

HW 6

1) Let $A = \begin{bmatrix} 1 & 1 \\ -2 & 3 \end{bmatrix}$. Will the Jacobi iterative method converge on A , regardless of the starting vector $x^{(0)}$?

For convergence, we need all eigenvalues λ_i of $(I - D^{-1}A)$ to be less than 1 in absolute value, ie

$$|\lambda_i| < 1 \quad \text{for } i \in \{1, 2, \dots, n\}$$

We have $A = \begin{bmatrix} 1 & 1 \\ -2 & 3 \end{bmatrix}$, $D = \begin{bmatrix} 1 & 0 \\ 0 & 3 \end{bmatrix}$

$$D^{-1} = \begin{bmatrix} 1 & 0 \\ 0 & 3 \end{bmatrix} \sim \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1/3 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1/3 \end{bmatrix}$$

$$\begin{aligned} I - D^{-1}A &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \begin{bmatrix} 1 & 0 \\ 0 & 1/3 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ -2 & 3 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \begin{bmatrix} 1 & 1 \\ -2/3 & 1 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ 2/3 & 0 \end{bmatrix} \end{aligned}$$

$$\det(\lambda I - (I - D^{-1}A)) = \det \begin{bmatrix} \lambda & -1 \\ 2/3 & \lambda \end{bmatrix} = \lambda^2 + 2/3 = 0$$

$$\Rightarrow \lambda^2 = -2/3 \Rightarrow \lambda_1 = \sqrt{2}/3 i, \lambda_2 = -\sqrt{2}/3 i$$

$$|\lambda_1| = \sqrt{2}/3 \quad |\lambda_2| = \sqrt{2}/3, \quad |\lambda_1| = |\lambda_2| < 1$$

Yes, the method will converge on A regardless of $x^{(0)}$.

2) Let $B = \begin{bmatrix} 1 & 1 \\ 1 & 3 \end{bmatrix}$ and $b = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$. Take 5 steps of Gauss-Seidel

to attempt to solve $Bx = b$, starting with $x^{(0)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$; write down $x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}, x^{(5)}$. Based on what you see, do you think Gauss-Seidel will converge?

Let $x^{(0)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$. We know that

$$b_{11}x_1 + b_{12}x_2 = 1$$

$$b_{21}x_1 + b_{22}x_2 = 1$$

$$\Rightarrow x_1^{(k+1)} = \frac{1}{b_{11}} (1 - b_{12}x_2^{(k)})$$

$$x_2^{(k+1)} = \frac{1}{b_{22}} (1 - b_{21}x_1^{(k+1)})$$

Let's plug in our initial guess for $k=0$:

$$x_1^{(1)} = \frac{1}{b_{11}} (1 - b_{12}x_2^{(0)})$$

$$= \frac{1}{1} (1 - (1)(0)) = 1 - 0 = 1$$

$$x_2^{(1)} = \frac{1}{b_{22}} (1 - b_{21}x_1^{(1)})$$

$$= \frac{1}{3} (1 - (1)(1)) = \frac{1}{3}(1 - 1) = 0$$

$$\Rightarrow x^{(1)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

And repeat:

$$x_1^{(2)} = \frac{1}{b_{11}} (1 - b_{12} x_2^{(1)}) \\ = \frac{1}{1} (1 - (1)(0)) = 1(1-0) = 1$$

$$x_2^{(2)} = \frac{1}{b_{22}} (1 - b_{21} x_1^{(2)}) \\ = \frac{1}{3} (1 - (1)(1)) = \frac{1}{3}(1-1) = 0$$

$$\Rightarrow x^{(2)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$x_1^{(3)} = \frac{1}{b_{11}} (1 - b_{12} x_2^{(2)}) \\ = \frac{1}{1} (1 - (1)(0)) = 1(1-0) = 1$$

$$x_2^{(3)} = \frac{1}{b_{22}} (1 - b_{21} x_1^{(3)}) \\ = \frac{1}{3} (1 - (1)(1)) = \frac{1}{3}(1-1) = 0$$

$$\Rightarrow x^{(3)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$x_1^{(4)} = \frac{1}{b_{11}} (1 - b_{12} x_2^{(3)}) \\ = \frac{1}{1} (1 - (1)(0)) = 1$$

$$x_2^{(4)} = \frac{1}{b_{22}} (1 - b_{21} x_1^{(4)}) \\ = \frac{1}{3} (1 - (1)(1)) = \frac{1}{3}(1-1) = 0$$

$$\Rightarrow \mathbf{x}^{(4)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\begin{aligned} x_1^{(5)} &= \frac{1}{b_{11}} (1 - b_{12} x_2^{(4)}) \\ &= \frac{1}{1} (1 - (1)(0)) = 1(1-0) = 1 \end{aligned}$$

$$\begin{aligned} x_2^{(5)} &= \frac{1}{b_{22}} (1 - b_{21} x_1^{(5)}) \\ &= \frac{1}{3} (1 - (1)(1)) = 0 \end{aligned}$$

$$\Rightarrow \mathbf{x}^{(5)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

Yes, Gauss-Seidel will converge, since we have found the same result for all 5 iterations.

3) A matrix $A \in \mathbb{C}^{n \times n}$ is called skew-Hermitian if $A^* = -A$.

Show that skew-Hermitian matrices have only purely imaginary eigenvalues, i.e., any eigenvalue of a skew-Hermitian matrix has the form $\lambda = ai$ with $a \in \mathbb{R}$.

Hint: Show that iA is Hermitian, then use what we know about Hermitian matrices.

We know that $A^* = (\bar{A})^T = -A$ (if A is skew-Hermitian), that $B^* = B$ if B is Hermitian, and for any $c \in \mathbb{C}, c = a+ib, \bar{c} = a-ib$.

Let $B = iA$. Then

$$\begin{aligned}B^* &= (\bar{B})^T = (\bar{i}\bar{A})^T = (\bar{i}\bar{A})^T = (-i\bar{A})^T = -i(\bar{A})^T \\&= -i(A^*) = -i(-A) = iA = B\end{aligned}$$

Hence, $B = iA$ is Hermitian, so all eigenvalues of iA are real. Now,

$$\begin{aligned}Av &= \lambda v \\ \Rightarrow iAv &= i(\lambda v) \\ \Rightarrow (iA)v &= (i\lambda)v\end{aligned}$$

$\Rightarrow i\lambda$ is a real eigenvalue of iA

Suppose $\lambda = a$, $a \in \mathbb{R}$. Since $i\lambda$ is real, λ must have some imaginary component such that multiplying i by λ results in a real value. Hence, $\lambda = ai$ for some $a \in \mathbb{R}$.

4) Using one of the two methods of computing the SVD using eigen decompositions outlined in class, compute by hand the reduced SVD of the matrix $A = \begin{bmatrix} -1 & 1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}$.

$$A = U \Sigma V^T$$

$U (3 \times 3)$
 $\checkmark (2 \times 2)$
 $\Sigma (3 \times 2)$

$$A^T A = V \Sigma^T \Sigma V^T \Rightarrow A^T A V = V \Sigma^T \Sigma$$

$$\Rightarrow (\delta_1, v_1), (\delta_1, v_2), \dots, (\delta_r, v_r), (0, v_{r+1}), \dots, (0, v_n)$$

are eigen pairs of $A^T A$

$$A^T A = \begin{bmatrix} -1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} -1 & 1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$$

$$\det(A^T A - \lambda I) = \det \begin{bmatrix} 2-\lambda & -1 \\ -1 & 2-\lambda \end{bmatrix} = (2-\lambda)(2-\lambda) - (-1)(-1) \\ = \lambda^2 - 4\lambda + 3 = 0$$

$$\Rightarrow (\lambda-3)(\lambda-1)=0 \Rightarrow \lambda_1=3, \lambda_2=1$$

$$\Rightarrow \delta_1 = \sqrt{3}, \delta_2 = 1 \Rightarrow \varepsilon = \begin{bmatrix} \sqrt{3} & 0 \\ 0 & 1 \end{bmatrix}$$

$$(A^T A - 3I)x = 0$$

$$\Rightarrow \begin{bmatrix} 2-3 & -1 \\ -1 & 2-3 \end{bmatrix} = \begin{bmatrix} -1 & -1 \\ -1 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\Rightarrow -x_1 - x_2 = 0 \Rightarrow x_2 = -x_1 \Rightarrow x_1 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$\|x_1\|_2 = \sqrt{1^2 + (-1)^2} = \sqrt{2} \Rightarrow v_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$(A^T A - I)x = 0$$

$$\Rightarrow \begin{bmatrix} 2-1 & -1 \\ -1 & 2-1 \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\Rightarrow x_1 - x_2 = 0 \Rightarrow x_1 = x_2 \Rightarrow x_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow v_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\Rightarrow V = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}$$

We now have $\Sigma = \begin{bmatrix} \sqrt{3} & 0 \\ 0 & 1 \end{bmatrix}$, $V = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}$

$$A = U \Sigma V^T$$

$$\Rightarrow AV = V\Sigma \Rightarrow U = AV\Sigma^{-1}$$

$$\Sigma^{-1} = \begin{bmatrix} \sqrt{3} & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix} \sim \begin{bmatrix} 1 & 0 & 1/\sqrt{3} & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix} = \begin{bmatrix} \sqrt{3}/3 & 0 \\ 0 & 1 \end{bmatrix}$$

$$U = AV\Sigma^{-1} = \begin{bmatrix} -1 & 1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} \sqrt{3}/3 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} -\sqrt{6}/3 & 0 \\ \sqrt{6}/6 & \sqrt{2}/2 \\ -\sqrt{6}/6 & \sqrt{2}/2 \end{bmatrix}$$

\Rightarrow The reduced SVD of A is

$$A = \begin{bmatrix} -\sqrt{6}/3 & 0 \\ \sqrt{6}/6 & \sqrt{2}/2 \\ -\sqrt{6}/6 & \sqrt{2}/2 \end{bmatrix} \begin{bmatrix} \sqrt{3} & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \sqrt{2}/2 & -\sqrt{2}/2 \\ \sqrt{2}/2 & \sqrt{2}/2 \end{bmatrix}$$

$$U \quad \Sigma \quad V^T$$

5) Run the sample-inverse-PM.m code provided. Below is an explanation of what the code does.

The first line resets the pseudo-random number generator to default

lines 2-4 construct a random matrix A w/ prescribed eigenvalues (given in line 2)

Line 5 constructs a matrix B which shifts the eigenvalues of A by the same amount, while leaving the eigenvectors unchanged.

Line 6 inverts B and names the inverse C .

The remainder of the code applies the power method to the matrix C , starting from the vector of all ones, using 8 iteration

Run the code and examine it closely, then answer the following.

(a) What are the eigenvalues of B and C ? Start with the eigenvalues of A and see how they get changed when we transform A into B and then C .

The eigenvalues of A are $2, -1, -0.25, 0.5, 2$, and -5 .

For matrix B , subtract 0.25 from each eigenvalue, so the eigenvalues of B are $1.75, -1.25, -0.5, 0.25, 1.75$, and -5.25 . Now, C is the inverse of B , so the eigenvalues of C are the reciprocal of each of the eigenvalues of B , so the eigenvalues are $\frac{1}{1.75}, -\frac{1}{1.25}, -2, 4, \frac{1}{0.25}, -\frac{1}{5.25}$.

(b) Use reasoning, not MATLAB, to answer this question, and show your reasoning. By the end of the code, q is a very good approximation for an eigenvector of A . What is the corresponding eigenvalue λ ?

The vector q is formed by applying the power method to the matrix C , which is used to find an approximate dominant eigenvalue and associated eigenvector. Here, the corresponding eigenvalue to q is the largest eigenvalue (in magnitude) of the matrix C , which is the

inverse of the matrix B . So, the largest eigenvalue of C corresponds to the smallest eigenvalue of the matrix B , which has the same eigenvalues as A , but subtracting 0.25. Therefore, the corresponding eigenvalue to the matrix C is 4, to matrix B is 0.25, and to matrix A is 0.5. Hence, $\lambda = 0.5$.

(c) If you replaced 0.25 with -0.5 in line 5 of the script, which eigenvector of A would q approximate? What would be the corresponding eigenvalue of A ? Like in (b), use and show reasoning.

If we replaced 0.25 in line 5 with 0.5, the new eigenvalues of B would be 1.5, -1.5, -0.75, 0, 1.5, -5.5, and so the new eigenvalues of C would be $1/1.5, -1/1.5, -1/0.75, 0, 1/1.5, 1/5.5$. This changes the dominant eigenvalue of C from 4 to $-1/0.75$, which corresponds to -0.75 in B , which finally corresponds to $\lambda = -0.25$ in A .

(d) After running the code with 0.25 in line 5 of the script, show how good an approximation the vector q is to an eigenvector of A by typing in MATLAB $\text{norm}((A - \lambda * \text{eye}(6)) * q)$, where λ is the number you obtained in (b). This completes the 2-norm of the vector in parentheses; if it is small, the vector q is a very good approximation to an eigenvector with eigenvalue λ .

```

5 A = Q*D*Q';
6 B = A - 0.25*eye(6);
7 C = inv(B);
8 [V,D] = eig(A);
9 disp([V,D])
10
11 q = ones(6, 1); s = 1;
12 for j = 1:12
13     q_old = q;
14     q_new = C*q_old;
15     [~, ind] = max(abs(q_new));
16     s = q_new(ind(1));
17     q = q_new/s;
18 end
19 norm((A - 0.5*eye(6))*q)
20

```

Command Window

```

ans =
3.5804e-04

```

The 2-norm computed is very, very small, so we can conclude that q is a very good approximation to an eigenvector with eigenvalue $\lambda = 0.5$.

HW 5

1) Let A be an $n \times n$ real matrix and denote by $\|A\|_F$ the Frobenius norm of A . Recall that the Frobenius norm has the following equivalent definition:

$$\|A\|_F^2 = \sum_{ij} |a_{ij}|^2 = \text{trace}(A^T A).$$

a) Show that $\|A\|_F = \|UA\|_F$ for any real orthogonal $n \times n$ matrix U .

We have that $\|A\|_F = \text{trace}(A^T A)$, hence,

$$\begin{aligned}\|UA\|_F &= \text{trace}((UA)^T (UA)) && ((AB)^T = B^T A^T) \\ &= \text{trace}(\underbrace{A^T U^T}_{I}, \text{ since } U \text{ is } n \times n \text{ orthogonal} U A) \\ &= \text{trace}(A^T A) = \|A\|_F\end{aligned}$$

b) Show that $\|A\|_F = \|AV\|_F$ for any real orthogonal $n \times n$ matrix V .

Similarly, $\|AV\|_F = \text{trace}((AV)^T (AV))$

$$= \text{trace}(V^T A^T A V)$$

Since V is $n \times n$ orthogonal, for any eigenvalue λ of a matrix $M_{n \times n}$, λ is also an eigenvalue of $V^T M V$. We have that $A^T A$ is $n \times n$. Thus,

$$\begin{aligned}&\text{trace}(V^T A^T A V) \\ &= \text{trace}(A^T A) = \|A\|_F.\end{aligned}$$

c) Conclude that $\|A\|_F = \sqrt{\sum_{i=1}^n \sigma_i^2}$, where $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$ are the singular values of A.

We have that $\|A\|_F^2 = \sum_{ij} |a_{ij}|^2$, which is the sum of the eigenvalues of $A^T A$. Now, since the eigenvalues of $A^T A$ are the squares of the singular values of A ($|a_{ij}|^2 = \sigma_i^2$), we have $\|A\|_F^2 = \sum_i^n \sigma_i^2$,

where n is the number of columns of A. Hence,

$$\sqrt{\|A\|_F^2} = \sqrt{\sum_i^n \sigma_i^2} = \|A\|_F .$$

2) Work this problem out by hand. Let

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} .$$

(carry out the power method with starting vector $q_0 = [a, b]^T$, where $a, b \geq 0$ and $a > b$. Calculate the first few terms; what do you notice?)

$$q_0 = \begin{bmatrix} a \\ b \end{bmatrix}, \quad \|q_0\|_\infty = a$$

$$q_1 = \frac{Aq_0}{\|Aq_0\|_\infty} = \frac{1}{a} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \frac{1}{a} \begin{bmatrix} b \\ a \end{bmatrix} = \begin{bmatrix} b/a \\ 1 \end{bmatrix}$$

$$q_2 = \frac{Aq_1}{\|Aq_1\|_\infty} = \frac{1}{1} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} b/a \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ b/a \end{bmatrix}$$

$$q_3 = \frac{Aq_2}{\|Aq_2\|_\infty} = \frac{1}{1} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ b/a \end{bmatrix} = \begin{bmatrix} b/a \\ 1 \end{bmatrix}$$

$$q_4 = \begin{bmatrix} 1 \\ b/a \end{bmatrix}, q_5 = \begin{bmatrix} b/a \\ 1 \end{bmatrix}, \dots$$

I notice that after the first term q_0 , the terms oscillate back and forth such that $q_i = \begin{bmatrix} b/a \\ 1 \end{bmatrix}$ where i is odd and $q_j = \begin{bmatrix} 1 \\ b/a \end{bmatrix}$ where j is even. The sequence fails to converge.

Explain why the sequence fails to converge. What is the problem with the convergence argument we had in lecture? Which of the conditions is violated?

We have that

$$\det(A - \lambda I) = \begin{vmatrix} 0-\lambda & 1 \\ 1 & 0-\lambda \end{vmatrix} = \lambda^2 - 1 = 0$$

$$\Rightarrow (\lambda - 1)(\lambda + 1) = 0 \Rightarrow \lambda_1 = 1, \lambda_2 = -1$$

$$\Rightarrow |\lambda_1| = |\lambda_2|.$$

The sequence fails to converge because the matrix A does not have a dominant eigenvalue.

3) Let A be an $n \times n$ complex matrix with eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$ ordered in decreasing order of magnitude, that is, $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|$. Let v_1, \dots, v_n be the corresponding eigenvectors.

Assume that $\lambda_1 = 5, \lambda_2 = 3, \lambda_3 = 2$, and consider performing the power method on A starting with a vector q which does not depend on the eigenvector v_1 :

$$q = v_2 + \sum_{i=3}^n c_i v_i, \text{ for some } c_i \in \mathbb{C}, 3 \leq i \leq n.$$

If we start the power method on A with the vector q , will the method converge? Why or why not? Explain.

Let us analyze the convergence of the power method on A with the given q . We have $q = v_2 + \sum_{i=3}^n c_i v_i$, and

$$Av_i = \lambda_i v_i. \text{ Hence,}$$

$$\begin{aligned} Aq &= A(v_2 + \sum_{i=3}^n c_i v_i) \\ &= Av_2 + A \sum_{i=3}^n c_i v_i \\ &= \lambda_2 v_2 + \sum_{i=3}^n c_i \lambda_i v_i \end{aligned}$$

$$\begin{aligned} A^2 q &= A^2(v_2 + \sum_{i=3}^n c_i v_i) \\ &= A^2 v_2 + A^2 \sum_{i=3}^n c_i v_i \\ &= \lambda_2^2 v_2 + \sum_{i=3}^n c_i \lambda_i^2 v_i \end{aligned}$$

$$\begin{aligned} &\vdots \\ &\vdots \\ &\vdots \end{aligned}$$

$$\begin{aligned}
 A^j q &= \lambda_2^j v_2 + \sum_{i=3}^n c_i \lambda_i^j v_i \\
 &= \lambda_2^j v_2 + c_3 \lambda_3^j v_3 + \dots + c_n \lambda_n^j v_n \\
 &= \lambda_2^j \left(v_2 + c_3 \frac{\lambda_3^j}{\lambda_2^j} v_3 + \dots + c_n \frac{\lambda_n^j}{\lambda_2^j} v_n \right)
 \end{aligned}$$

As j tends to ∞ , tend to 0.

$$\begin{aligned}
 \Rightarrow A^j q &= \lambda_2^j v_2 \\
 &= 3^j v_2 .
 \end{aligned}$$

Yes, the power method on A with vector q will converge, just not to the first eigenvector v_1 , but to the second, v_2 . This is a consequence of the fact that the initial q does not have a component along v_1 , but has components along v_2, v_3, \dots, v_n . Since $|\lambda_2| > |\lambda_i|$ for $i = 3, \dots, n$, the power method with scaling gives us a vector v aligned with v_2 with corresponding eigenvalue $\lambda_2 = 3$, and $\|v\|_\infty = 1$.

4) Let $A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$, the 2×2 identity matrix.

(a) Calculate the characteristic polynomial of A and show it has one eigenvalue $\lambda = 1$ with algebraic multiplicity 2.

$$\det(A - \lambda I) = \det\left(\begin{bmatrix} 1-\lambda & 0 \\ 0 & 1-\lambda \end{bmatrix}\right) = (1-\lambda)(1-\lambda) = 0$$

$\Rightarrow \lambda = 1$ is a root of $\det(A - \lambda I) = 0$ twice, hence A has one eigenvalue with algebraic multiplicity 2.

b) we now perturb one coefficient of the characteristic polynomial slightly and consider the equation

$$\lambda^2 - 2\lambda + (1-\varepsilon) = 0,$$

where $0 < \varepsilon \ll 1$. Solve the equation for the roots $\hat{\lambda}_1$ and $\hat{\lambda}_2$ in terms of ε .

We use the quadratic equation:

$$\lambda = \frac{2 \pm \sqrt{4 - 4(1-\varepsilon)}}{2} = \frac{2 \pm \sqrt{4\varepsilon}}{2}$$

$$\Rightarrow \hat{\lambda}_1 = \frac{2 + \sqrt{4\varepsilon}}{2} = \frac{2 + 2\sqrt{\varepsilon}}{2} = 1 + \sqrt{\varepsilon}$$

$$\hat{\lambda}_2 = \frac{2 - \sqrt{4\varepsilon}}{2} = \frac{2 - 2\sqrt{\varepsilon}}{2} = 1 - \sqrt{\varepsilon}$$

c) Show that if $\varepsilon = 10^{-12}$, $|\hat{\lambda}_1 - \lambda|$ and $|\hat{\lambda}_2 - \lambda|$ are one million times bigger than ε .

Let $\varepsilon = 10^{-12}$. Then

$$|\hat{\lambda}_1 - \lambda| = |1 + \sqrt{\epsilon} - 1| = \sqrt{\epsilon} = 10^{-6} = 10^6 \epsilon$$

$$|\hat{\lambda}_2 - \lambda| = |1 - \sqrt{\epsilon} - 1| = \sqrt{\epsilon} = 10^{-6} = 10^6 \epsilon$$

d) Conclude that, given a relative change $\epsilon \ll 1$ in one of the coefficients of the polynomial, the relative change in the eigenvalues can grow arbitrarily large; therefore, the computation is numerically unstable.

In part (b), we observed that the relative change in one of the coefficients by ϵ changed $\lambda = 1$ to $\hat{\lambda}_1 = 1 + \sqrt{\epsilon}$ $\hat{\lambda}_2 = 1 - \sqrt{\epsilon}$. Then, in part (c), we saw how this change in the coefficient drastically changed the difference between our initial and final eigenvalues, a change from $\lambda = 1$ to $\hat{\lambda}_1 = 1 + 10^{-6}$ and $\hat{\lambda}_2 = 1 - 10^{-6}$. Not only did we change the numerical value of the initial eigenvalue, but we changed the overall number of eigenvalues and their algebraic multiplicity. This can lead to major complications if the eigenvalues have a high multiplicity, or if their values are very close to one another, since the relative change between them will grow immensely. Hence, the computation is quite numerically unstable.

5) Let $A = \begin{bmatrix} 2 & 2 & 2 \\ 4 & 0 & 2 \\ 2 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix}$.

- Use the MATLAB function rank to figure the rank of A.
- Use MATLAB's pinv to compute the pseudo inverse of A.
- Use MATLAB's svd command to find the singular values of AA^+ and A^+A , and calculate AA^+ and A^+A .

Command Window

```
>> A = [2 2 2; 4 0 2; 2 0 1; 0 1 1]

A =

2      2      2
4      0      2
2      0      1
0      1      1

>> rank(A)

ans = (a)
3

>> A_pinv = pinv(A)

A_pinv =

0.5000    0.0000    0.0000   -1.0000
1.0000   -0.4000   -0.2000   -1.0000
-1.0000    0.4000    0.2000    2.0000 (b)

>> AA_pinv = A*A_pinv

AA_pinv =

1.0000    0.0000    0.0000   -0.0000
-0.0000    0.8000    0.4000      0
-0.0000    0.4000    0.2000      0
 0.0000   -0.0000   -0.0000    1.0000

>> [U,S,V] = svd(AA_pinv)
```

U =

-1.0000	-0.0000	0.0000	0
0.0000	-0.8698	-0.2084	-0.4472
0.0000	-0.4349	-0.1042	0.8944
-0.0000	0.2330	-0.9725	0

S =

1.0000	0	0	0
0	1.0000	0	0
0	0	1.0000	0
0	0	0	0

V =

-1.0000	0	0	0
0	-0.8698	-0.2084	0.4472
0	-0.4349	-0.1042	-0.8944
0	0.2330	-0.9725	-0.0000

>> A_pinvA = A_pinv*A

A_pinvA =

1.0000	-0.0000	0.0000
0.0000	1.0000	0
-0.0000	0	1.0000

>> [U,S,V] = svd(A_pinvA)

U =

-1.0000	0.0000	-0.0000
-0.0000	-0.6919	0.7220
0.0000	0.7220	0.6919

S =

1.0000	0	0
0	1.0000	0
0	0	1.0000

V =

-1.0000	0	0
0	-0.6919	0.7220
0	0.7220	0.6919

(c)

HW 4

1) Let $A \in \mathbb{R}^{n \times m}$, $n = m$, $\text{rank}(A) = m$. Compute the SVD of the matrix $(A^T A)^{-1} A^T$ in terms of the SVD of $A = U \Sigma V^T$, and explain what the parts are.

a) Based on the definitions in class, conclude that the pseudoinverse of A has the format $A^+ = (A^T A)^{-1} A^T$.

We have $A = U \Sigma V^T$, where A is $n \times m$, U is orthogonal and $n \times n$, Σ is diagonal and $n \times m$, and V is orthogonal and $m \times m$. $\text{Rank}(A) = m$, so we have m nonzero singular values in Σ .

$$\begin{aligned} \text{Now, } A^+ &= (U \Sigma V^T)^T = (V^T)^T \Sigma^T U^T \\ &= V \Sigma^T U^T \end{aligned}$$

$$\begin{aligned} A^T A &= (V \Sigma^T U^T)(U \Sigma V^T) & (U^T U = I) \\ &= V \Sigma^T \Sigma V^T \end{aligned}$$

We know that $A^T = \begin{bmatrix} & \\ & \\ & \end{bmatrix}_{m \times n}$, and $\text{rank}(A) = m$. Hence, $A^T A$ is $m \times m$ with full rank. Thus, $(A^T A)^{-1}$ exists.

$$\begin{aligned} (A^T A)^{-1} &= (V \Sigma^T \Sigma V^T)^{-1} \\ &= (V^T)^{-1} (\Sigma^T \Sigma)^{-1} (V)^{-1} & (V \text{ is orthogonal}; \\ && V^T = V^{-1}) \\ &= V (\Sigma^T \Sigma)^{-1} V^T \end{aligned}$$

$\left[\sigma_1 \right]$

Since $n \geq m$, ε looks like $\varepsilon = \begin{bmatrix} \cdot & \sigma_m \\ \vdots & 0 \\ 0 \end{bmatrix}$ (more rows than cols).

$$\therefore \varepsilon^T \varepsilon = \begin{bmatrix} \sigma_1 & & \\ \vdots & \ddots & \\ \sigma_m & & \end{bmatrix}_{m \times n} \begin{bmatrix} \sigma_1 & & \\ \vdots & \ddots & \\ \sigma_m & & \\ 0 & & \end{bmatrix}_{n \times m} = \begin{bmatrix} \sigma_1^2 & & \\ \vdots & \ddots & \\ \sigma_m^2 & & \\ 0 & & \end{bmatrix}_{m \times m} \quad \left((\varepsilon^T \varepsilon)^{-1} \text{ exists, } \text{since } \varepsilon^T \varepsilon \text{ has full rank} \right)$$

$$\therefore (\varepsilon^T \varepsilon)^{-1} = \begin{bmatrix} \sigma_1^{-2} & & \\ \vdots & \ddots & \\ \sigma_m^{-2} & & \end{bmatrix}^{-1} = \begin{bmatrix} \sigma_1^{-2} & & \\ \vdots & \ddots & \\ \sigma_m^{-2} & & \end{bmatrix}$$

$$\begin{aligned} (A^T A)^{-1} A^T &= V (\varepsilon^T \varepsilon)^{-1} V^T (V \varepsilon^T U^T) \quad (V^T V = I) \\ &= V (\varepsilon^T \varepsilon)^{-1} \varepsilon^T U^T \end{aligned}$$

$$(\varepsilon^T \varepsilon)^{-1} \varepsilon^T = \begin{bmatrix} \sigma_1^{-2} & & \\ \vdots & \ddots & \\ \sigma_m^{-2} & & \end{bmatrix}_{m \times m} \begin{bmatrix} \sigma_1 & & \\ \vdots & \ddots & \\ \sigma_m & & \\ 0 & & \end{bmatrix}_{n \times n} = \begin{bmatrix} \sigma_1^{-1} & & \\ \vdots & \ddots & \\ \sigma_m^{-1} & & \\ 0 & & \end{bmatrix}_{m \times n} = \varepsilon^+$$

Hence, the SVD of $(A^T A)^{-1} A^T$ is $V \varepsilon^+ U^T$, where ε^+ is the pseudo inverse of ε :

$$\varepsilon = \begin{bmatrix} \sigma_1 & & \\ \vdots & \ddots & \\ \sigma_m & & \\ 0 & & \end{bmatrix}_{n \times m} \quad \varepsilon^+ = \begin{bmatrix} \sigma_1^{-1} & & \\ \vdots & \ddots & \\ \sigma_m^{-1} & & \\ 0 & & \end{bmatrix}_{m \times n}$$

$$\Sigma \Sigma^+ = \begin{bmatrix} I_m & 0_{m \times (n-m)} \\ - & 0_{(n-m) \times m} & 0_{(n-m) \times (n-m)} \end{bmatrix}$$

Now, using the fact that $\Sigma^+ = (\Sigma^\top \Sigma)^{-1} \Sigma^\top$, and $(A^\top A)^{-1}$ exists, we can conclude that $A^+ = (A^\top A)^{-1} A^\top = V \Sigma^+ V^\top$.

b) Use (a) to calculate $A^+ A$.

$$\begin{aligned} A^+ A &= (V \Sigma^+ V^\top)(U \Sigma V^\top) && (V^\top V = I_n) \\ &= V \underbrace{\Sigma^+ \Sigma}_{\begin{bmatrix} I_m & \\ - & 0 \end{bmatrix}} V^\top \\ &= V V^\top = I_m \end{aligned}$$

c) Use the SVDs of A and A^+ to calculate $A A^+$, simplifying as much as possible

$$\begin{aligned} A A^+ &= (V \Sigma V^\top)(V \Sigma^+ V^\top) \\ &= U \underbrace{\Sigma \Sigma^+}_{\begin{bmatrix} I_m & \\ - & 0 \end{bmatrix}} U^\top \\ &= U U^\top \end{aligned}$$

$$2) \text{ Let } A = U \cdot \begin{bmatrix} 3 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{bmatrix} \cdot V, \quad B = P \cdot \begin{bmatrix} 4 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \cdot Q,$$

where $U \in \mathbb{R}^{3 \times 3}$, $V \in \mathbb{R}^{3 \times 3}$, $P \in \mathbb{R}^{3 \times 3}$, $Q \in \mathbb{R}^{4 \times 4}$ are all given orthogonal matrices. Compute $\|A\|_2$, $\|A^{-1}\|_2$, $\kappa_2(A)$, $\|B\|_2$ and the pseudo inverse B^+ (the answer will be in terms of P and Q).

$$\|A\|_2 = \text{largest singular value of } A = \sigma_1$$

$$\|A^{-1}\|_2 = \text{largest singular value of } A^{-1}$$

$$A^{-1} = V \Sigma^{-1} U^T, \quad \Sigma^{-1} = \begin{bmatrix} \sigma_1^{-1} & & \\ & \ddots & \\ & & \sigma_n^{-1} \end{bmatrix}$$

$$\kappa_2(A) = \frac{\sigma_1}{\sigma_n}$$

$$\Rightarrow A = U \begin{bmatrix} 3 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{bmatrix} V, \quad A^{-1} = V \begin{bmatrix} \frac{1}{3} & 0 & 0 \\ 0 & \frac{1}{2} & 0 \\ 0 & 0 & 1 \end{bmatrix} U$$

$$\|A\|_2 = 3$$

$$\Rightarrow \|A^{-1}\|_2 = \frac{1}{3} = 1$$

$$\kappa_2(A) = \frac{3}{1} = 3$$

$$\Rightarrow B = P \begin{bmatrix} 4 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} Q, \quad \|B\|_2 = 4$$

$$B^+ = Q \left((\varepsilon^\top \varepsilon)^{-1} \varepsilon^\top \right) P$$

$$\varepsilon^\top = \begin{bmatrix} 4 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad \varepsilon^\top \varepsilon = \begin{bmatrix} 4 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 4 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 16 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$$(\varepsilon^\top \varepsilon)^{-1} = \begin{bmatrix} 1/16 & 0 & 0 & 0 \\ 0 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad (\varepsilon^\top \varepsilon)^{-1} \varepsilon^\top = \begin{bmatrix} 1/16 & 0 & 0 & 0 \\ 0 & 1/4 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 4 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 1/4 & 0 & 0 \\ 0 & 1/2 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$B^+ = Q \begin{bmatrix} 1/4 & 0 & 0 \\ 0 & 1/2 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} P$$

$4 \times 4 \quad 4 \times 3$

3) Let $A \in \mathbb{R}^{n \times m}$, $n \geq m$.

a) Use the SVD of A to deduce the SVD of $A^T A$.

$$A = \begin{bmatrix} & \\ n & \\ & m \end{bmatrix} = U \Sigma V^T = \begin{bmatrix} & \\ n \times n & \\ U & \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r & 0 \\ & & & 0 \end{bmatrix} \begin{bmatrix} & \\ & \\ & m \times m \\ & V^T \end{bmatrix}$$

$$A^T = (U \Sigma V^T)^T = (V^T)^T \Sigma^T U^T = V \Sigma^T U^T$$

$$\begin{aligned} A^T A &= (V \Sigma^T U^T)(U \Sigma V^T) \\ &= V \Sigma^T \Sigma V^T \end{aligned}$$

b) If $m=n$ and A is full rank, use (a) to show that $\|A^T A\|_2 = \|A\|_2^2$ and that $K_2(A^T A) = K_2(A)^2$.

$$A = \begin{bmatrix} & \\ \text{rank}(A)=m & \\ n & \\ & m \end{bmatrix} = U \Sigma V^T = \begin{bmatrix} & \\ n \times n & \\ U & \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_m \end{bmatrix} \begin{bmatrix} & \\ & \\ & m \times m \\ & V \end{bmatrix}$$

(square)

We know that $\|A\|_2 = \text{largest singular value of } A = \sigma_1$, so $\|A\|_2^2 = \sigma_1^2$. Thus, $\|A^T A\|_2 = \text{largest singular value of } A^T A$.

We previously found that $A^T A = V \Sigma^T \Sigma V^T$, where

$$\Sigma^T \Sigma = \Sigma^2 = \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_m^2 \end{bmatrix}, \text{ since } \Sigma \text{ is diagonal. Hence,}$$

$$\left[\begin{array}{c} \vdots \\ \sigma_m^2 \end{array} \right]$$

$$\|A^T A\|_2 = \sigma_1^2 = \|A_2\|^2.$$

We also know that $K_2(A) = \frac{\sigma_1}{\sigma_m}$, and so $K_2(A)^2 = \frac{\sigma_1^2}{\sigma_m^2}$.

We have that $K_2(A^T A) = \frac{\sigma_1^2}{\sigma_m^2} = K_2(A)^2$.

4) Work this exercise using pencil and paper. You can use MATLAB to check your work. Let A be the following exterior (or outer) product of two vectors:

$$A = \begin{bmatrix} 1 \\ 2 \\ -1 \end{bmatrix} \cdot [1 \ 1]$$

Note that A is 3×2 . Answer the questions below; you can do so without ever forming A explicitly.

a) What is the rank r of A ?

Since we are multiplying a 3×1 vector with a 1×2 vector, the result has dimension 3×2 . Since the 1×2 vector is $[1 \ 1]$, the result will be a matrix of 2 columns with the same 3 entries in each column. Thus, the result is a matrix with linearly dependent columns. Hence, $\text{rank}(A) = 1$.

b) Think about the "sum of rank one matrices" expression for the SVD of A , then consider the reduced SVD of A :

$A = U_r \Sigma_r V_r^T$. What are the sizes of U_r , Σ_r , and V_r ?

We have $\begin{bmatrix} A \\ 3 \times 2 \end{bmatrix} = \begin{bmatrix} U \\ 3 \times 3 \end{bmatrix} \begin{bmatrix} \Sigma \\ 3 \times 2 \end{bmatrix} \begin{bmatrix} V^T \\ 2 \times 2 \end{bmatrix}$,

$A = \sum_{i=1}^r \sigma_i u_i v_i^T = \text{sum of rank 1 matrices}$

Since $\text{rank}(A) = r = 1$, we have

$$A = \sigma_1 U_1 V_1^T \Rightarrow A = U_1 \Sigma_1 V_1^T$$

$\Rightarrow \Sigma_1$ has one singular value, so Σ_1 is 1×1 . Hence,
 U_1 must be 3×1 , and V_1^T must be 1×2 .

c) Use the fact that the columns of U_r and V_r are orthonormal to figure out U_r , V_r , and Σ_r .

We know that U_r is the left singular vector of A , meaning it is the normalized column space of A . We know that $\text{rank}(A) = 1$, so $\text{col}(A) = \left\{ \begin{bmatrix} 1 \\ 2 \\ -1 \end{bmatrix} \right\}$. Hence,

$$U_r = \frac{1}{\| \begin{bmatrix} 1 \\ 2 \\ -1 \end{bmatrix} \|} \begin{bmatrix} 1 \\ 2 \\ -1 \end{bmatrix} = \frac{1}{\sqrt{6}} \begin{bmatrix} 1 \\ 2 \\ -1 \end{bmatrix}.$$

Similarly, V is the right singular vector, so it is the normalized row space of A , where $\text{row}(A) = \left\{ \begin{bmatrix} 1 & 1 \end{bmatrix} \right\}$. Thus,

$$V_r = \frac{1}{\| \begin{bmatrix} 1 & 1 \end{bmatrix} \|} \begin{bmatrix} 1 & 1 \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}$$

Now, to find Σ_r , we have

$$A^T A = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & -1 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 2 \\ -1 \end{bmatrix} \begin{bmatrix} 1 & 1 \end{bmatrix} = \begin{bmatrix} 6 & 6 \\ 6 & 6 \end{bmatrix}$$

$$\sigma_1 = \sqrt{6} \quad \Rightarrow \quad \Sigma_1 = [\sqrt{6}]$$

5) Run the attached low-rank-approximation.m MATLAB code.

a) Explain line by line what the code does (you might need to google some of the commands).

```
% load image
```

```
A = imread('street2.jpg'); - Reads image specified by 'street2.jpg'  
A = rgb2gray(A); - A is an RGB image, so this converts the color image to  
B = double(A); - Converts the type of the entries in A from int to double,  
stores the matrix w/ double type entries in B
```

```
% compute SVD
```

```
size(B) - gets size of B  
r = rank(B) - gets rank of B and stores the number in variable r  
[U,S,V] = svd(B); - computes the SVD of B, such that  $B = U \cdot S \cdot V'$ 
```

```
% approximate image
```

```
ranks = [1 2 4 8 16 32 64 r]; creates a row vector of these specified  
values, stored in variable "ranks"
```

```
l = length(ranks); - computes the length (number of elements) of ranks
```

```
| for i = 1:l - for loop, iterates l times (# of elements in "ranks")
```

```
% compute rank i approximation
```

```
k = ranks(i); - stores the ith element of "ranks" in variable K
```

```
approxB = U(:,1:k)*S(1:k,1:k)*V(:,1:k)'; - corresponds to  $\sum_k U_k V_k$ ,  
meaning for each entry k in ranks, we are computing the rank-k approximation of B.
```

```
approxA = uint8(approxB); - converts (approxB) into an unsigned 8-bit int
```

```
% plot images
```

```
figure(1) - gets figure with Number property = 1  
subplot(2,4,i) - divides current figure into a 2x4 grid, creates axes for each iteration  
imshow(approxA); - displays grayscale image of approx(A) according to figure  
title(sprintf('rank %d approximation',k)) - titles each subplot using  
the print function, %d corresponds to the value of K during the particular iteration.  
end end loop
```

b) Explain what the algorithm, as a whole, does.

The algorithm takes an image, converts it to grayscale, and then creates several approximations of the image with rank $1, 2, \dots, 64$, and finally the actual rank. Essentially, the algorithm uses Singular Value Decomposition to compress/expand an image by using a certain number of singular values.

c) Note that the approximation gets better as we increase K . Even when $K = 32$, the resulting approximation looks reasonable. What is the advantage to use/store the $K=32$ approximation instead of the original image? What is the disadvantage?

We know that the matrices U_{32} , Σ_{32} , and V_{32}^T are much smaller than U_{480} , Σ_{480} , V_{480}^T . Hence, storing the $K=32$ approx. uses much less space than the $K=480$ approx. Also, when using the image, such as for a newspaper, since the $K=32$ approx. has fewer singular values, processing the image is much faster than the original.

The disadvantage is that since we are approximating using a much smaller rank, and although the approximation is reasonable, it will never have as many fine details as the original.

HW 3

SABELLA SAVLINO

i) Let $w_1 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$, $w_2 = \begin{bmatrix} 3 \\ 3 \\ -1 \\ -1 \end{bmatrix}$, $w_3 = \begin{bmatrix} 6 \\ 0 \\ 2 \\ 0 \end{bmatrix}$ be vectors in \mathbb{R}^4 .

Apply the classical GS process to find an orthonormal basis for the subspace spanned by w_1, w_2, w_3 .

$$i) q_1 = \frac{w_1}{\|w_1\|_2} = \frac{1}{2} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \end{bmatrix} \quad r_{11} = 2$$

$$ii) \tilde{q}_2 = w_2 - \langle q_1, w_2 \rangle q_1$$

$$= \begin{bmatrix} 3 \\ 3 \\ -1 \\ -1 \end{bmatrix} - \left((\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2}) \cdot (3, 3, -1, -1) \right) \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \end{bmatrix}$$

$$= \begin{bmatrix} 3 \\ 3 \\ -1 \\ -1 \end{bmatrix} - (2) \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \end{bmatrix} \quad r_{12} = 2$$

$$= \begin{bmatrix} 3 \\ 3 \\ -1 \\ -1 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \\ -2 \\ -2 \end{bmatrix}$$

$$q_2 = \frac{\tilde{q}_2}{\|\tilde{q}_2\|_2} = \frac{1}{4} \begin{bmatrix} 2 \\ 2 \\ -2 \\ -2 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ -\frac{1}{2} \\ -\frac{1}{2} \end{bmatrix} \quad r_{22} = 4$$

$$iii) \tilde{q}_3 = w_3 - \langle q_1, w_3 \rangle q_1 - \langle q_2, w_3 \rangle q_2$$

$$= \begin{bmatrix} 6 \\ 0 \\ 2 \\ 0 \end{bmatrix} - \left((\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2}) \cdot (6, 0, 2, 0) \right) \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \end{bmatrix}$$

$$= \begin{bmatrix} 6 \\ 0 \\ 2 \\ 0 \end{bmatrix} - (4) \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \end{bmatrix} \quad r_{13} = 4$$

$$= \begin{bmatrix} 6 \\ 0 \\ 2 \\ 0 \end{bmatrix} - \begin{bmatrix} 2 \\ 2 \\ 2 \\ 2 \end{bmatrix} = \begin{bmatrix} 4 \\ -2 \\ 0 \\ -2 \end{bmatrix}$$

$$\begin{bmatrix} 4 \\ -2 \\ 0 \\ -2 \end{bmatrix} - \left((\frac{1}{2}, \frac{1}{2}, -\frac{1}{2}, -\frac{1}{2}) \cdot (6, 0, 2, 0) \right) \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ -\frac{1}{2} \\ -\frac{1}{2} \end{bmatrix}$$

$$= \begin{bmatrix} 4 \\ -2 \\ 0 \\ -2 \end{bmatrix} - (2) \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ -\frac{1}{2} \\ -\frac{1}{2} \end{bmatrix} = \begin{bmatrix} 3 \\ -3 \\ 1 \\ -1 \end{bmatrix} \quad r_{23} = 2$$

$$q_3 = \frac{\tilde{q}_3}{\|q_3\|_2} = \frac{1}{\sqrt{20}} \begin{bmatrix} 3 \\ -3 \\ 1 \\ -1 \end{bmatrix} \quad r_{33} = \sqrt{20}$$

Orthonormal basis :

$$\left\{ \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \end{bmatrix}, \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ -\frac{1}{2} \\ -\frac{1}{2} \end{bmatrix}, \frac{1}{\sqrt{20}} \begin{bmatrix} 3 \\ -3 \\ 1 \\ -1 \end{bmatrix} \right\}$$

2) For $A = [w_1, w_2]$ with w_1, w_2 as above, find the minimizer to the least squares problem $\min_{x \in \mathbb{R}^2} \|b - Ax\|_2$ for $b = [1, 1, 0, 1]^T$, using the QR method discussed in class.

$$A = \begin{bmatrix} 1 & 3 \\ 1 & 3 \\ 1 & -1 \\ 1 & -1 \end{bmatrix}_{w_1 \quad w_2} \quad Q = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix}_{q_1 \quad q_2} \quad R = \begin{bmatrix} r_{11} & r_{12} \\ 0 & r_{22} \end{bmatrix} \quad b = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 1 \end{bmatrix}$$

$$\min_{x \in \mathbb{R}^2} \|b - Ax\| \Rightarrow x = R^{-1} Q^T b$$

$$R^{-1} = \frac{1}{8} \begin{bmatrix} 4 & -2 \\ 0 & 2 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & -\frac{1}{4} \\ 0 & \frac{1}{4} \end{bmatrix} \quad Q^T = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \end{bmatrix}$$

$$R^{-1} Q^T b = \begin{bmatrix} \frac{1}{2} & -\frac{1}{4} \\ 0 & \frac{1}{4} \end{bmatrix}_{2 \times 2} \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \end{bmatrix}_{2 \times 4} \begin{bmatrix} 1 \\ 1 \\ 0 \\ 1 \end{bmatrix}$$

$$= \begin{bmatrix} \frac{1}{8} & \frac{1}{8} & \frac{3}{8} & \frac{3}{8} \\ \frac{1}{8} & \frac{1}{8} & -\frac{1}{8} & -\frac{1}{8} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 0 \\ 1 \end{bmatrix}$$

$$= \begin{bmatrix} \frac{5}{8} \\ \frac{1}{8} \end{bmatrix} = x$$

3) Let Q be an orthogonal $n \times n$ matrix ($QQ^T = Q^TQ = I$)

and consider the induced 2-norm on $n \times n$ matrices A ,

$$\|A\|_2 = \max_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2}$$

with $\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$: the usual 2-norm on vectors of length n

(a) calculate $\|Q\|_2$ and $\kappa_2(Q) = \kappa_{\|\cdot\|_2}(Q)$

Multiplication by orthogonal matrices preserves 2-norm:

$$\begin{aligned}\|Qx\|_2 &= \langle Qx, Qx \rangle \\ &= (Qx)^T (Qx) \\ &= x^T \underbrace{Q^T Q}_I x \\ &= \langle x, x \rangle = \|x\|_2.\end{aligned}$$

$$\text{Hence, } \|Q\|_2 = \frac{\|Qx\|_2}{\|x\|_2} = \frac{\|x\|_2}{\|x\|_2} = 1.$$

$$\text{Now, } \kappa_2(Q) = \kappa_{\|\cdot\|_2}(Q) = \|Q\|_2 \|Q^{-1}\|_2.$$

$$\text{We have that } QQ^T = Q^TQ = I \Rightarrow Q = Q^{-1}.$$

$$\text{Hence, } \|Q\|_2 = \|Q^{-1}\|_2 = 1.$$

$$\text{Therefore, } \kappa_2(Q) = 1.$$

(b) Consider the perturbed system $(Q + \delta Q)\hat{x} = b$,

with $\hat{x} = x + \delta x$. What can we say about the relative

error $\|\delta x\|_2 / \|x\|_2$ of the solution, in terms of the induced 2-norm $\|\delta Q\|_2$ of the perturbation matrix δQ ?

By a previous Theorem, we have

$$\frac{\|\delta x\|_2}{\|\hat{x}\|_2} \leq K_2(Q) \frac{\|\delta Q\|_2}{\|Q\|_2}.$$

We found in part (a) that $K_2(Q) = 1$. Hence,

$$\frac{\|\delta x\|_2}{\|x\|_2} \leq \frac{\|\delta Q\|_2}{\|Q\|_2}.$$

Also, we find that $\|Q\|_2 = 1$. Thus,

$$\frac{\|\delta x\|_2}{\|x\|_2} \leq \|\delta Q\|_2.$$

We can say that the relative error $\frac{\|\delta x\|_2}{\|x\|_2}$ of the solution is bounded by the relative perturbation $\|\delta Q\|_2$ of the matrix Q .

- 4) Let M be an $n \times m$ matrix, $n \geq m$, and A be an $m \times m$ matrix such that $A = M^T M$.
 Also, suppose M is full-rank.
 If $M = QR$ is the reduced QR factorization of M , write

down the relationship b/w R and A. Specifically, given the properties of M and R, what is R to A. Justify your answer.

First, let's draw a picture.

$$n \begin{bmatrix} M \\ m \end{bmatrix} = n \begin{bmatrix} Q \\ m \end{bmatrix}^m \begin{bmatrix} R \\ m \end{bmatrix} \quad (n \geq m)$$

$\Rightarrow \text{rank } M = m$

$$m \begin{bmatrix} M^T \\ n \end{bmatrix} n \begin{bmatrix} M \\ m \end{bmatrix} = m \begin{bmatrix} A \\ m \end{bmatrix}$$

We know by the Gram-Schmidt process that Q consists of m orthonormal columns and R is an upper triangular matrix. Hence, Q is orthogonal, so $Q^T Q = I_m$. Since M is full rank and R is $m \times m$, R is invertible. Also, since M is full rank, A is symmetric and positive definite.

We have $A = M^T M$, $M = QR$

$$\Rightarrow A = (QR)^T (QR) \quad ((AB)^T = B^T A^T)$$

$$\Rightarrow A = R^T \underbrace{Q^T Q}_{{I_m}} R$$

$$\Rightarrow A = R^T R$$

Hence, $A = R^T R$.

Now, we have that R is upper triangular and invertible and A is symmetric and positive definite. Therefore, R is the Cholesky factor of A .

5) Use the MATLAB code for classical and modified Gram-Schmidt Process

Command Window

```
>> A = hilb(8);
>> [Q1, R1] = gs(A);
Unrecognized function or variable 'gs'.
```

Did you mean:

```
>> [Q1, R1] = classicalgs(A);
```

$Q =$

Columns 1 through 6

0.8091	-0.5467	0.2064	-0.0605	0.0144	-0.0028
0.4046	0.2825	-0.6766	0.4934	-0.2236	0.0714
0.2697	0.3736	-0.2413	-0.4274	0.6135	-0.3855
0.2023	0.3636	0.0543	-0.4519	-0.1109	0.5920
0.1618	0.3354	0.2226	-0.2124	-0.4481	0.0758
0.1349	0.3059	0.3161	0.0579	-0.3360	-0.4428
0.1156	0.2792	0.3669	0.2932	0.0292	-0.3318
0.1011	0.2559	0.3929	0.4813	0.4965	0.4289

Columns 7 through 8

0.0001	-0.0001
-0.0074	0.0068
0.0958	-0.0926
-0.3983	0.3931
0.5803	-0.5805
-0.0057	0.0097
-0.6117	0.6135
0.3482	-0.3513

$R =$

fx Columns 1 through 6

1.2359	0.7192	0.5214	0.4130	0.3435	0.2947
0	0.1499	0.1665	0.1612	0.1512	0.1407
0	0	0.0118	0.0192	0.0233	0.0254
0	0	0	0.0007	0.0015	0.0023
0	0	0	0	0.0000	0.0001
0	0	0	0	0	0.0000
0	0	0	0	0	0
0	0	0	0	0	0

Columns 7 through 8

0.2584	0.2303
0.1308	0.1220
0.0264	0.0267
0.0028	0.0032
0.0002	0.0002
0.0000	0.0000
0.0000	0.0000
0	0.0000

```
>> [Q2, R2] = modifiedgs(A);
>> checkQ = norm(Q1-Q2)
```

checkQ =

1.4138

```
>> checkR = norm(R1-R2)
```

checkR =

2.0354e-06

fx >> orthoQ1 = norm(Q1'*Q1-eye(8))

orthoQ1 =

1.0325

```
>> orthoQ2 = norm(Q2'*Q2-eye(8))
```

orthoQ2 =

4.3754e-07

>>

HW 2

1) a) verify that $\|x\|_1$ is a norm

i) $\|x\|_1 > 0$, $\forall x \neq 0$; $\|x\|_1 = 0$ when $x = 0$

$$\|x\|_1 = \sum_{i=1}^n |x_i| = |x_1| + |x_2| + \dots + |x_n| .$$

let $x = 0$. Then $\|x\|_1 = |0| = 0$.

let $x \neq 0$. Then $\|x\|_1 = |x_1| + |x_2| + \dots + |x_n|$ is always positive since we take the absolute value of each entry and add.

ii) $\forall x$, $\|cx\|_1 = |c| \cdot \|x\|_1$,

$$\|cx\|_1 = \sum_{i=1}^n |cx_i| = |cx_1| + |cx_2| + \dots + |cx_n|$$

$$\Rightarrow |c| \sum_{i=1}^n = |c| (|x_1| + |x_2| + \dots + |x_n|)$$

$$= |c| \|x\|_1 .$$

iii) $\|x+y\|_1 \leq \|x\|_1 + \|y\|_1$,

$$\|x+y\|_1 = \sum_{i=1}^n |x_i + y_i| = |x_1 + y_1| + |x_2 + y_2| + \dots + |x_n + y_n|$$

If only x_i or only y_i is negative, $|x_i + y_i| < |x_i| + |y_i|$.

If both x_i and y_i are the same sign, $|x_i + y_i| = |x_i| + |y_i|$.

Hence, $\|x+y\|_1 \leq \|x\|_1 + \|y\|_1$.

b) Verify that $\|x\|_\infty$ is a norm

i) $\forall x \neq 0, \|x\|_\infty > 0$; $\|x\|_\infty = 0$ when $x = 0$

Let $x \neq 0$. Then $\|x\|_\infty = \max_{i=1:n} \{|x_i|\}$, where the maximum $|x_i|$ is always positive. Hence, $\|x\|_\infty > 0$ when $x \neq 0$.
Let $x = 0$. Then $\|x\|_\infty = \max_{i=1:n} \{|x_i|\} = |0| = 0$.

ii) $\forall x \in \mathbb{R}^n, \|cx\|_\infty = |c| \cdot \|x\|_\infty$

$$\begin{aligned}\|cx\|_\infty &= \max_{i=1:n} \{|cx_i|\} = \max \{|cx_1|, |cx_2|, \dots, |cx_n|\} \\ &= \max \{|c|(|x_1|, |x_2|, \dots, |x_n|)\} = |c| \cdot \max \{|x_1|, \dots, |x_n|\} \\ &= |c| \cdot \max_{i=1:n} \{|x_i|\} = |c| \|x\|_\infty\end{aligned}$$

- MAX NORM : L_∞ norm
 $\|x\|_\infty = \max_{i=1:n} \{|x_1|, |x_2|, \dots, |x_n|\}$
 $= \max_{i=1:n} \{|x_i|\}$

for vector $x = [x_1, x_2, \dots, x_n]$

iii) $\|x+y\|_\infty \leq \|x\|_\infty + \|y\|_\infty$

$$\|x+y\|_\infty = \max_{i=1:n} \{|x_i+y_i|\} = \max \{|x_1+y_1| + |x_2+y_2| + \dots + |x_n+y_n|\}$$

If only x_i or only y_i is positive, $|x_i+y_i| < |x_i|+|y_i|$.

If x_i and y_i are the same sign, $|x_i+y_i| = |x_i|+|y_i|$.

Hence, for $\max_{i=1:n} \{|x_i+y_i|\}, \|x+y\|_\infty \leq \|x\|_\infty + \|y\|_\infty$.

2) Let $A = \alpha I$ be a multiple of the $n \times n$

identity matrix, with $\alpha \in \mathbb{R}, \alpha \neq 0$, and consider $\|\cdot\|$ to be an induced matrix norm. Calculate $\|A\|$, $\|A^{-1}\|$, $\det(A)$, and $K_{\|\cdot\|}(A)$.

$$A = \alpha I, \alpha \neq 0 \Rightarrow \|A\| = \|\alpha I\| = |\alpha| \|I\| \\ = |\alpha| \cdot 1 = |\alpha|$$

$$\|A^{-1}\| = \|(\alpha I)^{-1}\| = \frac{1}{|\alpha|} \|I\| = \frac{1}{|\alpha|} \cdot 1 = \frac{1}{|\alpha|}$$

$$\det(A) = \det(\alpha I) = \underbrace{\alpha \cdot \alpha \cdot \dots \cdot \alpha}_{\# \text{cols/rows}} = \alpha^n$$

$$K_{\|\cdot\|}(A) = \|A\| \|A^{-1}\| = |\alpha| \cdot \frac{1}{|\alpha|} = 1$$

3) Let A be a non-singular $n \times n$ matrix.

a) Show that, in any norm, $K(A) = K(A^{-1})$

$$K(A) = \|A\| \|A^{-1}\|$$

$$K(A^{-1}) = \|A^{-1}\| \|(A^{-1})^{-1}\| \\ = \|A^{-1}\| \|A\| = \|A\| \|A^{-1}\| = K(A)$$

b) By rewriting $Ax = b$ as $A^{-1}b = x$, use the proof template we did in class to show the so-called companion inequality:

$$\frac{\|\delta b\|}{\|b\|} \leq K(A) \frac{\|\delta x\|}{\|x\|}$$

Rewrite $Ax = b$ as $A^{-1}b = x$.

Let us perturb x such that b will also perturb

$$A^{-1}(b + \delta b) = (x + \delta x)$$

$$\cancel{A^{-1}b + A^{-1}\delta b = x + \delta x} \quad (A^{-1}b = x)$$

$$A^{-1} \delta b = \delta x$$

$$\text{or } \delta b = A \delta x$$

Take a norm on both sides

$$\|\delta b\| = \|A \delta x\| \left(\leq \|A\| \|\delta x\| \right) \dots (1)$$

↳ we know that

$$\|AB\| \leq \|A\|\|B\|$$

Now, consider $A^{-1}b = x$, take norm on both sides

$$\|x\| = \|A^{-1}b\| \leq \|A^{-1}\| \|b\|$$

$$\Rightarrow \cancel{\|A^{-1}\|} \|b\| \geq \|A^{-1}b\| = \|x\| \Rightarrow \|b\| \geq \frac{\|x\|}{\|A^{-1}\|}$$

$$\Rightarrow \frac{1}{\|b\|} \leq \frac{\|A^{-1}\|}{\|\lambda\|} \dots \dots \dots \dots \dots \dots \quad (2)$$

combine ① and ② :

$$\frac{\|\delta b\|}{\|b\|} \leq \frac{\|A\| \|\delta x\|}{\|b\|} \quad \left(\frac{1}{\|b\|} \leq \frac{\|A^{-1}\|}{\|x\|} \right)$$

$$\Rightarrow \frac{\|\delta b\|}{\|b\|} \leq \frac{\|A\| \|\delta x\|}{\|b\|} \leq \|A\| \|\delta x\| \frac{\|A^{-1}\|}{\|x\|}$$

$$\leq \underbrace{\|A\| \|A^{-1}\|}_{K(A)} \frac{\|\delta x\|}{\|x\|}$$

Hence,

$$\frac{\|\delta b\|}{\|b\|} \leq K(A) \frac{\|\delta x\|}{\|x\|}$$

4) Let $A = \begin{bmatrix} 3 & 0 & 2 \\ 2 & 0 & -2 \\ 0 & 1 & 1 \end{bmatrix}$ and let $b = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$.

a) Calculate the norms $\|A\|_1$, $\|A\|_\infty$, $\|A\|_F$.

$$\|A\|_1 = \max_{x \neq 0} \frac{\|Ax\|}{\|x\|_1} = \max_{i=1:n} \sum_{j=1}^n |a_{ij}| = \max\{5, 1, 5\} = 5$$

$$\|A\|_\infty = \max_{x \neq 0} \frac{\|Ax\|_\infty}{\|x\|_\infty} = \max_{i=1..n} \sum_{j=1}^n |a_{ij}| = \max\{5, 4, 2\} = 5$$

$$\|A\|_F = \sqrt{\sum_{i,j=1}^n a_{ij}} = \sqrt{3^2 + 0^2 + 2^2 + 2^2 + 0^2 + (-2)^2 + 0^2 + 1^2 + 1^2} = \sqrt{23}$$

b) Compute the condition numbers $K_1(A)$, $K_\infty(A)$, $K_F(A)$ with respect to the 1-norm, ∞ -norm, Frobenius norm respectively.

$$A^{-1} = \left[\begin{array}{ccc|ccc} 3 & 0 & 2 & 1 & 0 & 0 \\ 2 & 0 & -2 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 \end{array} \right] \xrightarrow{\frac{1}{3}R_1} \sim \left[\begin{array}{ccc|ccc} 1 & 0 & \frac{2}{3} & 1 & \frac{1}{3} & 0 & 0 \\ 2 & 0 & -2 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 & 0 \end{array} \right] R_2 - 2R_1$$

$$\sim \left[\begin{array}{ccc|ccc} 1 & 0 & \frac{2}{3} & 1 & \frac{1}{3} & 0 & 0 \\ 0 & 0 & -\frac{10}{3} & 1 & -\frac{2}{3} & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 & 0 \end{array} \right] R_2 \leftrightarrow R_3 \sim \left[\begin{array}{ccc|ccc} 1 & 0 & \frac{2}{3} & 1 & \frac{1}{3} & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & -\frac{10}{3} & 1 & -\frac{2}{3} & 1 & 0 \end{array} \right] \xrightarrow{-\frac{3}{10} \cdot R_3}$$

$$\sim \left[\begin{array}{ccc|ccc} 1 & 0 & \frac{2}{3} & 1 & \frac{1}{3} & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & \frac{1}{5} & -\frac{3}{10} & 0 & 0 \end{array} \right] R_2 - R_3 \sim \left[\begin{array}{ccc|ccc} 1 & 0 & 0 & 1 & \frac{1}{5} & \frac{1}{5} & 0 \\ 0 & 1 & 0 & 0 & -\frac{1}{5} & \frac{3}{10} & 1 \\ 0 & 0 & 1 & \frac{1}{5} & -\frac{3}{10} & 0 & 0 \end{array} \right] A^{-1} = \left[\begin{array}{ccc} \frac{1}{5} & \frac{1}{5} & 0 \\ -\frac{1}{5} & \frac{3}{10} & 1 \\ \frac{1}{5} & -\frac{3}{10} & 0 \end{array} \right]$$

$$K_1(A) = \|A\|_1, \|A^{-1}\|_1,$$

$$\|A\|_1 = \max\{\text{col sum}\} = \max\{\frac{1}{5}, \frac{1}{5}, 1\} = 1$$

$$K_1(A) = 5 \cdot 1 = 5$$

$$K_\infty(A) = \|A\|_\infty \|A^{-1}\|_\infty$$

$$\|A^{-1}\|_\infty = \max\{\text{row sum}\} = \max\{\frac{2}{5}, \frac{1}{10}, -\frac{1}{10}\} = \frac{1}{10}$$

$$K_\infty(A) = 5 \cdot \frac{1}{10} = \frac{55}{10} = \frac{11}{2} = 5.5$$

$$K_F(A) = \|A\|_F \|A^{-1}\|_F$$

$$\|A^{-1}\|_F = \sqrt{(\frac{1}{5})^2 + (\frac{1}{5})^2 + 0^2 + (-\frac{1}{5})^2 + (\frac{3}{10})^2 + 1^2 + (\frac{1}{5})^2 + (-\frac{3}{10})^2 + 0^2} = \frac{\sqrt{134}}{10}$$

$$K_F(A) = \sqrt{23} \cdot \frac{\sqrt{134}}{10} = \frac{\sqrt{3082}}{10}$$

c) Let $\delta b = \begin{bmatrix} \epsilon \\ 0 \\ 0 \end{bmatrix}$. Consider the linear systems $Ax = b$,

$A\hat{x} = b + \delta b$ and let $\delta x = \hat{x} - x$ ($x + \delta x = \hat{x}$). Estimate the relative error from above and below through condition numbers, for norms $1, \infty$, and Frobenius. (see Q2 for lower bound)

From a previous theorem, we have

$$\frac{\|\delta x\|}{\|x\|} \leq K_{\|\cdot\|}(A) \frac{\|\delta b\|}{\|b\|}$$

↓
condition no.
of A

relative in
error in b

relative error
in x

We have the system $Ax = b$,

$$\begin{bmatrix} 3 & 0 & 2 \\ 2 & 0 & -2 \\ 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} \Rightarrow \begin{array}{l} 3x_1 + 2x_3 = 0 \\ 2x_1 - 2x_3 = 1 \\ x_2 + x_3 = 1 \end{array} \Rightarrow \begin{array}{l} 3x_1 = -2x_3 \\ x_1 = -\frac{2}{3}x_3 \\ x_2 = 1 - x_3 \end{array}$$

$$2\left(-\frac{2}{3}x_3\right) - 2x_3 = 1 \Rightarrow -\frac{4}{3}x_3 - \frac{6}{3}x_3 = 1 \Rightarrow -\frac{10}{3}x_3 = 1 \Rightarrow x_3 = -\frac{3}{10}$$

$$x_1 = -\frac{2}{3}\left(-\frac{3}{10}\right) = \frac{6}{30} = \frac{1}{5}$$

$$x_2 + \left(-\frac{3}{10}\right) = 1 \Rightarrow x_2 = 1 + \frac{3}{10} = \frac{13}{10} \Rightarrow \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1/5 \\ 13/10 \\ -3/10 \end{bmatrix}$$

Now, let us perturb b such that $Ax = b \Rightarrow$

$$A\hat{x} = b + \delta b$$

$$\begin{bmatrix} 3 & 0 & 2 \\ 2 & 0 & -2 \\ 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} \hat{x}_1 \\ \hat{x}_2 \\ \hat{x}_3 \end{bmatrix} = \begin{bmatrix} \epsilon \\ 1 \\ 1 \end{bmatrix} \Rightarrow \begin{array}{l} 3\hat{x}_1 + 2\hat{x}_3 = \epsilon \\ 2\hat{x}_1 - 2\hat{x}_3 = 1 \\ \hat{x}_2 + \hat{x}_3 = 1 \end{array} \Rightarrow \begin{array}{l} 3\hat{x}_1 = \epsilon - 2\hat{x}_3 \\ \hat{x}_1 = \frac{\epsilon}{3} - \frac{2}{3}\hat{x}_3 \\ \hat{x}_2 = 1 - \hat{x}_3 \end{array}$$

$$\Rightarrow 2\left(\frac{\epsilon}{3} - \frac{2}{3}\hat{x}_3\right) - 2\hat{x}_3 = 1 \Rightarrow \frac{2}{3}\epsilon - \frac{4}{3}\hat{x}_3 - \frac{6}{3}\hat{x}_3 = 1$$

$$\Rightarrow -\frac{10}{3} \hat{x}_3 = 1 - \frac{2}{3} \epsilon \Rightarrow \hat{x}_3 = -\frac{3}{10} \left(1 - \frac{2}{3} \epsilon\right) = -\frac{3}{10} + \frac{1}{5} \epsilon$$

$$\hat{x}_1 = \frac{\epsilon}{3} - \frac{2}{3} \left(-\frac{3}{10} + \frac{1}{5} \epsilon\right) = \frac{\epsilon}{3} + \frac{1}{5} - \frac{2}{15} \epsilon = \frac{1}{5} + \frac{1}{5} \epsilon$$

$$\hat{x}_2 = 1 - \left(-\frac{3}{10} + \frac{1}{5} \epsilon\right) = 1 + \frac{3}{10} - \frac{1}{5} \epsilon = \frac{13}{10} - \frac{1}{5} \epsilon$$

$$\begin{bmatrix} \hat{x}_1 \\ \hat{x}_2 \\ \hat{x}_3 \end{bmatrix} = \begin{bmatrix} \frac{1}{5} + \frac{1}{5} \epsilon \\ \frac{13}{10} - \frac{1}{5} \epsilon \\ -\frac{3}{10} + \frac{1}{5} \epsilon \end{bmatrix}$$

So, we have

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} \frac{1}{5} \\ \frac{13}{10} \\ -\frac{3}{10} \end{bmatrix}, \quad \begin{bmatrix} \hat{x}_1 \\ \hat{x}_2 \\ \hat{x}_3 \end{bmatrix} = \begin{bmatrix} x_1 + \delta x_1 \\ x_2 + \delta x_2 \\ x_3 + \delta x_3 \end{bmatrix} = \begin{bmatrix} \frac{1}{5} + \frac{1}{5} \epsilon \\ \frac{13}{10} - \frac{1}{5} \epsilon \\ -\frac{3}{10} + \frac{1}{5} \epsilon \end{bmatrix}.$$

$$\text{Hence, } \delta x = \begin{bmatrix} \delta x_1 \\ \delta x_2 \\ \delta x_3 \end{bmatrix} = \begin{bmatrix} \frac{1}{5} \epsilon \\ -\frac{1}{5} \epsilon \\ \frac{1}{5} \epsilon \end{bmatrix}.$$

1-norm:

$$\frac{\|\delta x\|_1}{\|x\|_1} \leq k_1(A) \frac{\|\delta b\|_1}{\|b\|_1},$$

$$\|\delta x\|_1 = \max \{ \text{col sum of } \begin{bmatrix} \frac{1}{5} \epsilon \\ -\frac{1}{5} \epsilon \\ \frac{1}{5} \epsilon \end{bmatrix} \} = \max \{ \frac{1}{5} \epsilon \} = \frac{1}{5} \epsilon$$

$$\|x\|_1 = \max \{ \text{col sum of } \begin{bmatrix} \frac{1}{5} \\ \frac{13}{10} \\ -\frac{3}{10} \end{bmatrix} \} = \max \{ \frac{6}{5} \} = \frac{6}{5}$$

$$\|\delta b\|_1 = \max \{ \text{col sum } \begin{bmatrix} \epsilon \\ 0 \\ 0 \end{bmatrix} \} = \max \{ \epsilon \} = \epsilon$$

$$\|b\|_1 = \max \{ \text{col sum } \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \} = \max \{ 2 \} = 2$$

$$\Rightarrow \frac{(\frac{1}{5} \epsilon)}{(\frac{6}{5})} \leq (5) \frac{\epsilon}{2} \Rightarrow \frac{\epsilon}{6} \leq \frac{5\epsilon}{2}$$

relative error in x $\downarrow K_\infty(A)$ relative error in b

∞ -norm:

$$\frac{\|\delta x\|_\infty}{\|x\|_\infty} \leq K_\infty(A) \frac{\|\delta b\|_\infty}{\|b\|_\infty}$$

$$\|\delta x\|_\infty = \max \left\{ \text{row sum of } \begin{bmatrix} \frac{1}{5}\epsilon & -\frac{1}{5}\epsilon & \frac{1}{5}\epsilon \\ -\frac{1}{5}\epsilon & \frac{1}{5}\epsilon & \frac{1}{5}\epsilon \\ \frac{1}{5}\epsilon & \frac{1}{5}\epsilon & \frac{1}{5}\epsilon \end{bmatrix} \right\} = \max \left\{ \frac{1}{5}\epsilon, -\frac{1}{5}\epsilon, \frac{1}{5}\epsilon \right\} = \frac{\epsilon}{5}$$

$$\|x\|_\infty = \max \left\{ \text{row sum of } \begin{bmatrix} \frac{1}{5} \\ \frac{13}{10} \\ -\frac{3}{10} \end{bmatrix} \right\} = \max \left\{ \frac{1}{5}, \frac{13}{10}, -\frac{3}{10} \right\} = \frac{13}{10}$$

$$\|\delta b\|_\infty = \max \left\{ \text{row sum } \begin{bmatrix} \epsilon \\ 0 \\ 0 \end{bmatrix} \right\} = \max \{ \epsilon, 0, 0 \} = \epsilon$$

$$\|b\|_\infty = \max \left\{ \text{row sum } \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} \right\} = \max \{ 0, 1, 1 \} = 1$$

$$\Rightarrow \frac{\left(\frac{\epsilon}{5}\right)}{\left(\frac{13}{10}\right)} \leq \left(\frac{1}{2}\right) \frac{\epsilon}{1} \Rightarrow \frac{2\epsilon}{13} \leq \frac{\epsilon}{2}$$

$\underbrace{}_{\text{relative error in } x}$
 $\downarrow K_\infty(A)$
 $\underbrace{}_{\text{relative error in } b}$

Frobenius norm:

$$\frac{\|\delta x\|_F}{\|x\|_F} \leq K_F(A) \frac{\|\delta b\|_F}{\|b\|_F}$$

$$\|\delta x\|_F = \sqrt{\left(\frac{\epsilon}{5}\right)^2 + \left(-\frac{\epsilon}{5}\right)^2 + \left(\frac{\epsilon}{5}\right)^2} = \sqrt{\frac{\epsilon^2}{25} + \frac{\epsilon^2}{25} + \frac{\epsilon^2}{25}} = \frac{\sqrt{3}\epsilon}{5}$$

$$\|x\|_F = \sqrt{\left(\frac{1}{5}\right)^2 + \left(\frac{13}{10}\right)^2 + \left(-\frac{3}{10}\right)^2} = \sqrt{\frac{1}{25} + \frac{169}{100} + \frac{9}{100}} = \frac{\sqrt{182}}{10}$$

$$\| \delta b \|_F = \sqrt{\epsilon^2 + 0^2 + 0^2} = \sqrt{\epsilon^2} = \epsilon$$

$$\| b \|_F = \sqrt{0^2 + 1^2 + 1^2} = \sqrt{2}$$

$$\Rightarrow \frac{\left(\frac{\sqrt{3}\epsilon}{5}\right)}{\left(\frac{\sqrt{182}}{10}\right)} \leq \left(\frac{\sqrt{3082}}{10}\right) \frac{(\epsilon)}{\sqrt{2}} \Rightarrow \frac{\sqrt{54}\epsilon}{91} \leq \frac{\sqrt{1541}}{10}\epsilon$$

$\underbrace{}_{\substack{\text{relative error} \\ \text{in } x}}$
 \downarrow
 \downarrow
 $\underbrace{}_{\substack{\text{relative} \\ \text{error in } b}}$

5) Prove the theorem: let $\|\cdot\|$ denote a norm on \mathbb{R}^m , and also the norm on $\mathbb{R}^{m \times m}$ that it induces. Let $I \in \mathbb{R}^{m \times m}$ denote the identity matrix. Then

$$(a) \|I\| = 1$$

$$I = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ \vdots & 0 & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & \cdots & 1 \end{bmatrix}_{m \times m} \quad x = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}_{m \times 1} \quad Ix = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}_{m \times 1} \quad (\text{the result } Ix \text{ is simply } x)$$

$$\|I\| = \max_{x \neq 0} \frac{\|Ix\|}{\|x\|} = \frac{\|x\|}{\|x\|} = 1$$

$$(b) \|Ax\| \leq \|A\| \|x\| \text{ for all } A \in \mathbb{R}^{m \times m} \text{ and } x \in \mathbb{R}^m$$

We have

$$\|A\| = \max_{x \neq 0} \frac{\|Ax\|}{\|x\|} \quad (\text{the definition of an induced norm}).$$

We take the maximum value of $\frac{\|Ax\|}{\|x\|}$ in order to find $\|A\|$.

Hence, when we do not choose the maximum for $\frac{\|Ax\|}{\|x\|}$,

$$\frac{\|Ax\|}{\|x\|} \leq \|A\|.$$

When we multiply by $\|x\|$ on both sides, this implies

$$\|Ax\| \leq \|A\| \cdot \|x\|.$$

6) What is wrong with the following reasoning?

Changing b by 1% means multiplying it by 1.01.

If $Ax = b$, then $A(1.01x) = 1.01b$, by linearity. So a 1% change of b , from b to $1.01b$, causes a 1% change in x , from x to $1.01x$. So forget the whole story about condition numbers — a 1% change in b always causes a 1% change in x .

There are a number of things wrong, starting with the fact that we cannot ignore condition numbers. A system with

a small condition number of matrix is well-conditioned, meaning there is a smaller margin between the error in x and the error in b . A system w/ large condition number of matrix is ill-conditioned, meaning the margin of error between x and b is amplified. This is shown by the relationship

$$\frac{\|\delta x\|}{\|x\|} \leq K_{\|\cdot\|}(A) \frac{\|\delta b\|}{\|b\|}$$

↓

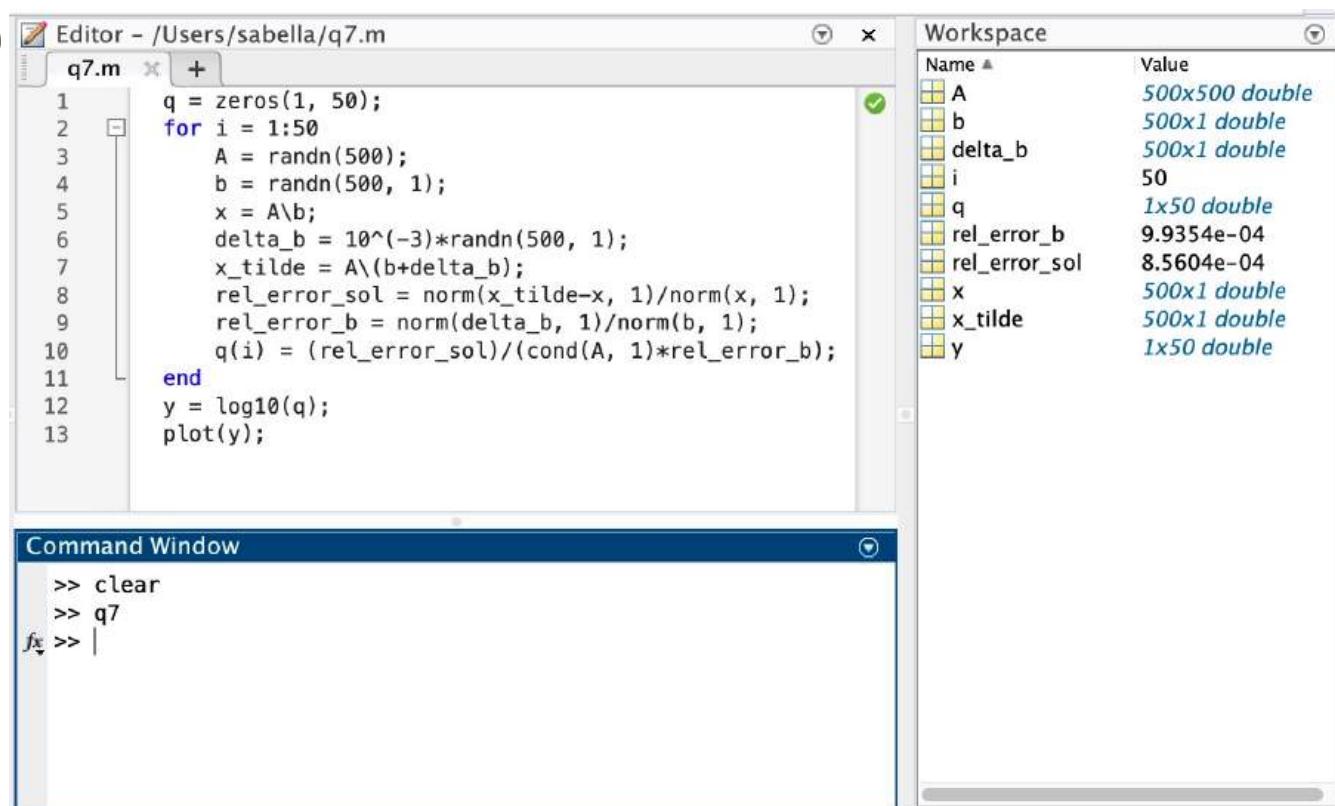
condition no.
of A

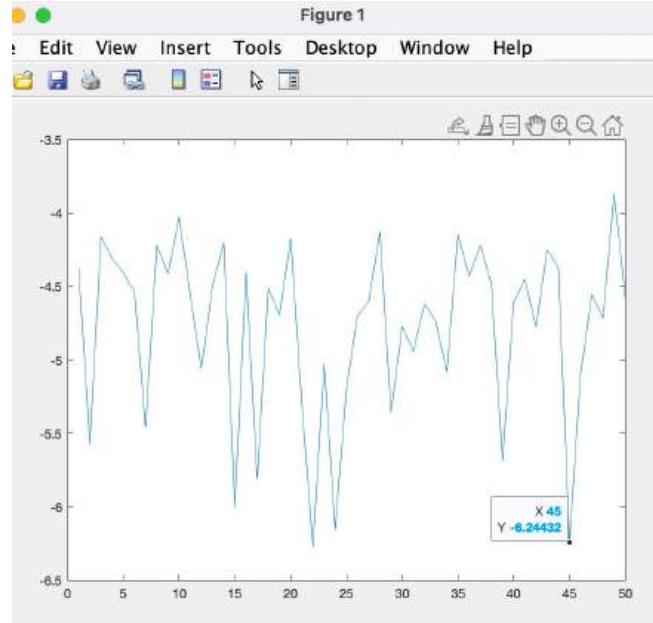
relative error
of x

relative error
in b

Hence, a 1% change of b does not necessarily cause a 1% change in x .

7)





$$\frac{\|\delta x\|}{\|x\|} \leq K_{1,1}(A) \frac{\|\delta b\|}{\|b\|}$$

relative error
in x

condition no.
of A

relative error
in b

$$q = \left(\frac{\|\delta x\|_1}{\|x\|_1} \right) / \left(K_1(A) + \frac{\|\delta b\|_1}{\|b\|_1} \right)$$

In the code, q represents the ratio

between the error in x (solution) and the perturbation of b . As shown in Figure 1, no value of $q(i)$ is greater than or equal to 1. Hence, the relative error in x is smaller than the perturbation in b (scaled by the condition number of A with respect to the 1-norm). So, our theory that well-conditioned systems lead to smaller errors in x holds.

1) Let L_1, L_2, L_3 be the matrices below:

$$L_1 = \begin{bmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad L_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ l_{31} & 0 & 1 \end{bmatrix}, \quad L_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & l_{32} & 1 \end{bmatrix}.$$

Calculate the matrix product $L_1 L_2 L_3$ and show its equal to

$$L = \begin{bmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ l_{31} & l_{32} & 1 \end{bmatrix}.$$

$$\hookrightarrow L_1 L_2 L_3 = \begin{bmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ l_{31} & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & l_{32} & 1 \end{bmatrix}$$

$$L_1 L_2 = \begin{bmatrix} 1+0+0 & 0+0+0 & 0+0+0 \\ l_{21}+0+0 & 0+1+0 & 0+0+0 \\ 0+0+l_{31} & 0+0+0 & 0+0+1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ l_{31} & 0 & 1 \end{bmatrix}$$

$$L_1 L_2 L_3 = \begin{bmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ l_{31} & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & l_{32} & 1 \end{bmatrix} = \begin{bmatrix} 1+0+0 & 0+0+0 & 0+0+0 \\ l_{21}+0+0 & 0+1+0 & 0+0+0 \\ l_{31}+0+0 & 0+0+l_{32} & 0+0+1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ l_{31} & l_{32} & 1 \end{bmatrix} = L = \begin{bmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ l_{31} & l_{32} & 1 \end{bmatrix}$$

2) Does the matrix $A = \begin{bmatrix} 6 & 0 & 6 \\ -3 & 2 & -5 \\ -4 & 1 & -4 \end{bmatrix}$ have an LU factorization? If so, compute it. If not, explain.

→ A is said to have LU factorization if there exists a lower triangular matrix L, upper triangular matrix U, such that

$$A = LU.$$

$$A = \begin{bmatrix} 6 & 0 & 6 \\ -3 & 2 & -5 \\ -4 & 1 & -4 \end{bmatrix} \quad R_1 \\ R_2 \\ R_3$$

$$R_2 + \frac{1}{2}R_1 \quad \begin{bmatrix} 1 & 0 & 0 \\ -\frac{1}{2} & 1 & 0 \\ -\frac{2}{3} & 0 & 1 \end{bmatrix} \quad \begin{bmatrix} 6 & 0 & 6 \\ 0 & 2 & -2 \\ 0 & 1 & 0 \end{bmatrix}$$

$$R_3 + \frac{2}{3}R_1$$

$$R_3 - \frac{1}{2}R_2 \quad \begin{bmatrix} 1 & 0 & 0 \\ -\frac{1}{2} & 1 & 0 \\ -\frac{2}{3} & \frac{1}{2} & 1 \end{bmatrix} \quad \begin{bmatrix} 6 & 0 & 6 \\ 0 & 2 & -2 \\ 0 & 0 & 1 \end{bmatrix}$$

$$L \qquad \qquad \qquad U$$

Yes, A has an LU factorization, as shown above.

3) Let $A = \begin{bmatrix} 0 & 1 & 1 & 2 \\ 2 & 0 & 0 & 0 \\ 1 & 1 & 1 & 3 \\ 1 & 0 & 2 & 4 \end{bmatrix}$. Does A have an LU factorization? If yes, compute. If no, explain why not, and then compute the PLU decomposition of A.

$$A = \begin{bmatrix} 0 & 1 & 1 & 2 \\ 2 & 0 & 0 & 0 \\ 1 & 1 & 1 & 3 \\ 1 & 0 & 2 & 4 \end{bmatrix}$$

The first element in row 1 is zero.

Hence we cannot use Gaussian Elimination, and so A does not have an LU factorization. We will compute the PLU decomposition of A.

$$\begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 & 1 & 2 \\ 2 & 0 & 0 & 0 \\ 1 & 1 & 1 & 3 \\ 1 & 0 & 2 & 4 \end{bmatrix} = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 1 & 1 & 2 \\ 1 & 1 & 1 & 3 \\ 1 & 0 & 2 & 4 \end{bmatrix} \quad R_3 - \frac{1}{2}R_1 \rightarrow R_3$$

$$R_4 - \frac{1}{2}R_1 \rightarrow R_4$$

$$L_{13} = \frac{1}{2}$$

$$L_{14} = \frac{1}{2}$$

$$\sim \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 1 & 1 & 2 \\ 0 & 1 & 1 & 3 \\ 0 & 0 & 2 & 4 \end{bmatrix} \quad R_3 - R_2 \rightarrow R_3$$

$$L_{23} = 1$$

$$\sim \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 1 & 1 & 2 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 2 & 4 \end{bmatrix}$$

swap
(Also swap
 L_{23} and L_{24})

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix} A = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 1 & 1 & 2 \\ 0 & 0 & 2 & 4 \\ 0 & 0 & 0 & 1 \end{bmatrix} = U$$

$$L = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ \boxed{\frac{1}{2}} & 0 & 1 & 0 \\ \boxed{\frac{1}{2}} & \boxed{1} & 0 & 1 \end{bmatrix}$$

Hence, the PLU decomposition of A is:

$$\begin{bmatrix} 0 & 1 & 1 & 2 \\ 2 & 0 & 0 & 0 \\ 1 & 1 & 1 & 3 \\ 1 & 0 & 2 & 4 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1/2 & 0 & 1 & 0 \\ 1/2 & 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 1 & 1 & 2 \\ 0 & 0 & 2 & 4 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

A P L U

4) Determine whether or not the following matrices have a Cholesky factorization; if they do, compute (by hand) the Cholesky factor R:

$$A = \begin{bmatrix} 1 & -2 & 0 \\ -2 & 13 & 6 \\ 0 & 6 & 5 \end{bmatrix}, \quad B = \begin{bmatrix} 4 & -4 & 0 \\ -4 & 4 & 0 \\ 0 & 0 & 5 \end{bmatrix}$$

In order to have a Cholesky factorization, a matrix A must be positive definite, that is, $A = A^T$ and $\forall x \neq 0$, $x^T A x > 0$. If λ is an eigenvalue of A with a corresponding eigenvector v ($v \neq 0$), $A v = \lambda v$. Hence,

$$v^T A v = v^T \lambda v = \lambda(v^T v) > 0,$$

since $\lambda > 0$ and $v^T v > 0$.

$$A \text{ is positive definite, since } A = \begin{bmatrix} 1 & -2 & 0 \\ -2 & 13 & 6 \\ 0 & 6 & 5 \end{bmatrix} = A^T,$$

$\det(A) = 1(65 - 36) + 2(-10 - 0) = 9 > 0$, and all pivots/eigenvalues are positive. Hence we proceed with

Cholesky factorization.

Define

$$R = \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ 0 & r_{22} & r_{23} \\ 0 & 0 & r_{33} \end{bmatrix}. \quad A = R^T R, \text{ so}$$

$$\begin{bmatrix} A \\ 1 & -2 & 0 \\ -2 & 13 & 6 \\ 0 & 6 & 5 \end{bmatrix} = \begin{bmatrix} r_{11} & 0 & 0 \\ r_{12} & r_{22} & 0 \\ r_{13} & r_{23} & r_{33} \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ 0 & r_{22} & r_{23} \\ 0 & 0 & r_{33} \end{bmatrix}$$

$$\begin{bmatrix} 1 & -2 & 0 \\ -2 & 13 & 6 \\ 0 & 6 & 5 \end{bmatrix} = \begin{bmatrix} r_{11}^2 & r_{11}r_{12} & r_{11}r_{13} \\ r_{12}r_{11} & r_{22}^2 + r_{12}^2 & r_{12}r_{13} + r_{22}r_{23} \\ r_{11}r_{13} & r_{12}r_{13} + r_{22}r_{23} & r_{13}^2 + r_{23}^2 + r_{33}^2 \end{bmatrix}$$

Row 1 $r_{11}^2 = 1 \Rightarrow r_{11} = 1$

$$r_{11}r_{12} = -2 \Rightarrow r_{12} = -2$$

$$r_{11}r_{13} = 0 \Rightarrow r_{13} = 0$$

Row 2 $r_{12}^2 + r_{22}^2 = 13 \Rightarrow (r_{12})^2 = 4 \Rightarrow r_{22}^2 = 13 - 4 = 9$
 $\Rightarrow r_{22} = 3$

$$r_{12}r_{13} + r_{22}r_{23} = 6 \Rightarrow (-2)(0) + (3)(r_{23}) = 6$$

 $\Rightarrow r_{23} = 2$

$$\text{Row 3} \quad r_{13}^2 + r_{23}^2 + r_{33}^2 = 5 \Rightarrow (0) + (2)^2 + r_{33}^2 = 5 \\ \Rightarrow r_{33} = 1$$

Here, $R = \begin{bmatrix} 1 & -2 & 0 \\ 0 & 3 & 2 \\ 0 & 0 & 1 \end{bmatrix}$.

Now, although $B = B^T$, B is not positive definite, since $\det(B) = 5(16 - 16) = 5(0) = 0 \not> 0$. Hence, B does not have a Cholesky factorization.

5) Let A be an $n \times n$ positive definite and symmetric matrix, and X an $n \times n$ invertible matrix. Show that $B = X^T A X$ is positive definite and symmetric.

A is positive definite, so $A = A^T$ and $x^T A x > 0$.

X is invertible, so X^{-1} exists.

Let $B = X^T A X$. Then

$$(B)^T = (X^T A X)^T \\ \Rightarrow B^T = (AX)^T (X^T)^T$$

$$\Rightarrow B^T = A^T X^T X$$

$$\Rightarrow B^T = X^T A^T X \quad (A = A^T)$$

$$\Rightarrow B^T = X^T A X = B.$$

Hence, B is symmetric.

Now, let $y \in \mathbb{R}^n$ be a vector in \mathbb{R}^n , $y \neq 0$. Then

$$\begin{aligned} y^T B y &= y^T (X^T A X) y \\ &= y^T X^T A X y \\ &= (X y)^T A (X y) \\ &= \langle X y, X y \rangle \end{aligned}$$

Let $z = X y$, since X is invertible and $X y = z$ always has a solution. Then

$$\begin{aligned} &\langle z, z \rangle \\ &= \sum_{i=1}^n (z_i)^2 > 0, \end{aligned}$$

since n^2 is always positive.

Hence, B is positive definite and symmetric.

6) Consider the system $Ax = b$, where A is upper triangular. The system above can be solved by backward substitution.

substitution, starting with solving for x_n , substituting into the second equation to get x_{n-1} , etc.

(a) Find the formula to compute x_i given x_{i+1}, \dots, x_n .

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ 0 & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & a_{nn} \end{bmatrix} x = \begin{bmatrix} x_1 \\ \vdots \\ x_i \\ \vdots \\ x_n \end{bmatrix} b = \begin{bmatrix} b_1 \\ \vdots \\ b_i \\ \vdots \\ b_n \end{bmatrix}$$

$$a_{nn} \cdot x_n = b_n \Rightarrow x_n = \frac{b_n}{a_{nn}}$$

$$a_{n-1,n-1}x_{n-1} + a_{n-1,n}x_n = b_{n-1} \Rightarrow x_{n-1} = \frac{b_{n-1} - a_{n-1,n}x_n}{a_{n-1,n-1}}$$

⋮

$$a_{ii} \cdot x_i + a_{i(i+1)}x_{i+1} + \dots + a_{in}x_n = b_i \\ \Rightarrow x_i = \frac{b_i - \sum_{j=i+1}^n a_{ij}x_j}{a_{ii}}$$

(b) Using the model code given for forward substitution, write a MATLAB code that performs backward substitution.

function $x = \text{backsub}(A, b)$

$n = \text{length}(b);$

$x = b;$

for $i = n : -1 : 1$

for $j = i+1 : n$

$x(i) = x(i) - A(i,j) * x(j);$

end

$x(i) = x(i) / A(i,i);$

end

7) Write a MATLAB code that does Gauss Elimination with partial pivoting (row switching). The input should be

- an $n \times n$ matrix A
- a column vector b of size $n \times 1$

and the output should be the solution to $Ax = b$.

function $x = \text{Gauss_Partial Pivot}(A, b)$

$n = \text{size}(A, 1);$ % get num rows of A

$Ab = [A, b];$

$x = zeros(n, 1);$

for $j = 1 : n - 1$

$[~, \text{maxRow}] = \max(\text{abs}(Ab(j:n, j)));$

$\text{maxRow} = \text{maxRow} + j - 1;$

if $\text{maxRow} \sim = j$

$Ab([j, \text{maxRow}], :) = Ab([\text{maxRow}, j], :);$

end

for $i = j + 1 : n$

$m = Ab(i, j) / Ab(j, j);$

$Ab(i, j:n+1) = Ab(i, j:n+1) - m * Ab(j, j:n+1);$

end

end

$x(n) = Ab(n, n+1) / Ab(n, n);$

```

for i = n-1 : -1 : 1
    x(i) = (Ab(i, n+1) - Ab(i, i+1:n) * x(i+1:n)) / Ab(i, i);
end
end

```

8) Write a MATLAB function with inputs

- a matrix A of size $m \times n$,
- a matrix B of size $n \times p$,

and which outputs the product matrix $A \times B$ (note that the size of this matrix should be $m \times p$). You code should check that the sizes are right and then do the multiplication using three nested "for" loops.

```
function X = matrix_product(A, B);
```

```
[m, n] = size(A);
```

```
[nB, p] = size(B);
```

```
if n ~= nB
```

```
    error ('incorrect sizes');
```

```
end
```

```
X = zeros(m, p);
```

```
for i = 1:m
```

```
    for j = 1:p
```

$$\sum_{j=1}^p 2n = p^2 + p$$

$$\sum_{i=1}^m p^2 + p = \sum_{i=1}^m p^2 + \sum_{i=1}^m p$$

```

for k = 1:n
    X(i,j) = X(i,j) + A(i,k) * B(k,j);
end
end
end

```

2 operations, n times

(a) Yes, $A \times B$ and $X = \text{matrix_product}(A, B)$ output the same result (I would include, but it's a lot).

(b) How many FLOPs does the code use? Find a formula in terms of m, n, p .

inner loop: 2 operations (+, *) n times = $2n$

$$\sum_{k=1}^n 2 = 2n$$

middle loop: iterates p times $\sum_{j=1}^p 2n = p^2 + p$

outer loop: iterates m times $\sum_{i=1}^m p^2 + p =$
 $\sum_{i=1}^m p^2 + \sum_{i=1}^m p = \frac{m(m+1)(2m+1)}{6} + \frac{m(m+1)}{2} = \frac{1}{3}m^3 + m^2 + \frac{2}{3}m$

$$\sum_{i=1}^m \left(\sum_{j=1}^p \left(\sum_{k=1}^n 2 \right) \right) = \frac{1}{3}m^3 + m^2 + \frac{2}{3}m$$

9) Use for and if loops to write two MATLAB functions

that take as an input

- $n \times n$ matrix A
- $n \times n$ matrix B
- $n \times 1$ vector x.

Check that the inputs are square matrices of the same size. Then have the first function compute ABx through $(AB)x$ and the second through $A(Bx)$. Then

(a) Take a screenshot of your first function

```
Editor - /Users/sabella/compute_ABx.m
compute_ABx.m + [x]

1 function first = compute_ABx(A, B, x)
2     [mA, nA] = size(A);
3     [mB, nB] = size(B);
4     if mA ~= nA || mB ~= nB || mA ~= mB
5         error('matrices must be square and same size');
6     end
7     AB = zeros(mA, nB);
8     for i = 1 : mA
9         for j = 1 : nB
10            for k = 1 : nA
11                AB(i, j) = AB(i, j) + A(i, k)*B(k, j); 2 op * n ] ] * n
12            end
13        end
14    end
15    first = zeros(mA, 1);
16    for i = 1 : mA
17        for j = 1 : nB
18            first(i) = first(i) + AB(i, j)*x(j); 2 op * n ] * n
19        end
20    end
21 end
```

Handwritten annotations:

- Red bracket spanning the innermost loop (k): $2 \text{ op} * n$
- Red bracket spanning the outer loops (i, j): $* n$
- Red bracket spanning the entire inner loop structure: $O(n^3)$
- Red bracket spanning the inner loops (j, k): $2 \text{ op} * n$
- Red bracket spanning the outer loops (i, j): $* n$
- Red bracket spanning the entire inner loop structure: $O(n^2)$

(b) Take a screenshot of your second function

```

1 function second = compute_A_Bx(A, B, x)
2     [mA, nA] = size(A);
3     [mB, nB] = size(B);
4     if mA ~= nA || mB ~= nB || mA ~= mB
5         error('matrices must be square and same size');
6     end
7     Bx = zeros(mB, 1);
8     for i = 1 : mB
9         for j = 1 : nB
10            Bx(i) = Bx(i) + B(i, j)*x(j);
11        end
12    end
13    second = zeros(mA, 1);
14    for i = 1 : mA
15        for j = 1 : nA
16            second(i) = second(i) + A(i, j)*Bx(j);
17        end
18    end
19 end
20
21

```

$$\left[2 \text{operations} * n \right] * n = O(n^2)$$

$$\left[2 \text{op} * n \right] * n = O(n^2)$$

(c) calculate the number of flops for both approaches, theoretically (in terms of n) and compare: which one is better to use?

Function 1: $\sum_{i=1}^n \left(\sum_{j=1}^n 2n \right)$ (2 operations inside j-loop
- the two operations iterate n times)

$$= \sum_{i=1}^n n^2 + n = \sum_{i=1}^n n^2 + \sum_{i=1}^n n$$

$$= \frac{n(n+1)(2n+1)}{6} + \frac{n(n+1)}{2} = \frac{1}{6}(2n^3 + 3n^2 + n) + \frac{1}{2}(n^2 + n)$$

$$= \frac{1}{3}n^3 + \frac{1}{2}n^2 + \frac{1}{6}n + \frac{1}{2}n^2 + \frac{1}{2}n = \frac{1}{3}n^3 + n^2 + \frac{2}{3}n = O(n^3)$$

$$(\text{second loop}) \sum_{i=1}^n 2i = n^2 + n = O(n^2)$$

$$\text{Total: } \frac{1}{3}n^3 + n^2 + \frac{2}{3}n + n^2 + n = \frac{1}{3}n^3 + 2n^2 + \frac{5}{3}n = O(n^3)$$

$$\underline{\text{Function 2}}: \sum_{i=1}^n 2i = n^2 + n = O(n^2)$$

$$(\text{second loop}) \sum_{i=1}^n 2i = n^2 + n = O(n^2)$$

$$\text{Total: } 2n^2 + 2n = O(n^2)$$

The second function is better to use since it is faster and takes less operations.