.099669
k = 1, x_k = -2.4, max error = 2.5, f(x_k) = -1.176
k = 2, x_k = -1.15, max error = 1.25, f(x_k) = -0.85505
k = 3, x_k = -0.525, max error = 0.625, f(x_k) = -0.48345
k = 4, x_k = -0.2125, max error = 0.3125, f(x_k) = -0.20939
k = 5, x_k = -0.05625, max error = 0.15625, f(x_k) = -0.056191
k = 6, x_k = 0.021875, max error = 0.078125, f(x_k) = 0.021872
k = 7, x_k = -0.017188, max error = 0.039062, f(x_k) = -0.017186
k = 8, x_k = 0.0023437, max error = 0.019531, f(x_k) = 0.0023437
k = 9, x_k = -0.0074219, max error = 0.0097656, f(x_k) = -0.0074217
k = 10, x_k = -0.0025391, max error = 0.0048828, f(x_k) = -0.0025391
k = 11, x_k = -9.7656e-05, max error = 0.0024414, f(x_k) = -9.7656e-05
k = 12, x_k = 0.001123, max error = 0.0012207, f(x_k) = 0.001123
k = 13, x_k = 0.0005127, max error = 0.00061035, f(x_k) = 0.0005127
k = 14, x_k = 0.00020752, max error = 0.00030518, f(x_k) = 0.00020752
k = 15, x_k = 5.4932e-05, max error = 0.00015259, f(x_k) = 5.4932e-05
k = 16, x_k = -2.1362e-05, max error = 7.6294e-05, f(x_k) = -2.1362e-05

ans =

-2.1362e-05
```

For  $\epsilon = 10^{-8}$ :

```
>> bisection(@(x) atan(x), -4.9, 5.1, 1e-8, 0, 50)
k = 0, x_k = 0.1, max error = 5, f(x_k) = 0.099669
k = 1, x_k = -2.4, max error = 2.5, f(x_k) = -1.176
k = 2, x_k = -1.15, max error = 1.25, f(x_k) = -0.85505
k = 3, x_k = -0.525, max error = 0.625, f(x_k) = -0.48345
k = 4, x_k = -0.2125, max error = 0.3125, f(x_k) = -0.20939
k = 5, x_k = -0.05625, max error = 0.15625, f(x_k) = -0.056191
k = 6, x_k = 0.021875, max error = 0.078125, f(x_k) = 0.021872
k = 7, x_k = -0.017188, max error = 0.039062, f(x_k) = -0.017186
k = 8, x_k = 0.0023437, max error = 0.019531, f(x_k) = 0.0023437
k = 9, x_k = -0.0074219, max error = 0.0097656, f(x_k) = -0.0074217
k = 10, x_k = -0.0025391, max error = 0.0048828, f(x_k) = -0.0025391
k = 11, x_k = -9.7656e-05, max error = 0.0024414, f(x_k) = -9.7656e-05
k = 12, x_k = 0.001123, max error = 0.0012207, f(x_k) = 0.001123
k = 13, x_k = 0.0005127, max error = 0.00061035, f(x_k) = 0.0005127
k = 14, x_k = 0.00020752, max error = 0.00030518, f(x_k) = 0.00020752
k = 15, x_k = 5.4932e-05, max error = 0.00015259, f(x_k) = 5.4932e-05
k = 16, x_k = -2.1362e-05, max error = 7.6294e-05, f(x_k) = -2.1362e-05
k = 17, x_k = 1.6785e-05, max error = 3.8147e-05, f(x_k) = 1.6785e-05
k = 18, x_k = -2.2888e-06, max error = 1.9073e-05, f(x_k) = -2.2888e-06
k = 19, x_k = 7.2479e-06, max error = 9.5367e-06, f(x_k) = 7.2479e-06
k = 20, x_k = 2.4796e-06, max error = 4.7684e-06, f(x_k) = 2.4796e-06
k = 21, x_k = 9.5367e-08, max error = 2.3842e-06, f(x_k) = 9.5367e-08
k = 22, x_k = -1.0967e-06, max error = 1.1921e-06, f(x_k) = -1.0967e-06
k = 23, x_k = -5.0068e-07, max error = 5.9605e-07, f(x_k) = -5.0068e-07
k = 24, x_k = -2.0266e-07, max error = 2.9802e-07, f(x_k) = -2.0266e-07
k = 25, x_k = -5.3644e-08, max error = 1.4901e-07, f(x_k) = -5.3644e-08
k = 26, x_k = 2.0862e-08, max error = 7.4506e-08, f(x_k) = 2.0862e-08
k = 27, x_k = -1.6391e-08, max error = 3.7253e-08, f(x_k) = -1.6391e-08
k = 28, x_k = 2.2352e-09, max error = 1.8626e-08, f(x_k) = 2.2352e-09
k = 29, x_k = -7.0781e-09, max error = 9.3132e-09, f(x_k) = -7.0781e-09
```

```
ans =
-7.0781e-09
```

(2) The maximum error  $M_k$  after  $k$  iterations of the bisection method is given by

$$M_k = \frac{b-a}{2^{k+1}} \quad (1)$$

To obtain a maximum error less than  $\varepsilon > 0$ , we need that  $k$  satisfies the inequality

$$M_k < \varepsilon \iff \frac{b-a}{2^{k+1}} < \varepsilon \quad (2)$$

Thus, we need

$$k > \log_2 \left( \frac{b-a}{2\varepsilon} \right) \quad (3)$$

Since  $k$  must be an integer, the least number of iterations needed to guarantee an error no greater than  $\varepsilon$  is given by the ceiling of the left side of (3), that is, the smallest integer greater than LHS(3):

$$k = \left\lceil \log_2 \left( \frac{b-a}{2\varepsilon} \right) \right\rceil \quad (4)$$

For  $[a, b] = [-4.9, 5.1]$  and  $\varepsilon = 10^{-2}$ , this gives  $k = \lceil 8.9658 \rceil = 9$ ; for  $\varepsilon = 10^{-4}$ , it gives  $k = \lceil 15.6096 \rceil = 16$ ; and for  $\varepsilon = 10^{-8}$  it gives  $k = \lceil 28.8974 \rceil = 29$ . These are exactly the number of iterations that were executed in the numerical experiments.

## Problem 2.

(a) See `fixed.m` – also copied here for convenience.

```
1 function result = fixed(g, x0, epsilon, epsilon_f, max_it)
2   x_k = x0;
3   x_next = g(x_k);
4   fprintf("k = 0, x_k = %.5g, error = unknown, f(x_k) = %.5g\n", x_k, x_next);
5
6   for k = 1:max_it
7     x_k = x_next;
8     x_next = g(x_k);
9
10    fprintf( ...
11      "k = %d, x_k = %.5g, error = %.5g, f(x_k) = %.5g\n", ...
12      k, x_k, abs(x_next - x_k), x_next ...
13    )
14
15    if abs(x_next - x_k) < epsilon || abs(x_next) < epsilon_f
16      break;
17    end
18  end
19
20  result = x_k;
```

The following outputs are copied from MATLAB. For  $x_0 = 5$ :

```
>> fixed(@(x) x - atan(x), 5, 0, 0, 10)
k = 0, x_k = 5, error = unknown, f(x_k) = 3.6266
k = 1, x_k = 3.6266, error = 1.3017, f(x_k) = 2.3249
k = 2, x_k = 2.3249, error = 1.1646, f(x_k) = 1.1603
k = 3, x_k = 1.1603, error = 0.85945, f(x_k) = 0.30082
k = 4, x_k = 0.30082, error = 0.29221, f(x_k) = 0.008611
k = 5, x_k = 0.008611, error = 0.0086108, f(x_k) = 2.1282e-07
k = 6, x_k = 2.1282e-07, error = 2.1282e-07, f(x_k) = 3.2028e-21
k = 7, x_k = 3.2028e-21, error = 3.2028e-21, f(x_k) = 0
k = 8, x_k = 0, error = 0, f(x_k) = 0
k = 9, x_k = 0, error = 0, f(x_k) = 0
k = 10, x_k = 0, error = 0, f(x_k) = 0

ans =

0
```

For  $x_0 = -5$ :

```
>> fixed(@(x) x - atan(x), -5, 0, 0, 10)
k = 0, x_k = -5, error = unknown, f(x_k) = -3.6266
k = 1, x_k = -3.6266, error = 1.3017, f(x_k) = -2.3249
k = 2, x_k = -2.3249, error = 1.1646, f(x_k) = -1.1603
k = 3, x_k = -1.1603, error = 0.85945, f(x_k) = -0.30082
k = 4, x_k = -0.30082, error = 0.29221, f(x_k) = -0.008611
k = 5, x_k = -0.008611, error = 0.0086108, f(x_k) = -2.1282e-07
k = 6, x_k = -2.1282e-07, error = 2.1282e-07, f(x_k) = -3.2028e-21
k = 7, x_k = -3.2028e-21, error = 3.2028e-21, f(x_k) = 0
k = 8, x_k = 0, error = 0, f(x_k) = 0
k = 9, x_k = 0, error = 0, f(x_k) = 0
k = 10, x_k = 0, error = 0, f(x_k) = 0

ans =

0
```

For  $x_0 = 1$ :

```
>> fixed(@(x) x - atan(x), 1, 0, 0, 10)
k = 0, x_k = 1, error = unknown, f(x_k) = 0.2146
k = 1, x_k = 0.2146, error = 0.2114, f(x_k) = 0.0032063
k = 2, x_k = 0.0032063, error = 0.0032063, f(x_k) = 1.0987e-08
k = 3, x_k = 1.0987e-08, error = 1.0987e-08, f(x_k) = 0
k = 4, x_k = 0, error = 0, f(x_k) = 0
k = 5, x_k = 0, error = 0, f(x_k) = 0
k = 6, x_k = 0, error = 0, f(x_k) = 0
k = 7, x_k = 0, error = 0, f(x_k) = 0
k = 8, x_k = 0, error = 0, f(x_k) = 0
k = 9, x_k = 0, error = 0, f(x_k) = 0
k = 10, x_k = 0, error = 0, f(x_k) = 0

ans =

0
```

For  $x_0 = -1$ :

```
>> fixed(@(x) x - atan(x), -1, 0, 0, 10)
k = 0, x_k = -1, error = unknown, f(x_k) = -0.2146
k = 1, x_k = -0.2146, error = 0.2114, f(x_k) = -0.0032063
k = 2, x_k = -0.0032063, error = 0.0032063, f(x_k) = -1.0987e-08
k = 3, x_k = -1.0987e-08, error = 1.0987e-08, f(x_k) = 0
k = 4, x_k = 0, error = 0, f(x_k) = 0
k = 5, x_k = 0, error = 0, f(x_k) = 0
k = 6, x_k = 0, error = 0, f(x_k) = 0
k = 7, x_k = 0, error = 0, f(x_k) = 0
k = 8, x_k = 0, error = 0, f(x_k) = 0
k = 9, x_k = 0, error = 0, f(x_k) = 0
k = 10, x_k = 0, error = 0, f(x_k) = 0

ans =

0
```

For  $x_0 = 0.1$ :

```
>> fixed(@(x) x - atan(x), 0.1, 0, 0, 10)
k = 0, x_k = 0.1, error = unknown, f(x_k) = 0.00033135
k = 1, x_k = 0.00033135, error = 0.00033135, f(x_k) = 1.2126e-11
k = 2, x_k = 1.2126e-11, error = 1.2126e-11, f(x_k) = 0
k = 3, x_k = 0, error = 0, f(x_k) = 0
k = 4, x_k = 0, error = 0, f(x_k) = 0
k = 5, x_k = 0, error = 0, f(x_k) = 0
k = 6, x_k = 0, error = 0, f(x_k) = 0
k = 7, x_k = 0, error = 0, f(x_k) = 0
k = 8, x_k = 0, error = 0, f(x_k) = 0
k = 9, x_k = 0, error = 0, f(x_k) = 0
k = 10, x_k = 0, error = 0, f(x_k) = 0

ans =

0
```

- (b) It appears that the algorithm converges to 0 from all the initial guesses that I experimented with. The ones that start closer to 0 are a few iterations ahead of the ones that start farther from 0.

First, set  $G = [-R, R]$ , where  $R > 0$  is large enough that  $x_0 \in [-R, R]$ . By the Fundamental Theorem of Calculus, if  $x \in G$ , then

$$|g(x)| = |x - \tan^{-1}(x)| = \left| \int_0^x \left( 1 - \frac{1}{1+t^2} \right) dt \right| \leq |x| \leq R \quad (5)$$

Therefore,  $g(G) \subseteq G$ . Furthermore,  $g$  is  $L$ -Lipschitz on  $[-R, R]$  with  $L = 1 - \frac{1}{1+R^2} < 1$  because

$$g'(x) = 1 - \frac{1}{1+x^2} \leq 1 - \frac{1}{1+R^2} \quad (6)$$

if  $x \in G$ . Therefore,  $g$  is a contraction on  $G$ , so the fixed point method must converge for any  $x_0 \in G$ . Since  $G = [-R, R]$ , and  $R > 0$  was arbitrary, it follows that the fixed point method should converge for all initial guesses.

Second, if the initial guess  $x_0$  is farther from the fixed point  $z = 0$ , then the error bound

$$|x_k - z| \leq \frac{L^k}{1-L} |x_1 - x_0| \quad (7)$$

is looser as  $L$  gets bigger, and we need to choose a bigger  $L$  when  $x_0$  is farther from 0 because we need to choose  $R$  large enough so that  $x_0 \in [-R, R]$  in order for the fixed point theorem to apply with the initial guess  $x_0$ . The looser bound for  $x_0$  farther from 0 suggests that the algorithm will require more iterations when  $x_0$  is farther from 0.

---

**Problem 3.**


---

First, we need to show that  $g(G) \subseteq G$ . Note that

$$g'(x) = \frac{1}{3} \left( x^2 - 2x - \frac{5}{4} \right), \quad g''(x) = \frac{2}{3}(x - 1) \quad (8)$$

The roots of  $g'$  are  $1 \pm \frac{1}{2}\sqrt{4+5} = 1 \pm \frac{3}{2}$ , and the only root of  $g''$  is 1. Since  $1 \pm \frac{3}{2} \notin [0, 2]$ , the Extreme Value Theorem implies that

$$\max_{x \in G} g(x) = \max\{g(0), g(2)\} = \max\left\{\frac{4}{3}, \frac{1}{18}\right\} \quad (9)$$

$$\min_{x \in G} g(x) = \min\{g(0), g(2)\} = \min\left\{\frac{4}{3}, \frac{1}{18}\right\} = \frac{1}{18} \quad (10)$$

Therefore,  $g(G) \subseteq \left[\frac{1}{18}, \frac{4}{3}\right] \subseteq G$ . Furthermore, the Extreme Value Theorem also implies that

$$\max_{x \in G} g'(x) = \max\{g'(0), g'(2), g'(1)\} = \max\left\{-\frac{5}{12}, -\frac{9}{12}\right\} = -\frac{5}{12} \quad (11)$$

$$\min_{x \in G} g'(x) = \min\{g'(0), g'(2), g'(1)\} = \min\left\{-\frac{5}{12}, -\frac{9}{12}\right\} = -\frac{9}{12} \quad (12)$$

Therefore,  $|g'| \leq \frac{9}{12}$  on  $G$ , so  $g$  is  $L$ -Lipschitz on  $G$  with  $L = \frac{9}{12} < 1$ . By the Contraction Mapping Theorem, there is a unique fixed point  $z$  of  $g$  on  $G$ , and for any  $x_0 \in G$ , the sequence  $\{x_k\}_{k=0}^{\infty}$  defined recursively by  $x_{k+1} = g(x_k)$  converges to  $z$  as  $k \rightarrow \infty$ .