

# Assignment 1

September 5, 2023 1:01 PM

1. Given  $x \in \mathbb{R}$ , truncation to  $k$  digits is equivalent to truncating

to  $k-1$  digits after the decimal due to the machine representation of the mantissa.

The truncation of  $x$  (denoted  $\tilde{x}$ ) can thus be written as  $\tilde{x} = \frac{\lfloor \beta^{k-1} x \rfloor}{\beta^{k-1}}$ .

It should also be noted that  $x \leq \frac{\lceil \beta^{k-1} x \rceil}{\beta^{k-1}}$ . Therefore  $x - \tilde{x} \leq \frac{\lceil \beta^{k-1} x \rceil}{\beta^{k-1}} - \frac{\lfloor \beta^{k-1} x \rfloor}{\beta^{k-1}} = \frac{1}{\beta^{k-1}} (\lceil \beta^{k-1} x \rceil - \lfloor \beta^{k-1} x \rfloor)$ .

Since  $\forall a \in \mathbb{R} \lceil a \rceil - \lfloor a \rfloor \leq 1$  and  $\forall x \in \mathbb{R} x - \tilde{x} \leq \beta^{-k}$ ,  $\epsilon_{mach} = \beta^{-k}$ .  $\square$

$$2. x_1 + dx_2 = 1 \text{ \& } dx_1 + x_2 = 0 \Rightarrow \begin{pmatrix} 1 & d \\ d & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \Rightarrow \begin{pmatrix} 1 & d \\ 0 & 1-d^2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ -d \end{pmatrix} \Rightarrow \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 + \frac{d^2}{1-d^2} \\ -\frac{d}{1-d^2} \end{pmatrix}.$$

This means the problem is well posed. Now for the condition number:

$$\text{cond} = \frac{\|f(d+\Delta d) - f(d)\|_\infty}{\|f(d)\|_\infty} \frac{\|\Delta d\|_\infty}{\|\Delta d\|_\infty}$$

$$= \left\| \begin{pmatrix} \frac{1}{1-(d+\Delta d)^2} \\ \frac{-(d+\Delta d)}{1-(d+\Delta d)^2} \end{pmatrix} - \begin{pmatrix} \frac{1}{1-d^2} \\ \frac{-d}{1-d^2} \end{pmatrix} \right\|_\infty \frac{|\Delta d|}{|\Delta d| \left( \frac{1}{1-d^2} \right)}$$

$$= \frac{|\Delta d| (1-d^2)}{|\Delta d|} \left\| \frac{1}{(1-(d+\Delta d)^2)(1-d^2)} \begin{pmatrix} 1-d^2-1+(d+\Delta d)^2 \\ -(-1-d^2)(d+\Delta d)+d(1-(d+\Delta d)^2) \end{pmatrix} \right\|_\infty$$

$$= \frac{|\Delta d|}{|\Delta d| (1-(d+\Delta d)^2)} \left\| \begin{pmatrix} 2d\Delta d + \Delta d^2 \\ -(d+\Delta d-d^3-d^2\Delta d) + (d-d^3-2d^2\Delta d-d^2\Delta d^2) \end{pmatrix} \right\|_\infty$$

$$= \frac{|\Delta d|}{|\Delta d| (1-(d+\Delta d)^2)} \left\| \begin{pmatrix} (2d+\Delta d)\Delta d \\ -\Delta d - d^2\Delta d - d^2\Delta d^2 \end{pmatrix} \right\|_\infty$$

$$= \frac{|\Delta d|}{1-(d+\Delta d)^2} \left\| \begin{pmatrix} 2d+\Delta d \\ 1+d^2(1+\Delta d) \end{pmatrix} \right\|_\infty$$

$$\text{Taking } \lim_{\Delta d \rightarrow 0}, \text{ since } |1+d^2| \geq |2d|, \text{ cond} = \frac{|d|(1+d^2)}{1-d^2}. \square$$

Since cond is symmetric & monotonic on  $(0,1)$ , binary search can give  $|d| < 0.9125...$  (see q2.m)

3. a. See ForwardSubstitute.m

b. 1 function BackwardSubstitution(U,b)

```

2   x_n ← b_n / U_nn
3   for i = n-1 to 1 do
4       s ← b_i
5       for j = i+1 to n do
6           s ← s - U_ij x_j
7       x_i ← s / U_ii
8   return x
    
```

```

7 | L x_i ← S/A_{ii}
8 | return x

```

c. Line 2: 1 division = 1 FLOPs

Line 6:  $(n-i-1)$  multiplications & subtractions =  $2(n-i-1)$  FLOPs

Line 7: 1 division = 1 FLOPs

Total FLOPs: 
$$T = 1 + \sum_{i=1}^{n-1} (2n-2i+1) = 1 + (2n+1)(n-1) - n(n-1)$$

$$= n^2 \in O(n^2)$$

d. See Backward Substitution.m

e. See LU Decomposition.m

f. See LUSolve.m

4. a. For brevity let  $p_i \equiv p(t_i)$ ,  $q_i \equiv q(t_i)$ ,  $b_i \equiv b(t_i)$  for  $i = 0, \dots, n+1$ . Using (2) & (3), (1) becomes:

$$\frac{x_{i+1} - 2x_i + x_{i-1}}{h^2} = p_i \frac{x_{i+1} - x_{i-1}}{2h} + q_i x_i + b_i, x_0 = \alpha, x_{n+1} = \beta.$$

$$x_{i+1} \left( \frac{1}{h^2} - \frac{p_i}{2h} \right) - x_i \left( \frac{2}{h^2} + q_i \right) + x_{i-1} \left( \frac{1}{h^2} + \frac{p_i}{2h} \right) - b_i = 0$$

$$x_{i+1} \left( \frac{1}{2} p_i - 1 \right) + x_i (h^2 q_i + 2) - x_{i-1} \left( \frac{1}{2} p_i + 1 \right) = -h^2 b_i \implies \begin{cases} a_{i,i+1} = -\frac{1}{2} p_i - 1 \\ a_{i,i} = h^2 q_i + 2 \\ a_{i,i-1} = \frac{1}{2} p_i - 1 \end{cases}$$

b.  $L = \|P\|_\infty = \max_{t \in [0,1]} |p(t)|$ ,  $hL < 2$ .  $A$  is strictly row diagonal dominant when

$$|a_{i,i}| > \sum_{j \neq i} |a_{i,j}| = |a_{i,i-1}| + |a_{i,i+1}|$$

$$|h^2 q_i + 2| > \left| -\frac{1}{2} p_{i+1} - 1 \right| + \left| \frac{1}{2} p_{i-1} - 1 \right| = 2 \quad \left( \text{since } \max_{t \in [0,1]} |p(t)| = L, \frac{hL}{2} < 1, \text{ and } |t-1| + |t-1| = 2 \forall t \in [0,1] \right)$$

Since  $h^2 > 0$  and  $q_i > 0$  the above always holds, so  $A$  is strictly row diagonal dominant.  $\square$

c.  $L = 20$  so  $h < \frac{1}{10}$  and  $\boxed{n_0 = 10}$ . There does appear to be convergence as  $n \rightarrow \infty$  (see q4c.png)

$$d. B_j = I - DA \implies B_{j,i} = 0, B_{j,i,i-1} = \frac{\frac{1}{2} p_i - 1}{h^2 q_i + 2}, B_{j,i,i+1} = \frac{\frac{1}{2} p_{i+1} - 1}{h^2 q_i + 2}$$

$$\implies \|B_j\|_\infty = \max_{i \in \{1, \dots, n\}} \left( \frac{2}{h^2 q_i + 2} \right) \implies \boxed{\|B_j\|_\infty < \frac{2}{h^2 q_{\min} + 2}}$$

Since  $q_{\min} > 0$  &  $h > 0$ ,  $\frac{2}{h^2 q_i + 2} \in (0, 1)$  so it always converges in this case.  $\square$

$\lim_{h \rightarrow 0} \|B_j\|_\infty = 1$  so it gets very slow as  $h \rightarrow 0$ .

e. This took  $\boxed{102958}$  iterations ( $\|B_j\|_\infty > 0.999992$ ). It is way faster here. When  $q = 10^8(t^2+1)$ ,  $p = \frac{1}{20} \approx 0.05$  ( $\pi/2$ ) it took  $\boxed{7}$  iterations. This is due to the fact that  $h^2 q_{\min} \sim 10^8$  so  $\|B_j\|_\infty \approx 1/q$ . (see q4e1.png & q4e2.png)  
The LU method was slower here since it had to perform  $O(n^3)$  FLOPs while Jacobi effectively only needed  $7 \cdot O(n)$  operations.

5. First note that  $f(x) = \begin{pmatrix} 1 \\ x \\ x^2 \end{pmatrix} \cdot c$ . If  $y_i = f(x_i) \forall i \in \{1, \dots, n\}$ , then  $\begin{pmatrix} 1 & x_1 & x_1^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{pmatrix} \begin{pmatrix} c_0 \\ c_1 \\ c_2 \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$  is required.

Since  $M$  has more rows than columns, it must be overdetermined.  $\square$

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 & x_n & x_n^2 \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} \begin{pmatrix} y_n \end{pmatrix}$$

Since  $M$  has more rows than columns, it must be overdetermined.

b.  $A^T = (M^T M)^T = M^T (M^T)^T = M^T M = A$   $\square$

c.  $Mc = y \Rightarrow M^T M c = M^T y \Rightarrow \boxed{b = Ay}$

d. See [LeastSquaresQuadratic.m](#).

e. The best fit quadratic was  $y = -19.4257 + 5.7451x - 0.2841x^2$  with  $\|\tilde{y} - y\| = 2347726$ . (see [q5e.m](#)).

## Question 2

```
1 function q2()
2     fprintf("%d\n", BinarySearch(@(d) abs(d) * (1 + d^2) / (1 - d^2), 10, 0, 1, 1e-6));
3 end
```

Listing 1: Code for question 2 calls BinarySearch.m.

```
1 function x = BinarySearch(f, goal, min, max, tol)
2     x = (min+max)/2;
3     y = f(x);
4     while abs(y-goal) > tol
5         if y < goal
6             min = x;
7         else
8             max = x;
9         end
10        x = (min+max)/2;
11        y = f(x);
12    end
13 end
```

Listing 2: Code for Binary search called in q2.m.

## Question 3

```
1 function y = ForwardSubstitute(L, b)
2     % Inputs:
3     %   L: A square lower triangular matrix
4     %   b: a vector of the length as the columns of U
5     % Outputs:
6     %   y: A vector that is the solution to Ly = b.
7     n = size(L, 1);
8     % obtain the first element of y
9     y = b / L(1, 1);
10    % obtain the other elements of y by iterating downwards in L.
11    for k=2:n
12        y(k) = (b(k) - L(k, 1:k-1)*y(1:k-1));
13    end
14 end
```

Listing 3: Code used for question 3a.

```
1 function x = BackwardSubstitute(U, b)
2     % Inputs:
3     %   U: A square upper triangular matrix
4     %   b: a vector of the length as the columns of U
5     % Outputs:
6     %   x: A vector that is the solution to Ux = b.
7     n = size(U, 1);
8     % obtain the last element of x
9     x = b / U(n, n);
10    % obtain the other elements of x by iterating upwards in U.
11    for k=n-1:-1:1
12        x(k) = (b(k) - U(k, k+1:n) * x(k+1:n));
13    end
14 end
```

Listing 4: Code used for question 3d.

```

1 function A = LUDecomposition(A)
2 % Inputs:
3 %   A: A square matrix
4 % Outputs:
5 %   A: A square matrix whose upper and lower triangle (below the
6 %   diagonal) are the respective U and L such that A = LU.
7 n = size(A, 1);
8 for k=1:n-1
9     for i=k+1:n
10        A(i, k) = A(i, k) / A(k, k);
11        for j=k+1:n
12            A(i, j) = A(i, j) - A(i, k) * A(k, j);
13        end
14    end
15 end
16 end

```

Listing 5: Code used for question 3e.

```

1 function x = LUSolve(A, b)
2 % Inputs:
3 %   A: A square matrix
4 %   b: a vector of the length as the columns of A
5 % Outputs:
6 %   x: A vector that is the solution to Ax = b done by LU Decomposition.
7 A = LUDecomposition(A);
8 % Perform the LU Decomposition of A
9 L = tril(A, -1) + eye(size(A, 1));
10 U = triu(A);
11 % Solve Ly = b and Ux = y in order to solve Ax = b
12 x = BackwardSubstitute(U, ForwardSubstitute(L, b));
13 end

```

Listing 6: Code used for question 3f.

## Question 4

```

1 function q4c()
2 % Code to produce graphs for A1.4
3 n = arrayfun(@(x) 10 * 2^x, 0:5);
4 for k=0:length(n)-1
5     [T, X, ~] = ODEsolve(@(t) 100 * (t^2 + 1), ...
6         @(t) 20 * cos(pi * t), ...
7         @(t) sin(pi * t), ...
8         2, -1, 0, 1, n(k+1), "LU");
9     plot(T, X, 'DisplayName', "n = " + k);
10    hold on;
11 end
12 legend()
13 end

```

Listing 7: Code used for question 4c.

```

1 function q4e()
2 addpath 'C:\Users\rjust\f2023\math578'
3 % Code to produce graphs for A1.4e
4 % [T, X, iter] = ODEsolve(@(t) 1000 * (t^2 + 1), ...
5 %     @(t) 20 * cos(pi * t), ...
6 %     @(t) sin(pi * t), ...
7 %     2, -1, 0, 1, 2500, "J");
8 % [T, X, iter] = ODEsolve(@(t) 10^8 * (t^2 + 1), ...
9 %     @(t) (1/20) * cos(pi * t), ...
10 %     @(t) sin(pi * t), ...
11 %     2, -1, 0, 1, 2500, "J");
12 [T, X, iter] = ODEsolve(@(t) 10^8 * (t^2 + 1), ...
13     @(t) (1/20) * cos(pi * t), ...
14     @(t) sin(pi * t), ...
15     2, -1, 0, 1, 2500, "LU");
16 plot(T, X);
17 fprintf("%i\n", iter);
18 end

```

Listing 8: Code used for question 4e. It was run three times by uncommenting lines 3-6, 7-10, and 12-15 to produce the respective left and right subplots of Figure 2 (The last two blocks produce identical plots).

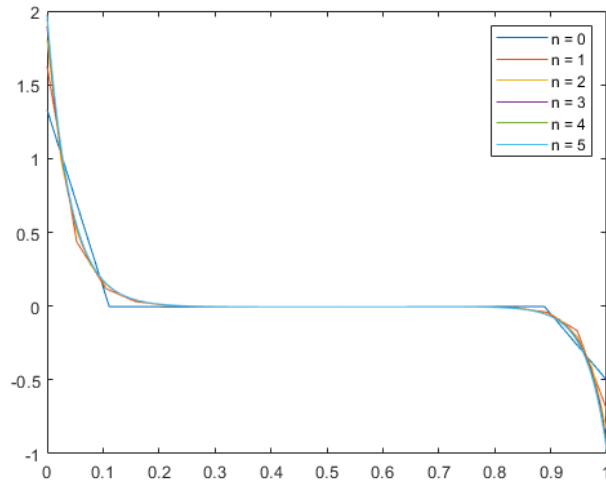


Figure 1: Graph of solution to the ODE for question 4c.

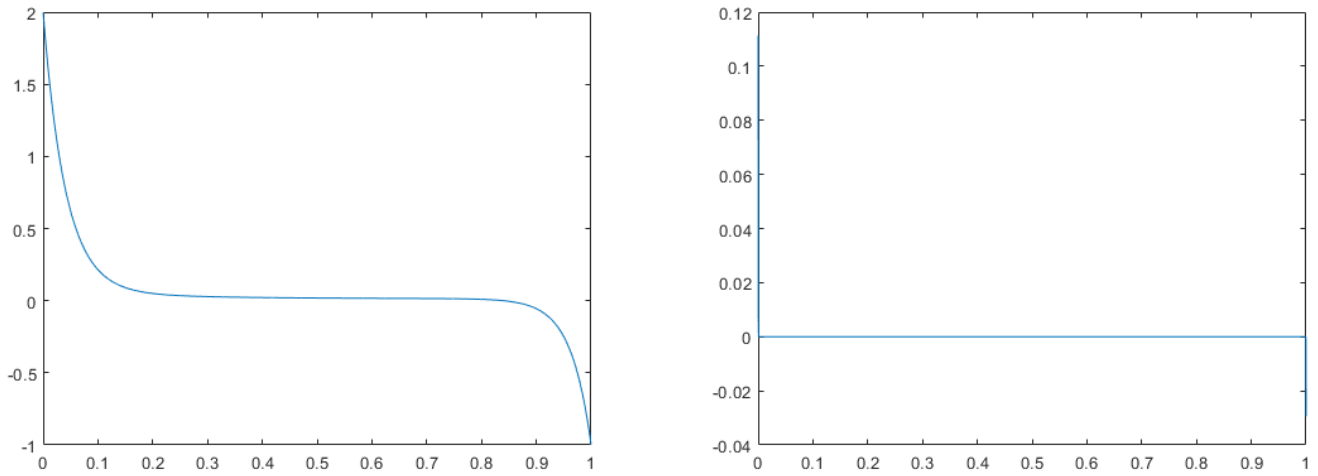


Figure 2: Graph of solution to the respective ODEs for question 4e. The new ODE given (the right subplot) presents much sharper transitions at  $x = 0$  and  $x = 1$ .

## Question 5

```
1 function c = LeastSquaresQuadratic(x, y)
2     % Inputs:
3     %   x: Vector of independant of the data
4     %   y: Vector of dependant of the data
5     % Outputs:
6     %   c: Vector representing the coefficients of the best fitting
7     %   quadratic to the data
8     % Form M from x
9     M = [ones(size(x)) x x.*x];
10
11     % Solve M^TMc = M^Ty
12     c = LUSolve(M.' * M, M.' * y);
13 end
```

Listing 9: Code used for question 5d.

```
1 function c = q5e()
2     % code to produce graphs and error margins for 1.5e
3     % form points
4     addpath 'C:\Users\rjust\f2023\math578'
5     x = [5, 5.5, 6.5, 8, 8.5, 10.8, 11.5, 13.7, 14.5, 15.9];
6     y = [1, 4, 7, 8, 9.5, 9.2, 9, 6, 3, 1];
7
8     % fitting
9     c = LeastSquaresQuadratic(x.', y. ');
10    fit_fn = @(X) c(1) + c(2)*X + c(3)*X^2;
11
12    % calculate residuals
13    fprintf("%d\n", sqrt(sum((y - arrayfun(fit_fn, x)).^2)));
14
15    % plot
16    scatter(x, y)
17    hold on;
18    t = linspace(min(x), max(x));
19    plot(t, arrayfun(fit_fn, t));
20 end
```

Listing 10: Code used for question 5e.