

1. (a) Let λ be an eigenvalue of A . Show that $|\lambda - a_{ii}| \leq \sum_{j \neq i} |a_{ij}|$ for some index i . (This is Gershgorin's theorem).

Proof. Let $Ax = \lambda x$, $x \neq 0$, pick index k s.t. $|x_k| = \max_i |x_i|$.

Consider the k -th row of the $Ax = \lambda x$:

$$\sum_{j=1}^n a_{kj}x_j = \lambda x_k \implies (\lambda - a_{kk})x_k = \sum_{j \neq k} a_{kj}x_j$$

Taking absolute values and using the triangle inequality, we have:

$$\begin{aligned} |\lambda - a_{kk}| |x_k| &= \left| \sum_{j \neq k} a_{kj}x_j \right| \leq \sum_{j \neq k} |a_{kj}| |x_j| \leq \sum_{j \neq k} |a_{kj}| |x_k| \\ |\lambda - a_{kk}| &\leq \sum_{j \neq k} |a_{kj}| \end{aligned}$$

□

- (b) Show that if A is strictly diagonally dominant (that is, $|a_{ii}| > \sum_{j \neq i} |a_{ij}|$ for all i), then A is invertible.

Proof. If $Ax = 0$, $x \neq 0$, pick index k s.t. $|x_k| = \max_i |x_i|$. Then,

$$\begin{aligned} |a_{kk}| |x_k| &\leq \sum_{j \neq k} |a_{kj}| |x_j| \leq \left(\sum_{j \neq k} |a_{kj}| \right) |x_k| \\ \implies |a_{kk}| &\leq \sum_{j \neq k} |a_{kj}| \end{aligned}$$

contradicting the strictly diagonally dominant condition. Thus, $Ax = 0 \Rightarrow x = 0$, so A is invertible. □

2. Consider the matrix

$$A = \begin{pmatrix} 2 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 2 \end{pmatrix}$$

- (a) What does Gershgorin's theorem say about the e-values of A ?

For the first row, we have $|\lambda - 2| \leq 1$, i.e. $\lambda \in [1, 3]$.

For the second row, we have $|\lambda - 2| \leq 2$, i.e. $\lambda \in [0, 4]$.

So, the union of all Gershgorin discs is $[0, 4]$. Thus, all eigenvalues of A lie in $[0, 4]$.

- (b) Let \mathbf{e}_i be the standard unit vector in the i th direction. Calculate the Rayleigh quotients of $\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_1 + \mathbf{e}_3, \mathbf{e}_1 - \mathbf{e}_2 + \mathbf{e}_3 - \mathbf{e}_4, \mathbf{e}_1 + \mathbf{e}_2 + \mathbf{e}_3 + \mathbf{e}_4$. Use this information to estimate λ_1 and λ_4 , the smallest and the largest e-values of A .

$$R_A(\mathbf{e}_1) = \frac{\mathbf{e}_1^T A \mathbf{e}_1}{\mathbf{e}_1^T \mathbf{e}_1} = 2$$

$$R_A(\mathbf{e}_2) = \frac{\mathbf{e}_2^T A \mathbf{e}_2}{\mathbf{e}_2^T \mathbf{e}_2} = 2$$

$$R_A(\mathbf{e}_1 + \mathbf{e}_3) = \frac{(\mathbf{e}_1 + \mathbf{e}_3)^T A (\mathbf{e}_1 + \mathbf{e}_3)}{(\mathbf{e}_1 + \mathbf{e}_3)^T (\mathbf{e}_1 + \mathbf{e}_3)} = \frac{4}{2} = 2$$

$$R_A(\mathbf{e}_1 - \mathbf{e}_2 + \mathbf{e}_3 - \mathbf{e}_4) = \frac{(\mathbf{e}_1 - \mathbf{e}_2 + \mathbf{e}_3 - \mathbf{e}_4)^T A (\mathbf{e}_1 - \mathbf{e}_2 + \mathbf{e}_3 - \mathbf{e}_4)}{(\mathbf{e}_1 - \mathbf{e}_2 + \mathbf{e}_3 - \mathbf{e}_4)^T (\mathbf{e}_1 - \mathbf{e}_2 + \mathbf{e}_3 - \mathbf{e}_4)} = \frac{14}{4} = 3.5$$

$$R_A(\mathbf{e}_1 + \mathbf{e}_2 + \mathbf{e}_3 + \mathbf{e}_4) = \frac{(\mathbf{e}_1 + \mathbf{e}_2 + \mathbf{e}_3 + \mathbf{e}_4)^T A (\mathbf{e}_1 + \mathbf{e}_2 + \mathbf{e}_3 + \mathbf{e}_4)}{(\mathbf{e}_1 + \mathbf{e}_2 + \mathbf{e}_3 + \mathbf{e}_4)^T (\mathbf{e}_1 + \mathbf{e}_2 + \mathbf{e}_3 + \mathbf{e}_4)} = \frac{2}{4} = 0.5$$

Thus, we estimate

$$0 < \lambda_1 \leq 0.5, \quad 3.5 \leq \lambda_4 \leq 4$$

- (c) Find $\|A\|_1, \|A\|_2, \|A\|_\infty, \|A\|_E, \rho(A), \text{tr}(A), \kappa_2(A)$.

$$\|A\|_1 = \max_{1 \leq j \leq 4} \sum_{i=1}^4 |a_{ij}| = 4$$

$$\|A\|_\infty = \max_{1 \leq i \leq 4} \sum_{j=1}^4 |a_{ij}| = 4$$

$$\|A\|_E = \sqrt{\sum_{i,j=1}^4 |a_{ij}|^2} = \sqrt{22}$$

$$\text{tr}(A) = \sum_{i=1}^4 a_{ii} = 8$$

$$\rho(A) = \max_{1 \leq i \leq 4} |\lambda_i| = \frac{5 + \sqrt{5}}{2}$$

$$\|A\|_2 = \sqrt{\rho(A^T A)} = \frac{5 + \sqrt{5}}{2}$$

$$\kappa_2(A) = \|A\|_2 \|A^{-1}\|_2 = \frac{\frac{5+\sqrt{5}}{2}}{\frac{3-\sqrt{5}}{2}} = 5 + 2\sqrt{5}$$

3. (a) Suppose $A\mathbf{x} = \mathbf{b}$. Let $\tilde{\mathbf{x}}$ be an approximation of the exact solution \mathbf{x} . The error is defined by $\mathbf{e} = \mathbf{x} - \tilde{\mathbf{x}}$, and the residual is $\mathbf{r} = \mathbf{b} - A\tilde{\mathbf{x}}$. Show that $A\mathbf{e} = \mathbf{r}$ and $\frac{\|\mathbf{e}\|}{\|\mathbf{x}\|} \leq \kappa(A) \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|}$.

Proof.

$$A\tilde{\mathbf{x}} = \mathbf{b} - \mathbf{r}, \quad A\mathbf{x} = \mathbf{b} \Rightarrow A(\mathbf{x} - \tilde{\mathbf{x}}) = \mathbf{r} \Rightarrow Ae = \mathbf{r}.$$

$$\|e\| = \|A^{-1}\mathbf{r}\| \leq \|A^{-1}\| \|\mathbf{r}\|, \quad \|\mathbf{b}\| = \|Ax\| \leq \|A\| \|\mathbf{x}\| \Rightarrow \frac{\|e\|}{\|\mathbf{x}\|} \leq \|A\| \|A^{-1}\| \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} = \kappa(A) \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|}.$$

□

- (b) It follows that if A is invertible, then $e = \mathbf{0}$ if and only if $\mathbf{r} = \mathbf{0}$. However, if A is ill-conditioned, the error e may be large even though the residual \mathbf{r} is small. This occurs in the example below. Show that $A\mathbf{x} = \mathbf{b}$. Find $\|e\|_\infty$ and $\|\mathbf{r}\|_\infty$ for the vectors $\tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_2$.

Find $\kappa_\infty(A)$. (Use exact arithmetics). $A = \begin{pmatrix} 1.2969 & 0.8648 \\ 0.2161 & 0.1441 \end{pmatrix}, \mathbf{b} = \begin{pmatrix} 0.8642 \\ 0.1440 \end{pmatrix}, \mathbf{x} = \begin{pmatrix} 2 \\ -2 \end{pmatrix}, \tilde{\mathbf{x}}_1 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \tilde{\mathbf{x}}_2 = \begin{pmatrix} 0.9911 \\ -0.4870 \end{pmatrix}$

Solution:

$$A = \begin{pmatrix} 1.2969 & 0.8648 \\ 0.2161 & 0.1441 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 0.8642 \\ 0.1440 \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} 2 \\ -2 \end{pmatrix}.$$

$$A\mathbf{x} = \begin{pmatrix} 1.2969 \cdot 2 + 0.8648(-2) \\ 0.2161 \cdot 2 + 0.1441(-2) \end{pmatrix} = \begin{pmatrix} 0.8642 \\ 0.1440 \end{pmatrix} = \mathbf{b}.$$

For $\tilde{\mathbf{x}}_1 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$:

$$\mathbf{e}_1 = \begin{pmatrix} 2 \\ -3 \end{pmatrix}, \quad \|\mathbf{e}_1\|_\infty = 3, \quad A\tilde{\mathbf{x}}_1 = \begin{pmatrix} 0.8648 \\ 0.1441 \end{pmatrix}, \quad \mathbf{r}_1 = \begin{pmatrix} -0.0006 \\ -0.0001 \end{pmatrix}, \quad \|\mathbf{r}_1\|_\infty = 6 \times 10^{-4}.$$

For $\tilde{\mathbf{x}}_2 = \begin{pmatrix} 0.9911 \\ -0.4870 \end{pmatrix}$:

$$\mathbf{e}_2 = \begin{pmatrix} 1.0089 \\ -1.5130 \end{pmatrix}, \quad \|\mathbf{e}_2\|_\infty = 1.5130, \\ A\tilde{\mathbf{x}}_2 = \begin{pmatrix} 0.86419999 \\ 0.14400001 \end{pmatrix}, \quad \mathbf{r}_2 = \begin{pmatrix} 10^{-8} \\ -10^{-8} \end{pmatrix}, \quad \|\mathbf{r}_2\|_\infty = 10^{-8}.$$

Condition number $\kappa_\infty(A)$.

$$\|A\|_\infty = \max\{1.2969 + 0.8648, 0.2161 + 0.1441\} = 2.1617,$$

$$\det(A) = 1.2969 \cdot 0.1441 - 0.8648 \cdot 0.2161 = 10^{-8},$$

$$A^{-1} = \frac{1}{10^{-8}} \begin{pmatrix} 0.1441 & -0.8648 \\ -0.2161 & 1.2969 \end{pmatrix}, \quad \|A^{-1}\|_\infty = \max\{1.0089, 1.5130\} \times 10^8 = 1.5130 \times 10^8,$$

$$\kappa_\infty(A) = \|A\|_\infty \|A^{-1}\|_\infty = 2.1617 \cdot 1.5130 \times 10^8 = 3.27065210 \times 10^8.$$

4. Compute the solution of the 2×2 system:

$$\begin{aligned} 10^{-3}x + y &= 5 \\ x - y &= 6 \end{aligned}$$

using standard Gaussian elimination and Gaussian elimination with pivoting. Conduct all computations with two significant digits (decimal). Compare and explain the results.

(A) Standard Gaussian elimination (no pivoting).

Start with (to two sig. figs.)

$$\left[\begin{array}{cc|c} 1.0 \times 10^{-3} & 1.0 & 5.0 \end{array} \right], \quad \left[\begin{array}{cc|c} 1.0 & -1.0 & 6.0 \end{array} \right].$$

Eliminate a_{21} using multiplier $m = \frac{1.0}{1.0 \times 10^{-3}} = 1.0 \times 10^3$.

$$(1.0 \times 10^3) \cdot \text{row}_1 = \left[\begin{array}{cc|c} 1.0 & 1.0 \times 10^3 & 5.0 \times 10^3 \end{array} \right].$$

Row operation (to two sig. figs.):

$$\text{row}_2 \leftarrow \text{row}_2 - (1.0 \times 10^3) \cdot \text{row}_1 \Rightarrow \left[\begin{array}{cc|c} 0.0 & -1.0 \times 10^3 & -5.0 \times 10^3 \end{array} \right].$$

Thus $-1.0 \times 10^3 y = -5.0 \times 10^3 \Rightarrow y = 5.0$. Back-substitute into the first equation:

$$1.0 \times 10^{-3}x + 5.0 = 5.0 \Rightarrow 1.0 \times 10^{-3}x = 0.0 \Rightarrow x = 0.0.$$

$$\boxed{(x, y)_{\text{no pivot}} = (0.0, 5.0)}.$$

(B) Gaussian elimination with partial pivoting.

Swap the rows so the first pivot is 1.0:

$$\left[\begin{array}{cc|c} 1.0 & -1.0 & 6.0 \end{array} \right], \quad \left[\begin{array}{cc|c} 1.0 \times 10^{-3} & 1.0 & 5.0 \end{array} \right].$$

Eliminate a_{21} with $m = \frac{1.0 \times 10^{-3}}{1.0} = 1.0 \times 10^{-3}$:

$$(1.0 \times 10^{-3}) \cdot \text{row}_1 = \left[\begin{array}{cc|c} 1.0 \times 10^{-3} & -1.0 \times 10^{-3} & 6.0 \times 10^{-3} \end{array} \right].$$

Row operation (rounded to two sig. figs.):

$$\text{row}_2 \leftarrow \text{row}_2 - (1.0 \times 10^{-3}) \cdot \text{row}_1 \Rightarrow \left[\begin{array}{cc|c} 0.0 & 1.0 & 5.0 \end{array} \right],$$

so $y = 5.0$. Back-substitute in the first (pivoted) row:

$$x - y = 6.0 \Rightarrow x = 11.0.$$

$$\boxed{(x, y)_{\text{pivot}} = (11.0, 5.0)}.$$

Comparison and explanation. The exact solution is $(x, y) = (10.989\dots, 4.989\dots)$. Without pivoting the tiny pivot 10^{-3} produces a huge multiplier 10^3 , forcing subtraction of nearly equal large numbers and catastrophic cancellation, yielding $x = 0.0$. Pivoting selects the well-scaled pivot 1.0, avoids amplification of rounding, and returns the correct two-digit result $(11.0, 5.0)$.

5. (a) Show that for any matrix A , there exists a permutation matrix P such that $PA = LU$, where L is a unit lower triangular matrix and U is an upper triangular matrix.

Hint: P is a composition, that is, a product of all permutation matrices used at every stage of the Gaussian elimination algorithm.

Proof. Let Gaussian elimination with (partial) pivoting act on A . At step k a permutation P_k swaps rows to place a nonzero pivot and a unit lower-triangular elimination matrix E_k zeroes entries below it. After m steps,

$$P_m \cdots P_1 A = LU, \quad E_m \cdots E_1 P_m \cdots P_1 A = U.$$

Set $P = P_m \cdots P_1$ and $L = (E_m \cdots E_1)^{-1}$ (unit lower triangular). Then $PA = LU$. \square

(b) For the following matrix A , find P, L and U . Use the LU factorization to solve $A\mathbf{x} = \mathbf{b}$.

$$A = \begin{pmatrix} 0 & 2 & -1 \\ 1 & 1 & 1 \\ 2 & 0 & 1 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} -2 \\ 1 \\ 0 \end{pmatrix}.$$

For $A = \begin{pmatrix} 0 & 2 & -1 \\ 1 & 1 & 1 \\ 2 & 0 & 1 \end{pmatrix}$, $\mathbf{b} = \begin{pmatrix} -2 \\ 1 \\ 0 \end{pmatrix}$, partial pivoting on column 1 then column 2 gives

$$P = \begin{pmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}, \quad L = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ \frac{1}{2} & \frac{1}{2} & 1 \end{pmatrix}, \quad U = \begin{pmatrix} 2 & 0 & 1 \\ 0 & 2 & -1 \\ 0 & 0 & 1 \end{pmatrix}, \quad \text{so } PA = LU.$$

Solve $A\mathbf{x} = \mathbf{b}$ via $LU\mathbf{x} = P\mathbf{b}$:

$$P\mathbf{b} = \begin{pmatrix} 0 \\ -2 \\ 1 \end{pmatrix}, \quad L\mathbf{y} = P\mathbf{b} \Rightarrow \mathbf{y} = \begin{pmatrix} 0 \\ -2 \\ 2 \end{pmatrix}, \quad U\mathbf{x} = \mathbf{y} \Rightarrow \mathbf{x} = \begin{pmatrix} -1 \\ 0 \\ 2 \end{pmatrix}.$$

6. Let M be an invertible matrix. Show that $\|\mathbf{x}\|_M := \|M\mathbf{x}\|_\infty$ defines a vector norm for which the subordinate matrix norm is $\|A\|_M := \|MAM^{-1}\|_\infty$.

Proof. 1) $\|\cdot\|_M$ is a norm.

- (Positivity) $\|\mathbf{x}\|_M = \|M\mathbf{x}\|_\infty \geq 0$, and $\|\mathbf{x}\|_M = 0 \iff M\mathbf{x} = 0 \iff \mathbf{x} = 0$ (since M is invertible).
- (Absolute homogeneity) $\|\alpha\mathbf{x}\|_M = \|M(\alpha\mathbf{x})\|_\infty = |\alpha| \|M\mathbf{x}\|_\infty = |\alpha| \|\mathbf{x}\|_M$.
- (Triangle inequality) $\|\mathbf{x} + \mathbf{y}\|_M = \|M(\mathbf{x} + \mathbf{y})\|_\infty \leq \|M\mathbf{x}\|_\infty + \|M\mathbf{y}\|_\infty = \|\mathbf{x}\|_M + \|\mathbf{y}\|_M$, using the triangle inequality of $\|\cdot\|_\infty$.

2) Subordinate matrix norm. Define the operator norm subordinate to $\|\cdot\|_M$ by

$$\|A\|_M := \sup_{x \neq 0} \frac{\|Ax\|_M}{\|x\|_M} = \sup_{x \neq 0} \frac{\|M\mathbf{A}\mathbf{x}\|_\infty}{\|M\mathbf{x}\|_\infty}.$$

Let $y = M\mathbf{x}$ (bijection since M is invertible). Then

$$\|A\|_M = \sup_{y \neq 0} \frac{\|MAM^{-1}y\|_\infty}{\|y\|_\infty} = \|MAM^{-1}\|_\infty.$$

Hence the subordinate matrix norm to $\|\cdot\|_M$ is exactly $\|A\|_M := \|MAM^{-1}\|_\infty$. \square

7. Consider the 2-point boundary value problem

$$-\phi''(x) + x^2\phi(x) = (1 + 4x + 2x^2 - x^4)e^x, \quad \phi(0) = 1, \phi(1) = 0.$$

Check that the solution is $\phi(x) = (1 - x^2)e^x$. Write a computer program to solve this problem using the second-order finite-difference scheme discussed in class. Solve the tridiagonal system using the LU factorization. Run the code with mesh size $h = 1/2^p$ for $p = 1, \dots, 14$. Output the results in the following format,

column1:	h
column2:	$\ u_h - \phi_h\ _\infty$
column3:	$\ u_h - \phi_h\ _\infty/h^2$

Describe and explain what happens to the error for small values of h .

Model problem and exact solution.

$$-\phi''(x) + x^2\phi(x) = (1 + 4x + 2x^2 - x^4)e^x, \quad \phi(0) = 1, \phi(1) = 0, \quad \phi(x) = (1 - x^2)e^x.$$

Second-order FD scheme (interior nodes $x_i = ih$, $i = 1, \dots, n$).

$$-\frac{u_{i-1} - 2u_i + u_{i+1}}{h^2} + x_i^2 u_i = f(x_i), \quad u_0 = 1, \quad u_{n+1} = 0,$$

$$A = \text{tridiag}\left(-\frac{1}{h^2}, \frac{2}{h^2} + x_i^2, -\frac{1}{h^2}\right), \quad b_i = f(x_i) + \frac{1}{h^2} \mathbf{1}_{\{i=1\}}.$$

Solve by LU. Compute $[L, U, P] = \text{lu}(A)$, then $U u = L^{-1} P b$. Form $u_h = [1; u; 0]$ and the exact vector $\phi_h = \phi(x)$. Report:

h	$\ u_h - \phi_h\ _\infty$	$\ u_h - \phi_h\ _\infty/h^2$
5.000000e-01	6.472527e-02	2.589011e-01
2.500000e-01	1.651983e-02	2.643173e-01
1.250000e-01	4.154789e-03	2.659065e-01
6.250000e-02	1.056844e-03	2.705520e-01
3.125000e-02	2.643184e-04	2.706620e-01
1.562500e-02	6.608632e-05	2.706896e-01
7.812500e-03	1.652200e-05	2.706965e-01
3.906250e-03	4.130527e-06	2.706982e-01
1.953125e-03	1.032633e-06	2.706986e-01
9.765625e-04	2.581585e-07	2.706988e-01
4.882812e-04	6.454092e-08	2.707042e-01
2.441406e-04	1.613556e-08	2.707097e-01
1.220703e-04	3.081321e-09	2.067840e-01
6.103516e-05	2.814358e-10	7.554734e-02

Behavior for small h : For h down to 2.44×10^{-4} , the ratio $\|u_h - \phi_h\|_\infty/h^2$ stays essentially constant at $\approx 2.707 \times 10^{-1}$, confirming the expected second-order convergence:

$$\|u_h - \phi_h\|_\infty \approx \alpha h^2, \quad \alpha \approx 2.707 \times 10^{-1}.$$

For the last two meshes, $h = 1.22 \times 10^{-4}$ and 6.10×10^{-5} , the ratio *drops* (to 2.07×10^{-1} and 7.55×10^{-2}). This is a finite-precision effect: the total error is well modeled by

$$E(h) \approx \underbrace{\alpha h^2}_{\text{truncation}} + \underbrace{\beta \varepsilon_{\text{mach}} h^{-2}}_{\text{roundoff/backward error}},$$

with $\varepsilon_{\text{mach}} \approx 2.22 \times 10^{-16}$ in double precision. When h is small, the two terms become comparable; if β has opposite sign to α , they partially cancel, producing the observed dip in $E(h)/h^2$. If you refine further, the roundoff term will eventually dominate and $E(h)$ will stagnate and then *increase* (so $E(h)/h^2$ rises), as the LU solve backward error scales like $O(\varepsilon_{\text{mach}} \|A\|) \sim O(\varepsilon_{\text{mach}} h^{-2})$.

Summary.

$$\begin{cases} \text{Second-order regime: } h \gtrsim 2.4 \times 10^{-4}, & E(h)/h^2 \approx 0.2707, \\ \text{Roundoff interaction: } h \lesssim 1.2 \times 10^{-4}, & \text{favorable cancellation lowers } E(h)/h^2, \\ \text{Prediction: } & \text{for even smaller } h, \text{ roundoff dominates } \Rightarrow E(h) \uparrow. \end{cases}$$