

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 $\forall x \in \mathbb{R} x - \tilde{x} \leq \beta^{1-k}$, $\epsilon_{mach} = \beta^{1-k}$. \square

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

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

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

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

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

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

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

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

Taking $\lim_{\Delta d \rightarrow 0}$, 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 $|d| \leq 0.9125...$ (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 $p_i \equiv p(t_i)$, $q_i \equiv (t_i)$, $b_i \equiv b(t_i)$ 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

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

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

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)

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

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

Since $q_{\min} > 0$ & $h > 0$, $\frac{2}{h^2 q_{\min} + 2} \in (0, 1)$ 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 102958 iterations ($\|B_J\|_\infty > 0.999992$). It is way faster here. When $q = 10^8(t^2+1)$, $p = \frac{1}{20} \approx \pi/6$ it took 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 $O(n)$ 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. \square

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

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](#)).