

Homework #6

21.6.) Proof: Let $A \in \mathbb{C}^{m \times m}$ such that $\sum_{j \neq k} |a_{jk}| < |a_{kk}|$ and let $\tilde{A} = L_1 P_1 A$ be the matrix obtained after the first step of the Gaussian elimination with partial pivoting algorithm. Note that $|a_{11}| > |a_{i1}|$ for all $2 \leq i \leq m$, so $P_1 = I$ and $\tilde{A} = L_1 A$. Show that if Gaussian elimination with partial pivoting is applied to A , no row interchanges take place. Observe that if we can show the $k \times k$ submatrix of \tilde{A} is strictly column diagonally dominant, then the largest element in the k th column of this submatrix must always lie at the k th row, k th column, which means no row swaps would be needed. Let this submatrix be called $\tilde{\tilde{A}}$. Then by algorithm 20.1 we see that $\tilde{\tilde{a}}_{j,im} = a_{j,im} - \left(\frac{a_{ji}}{a_{11}}\right)a_{1,im}$, so if $2 \leq k \leq m$ then $\tilde{\tilde{a}}_{jk} = a_{jk} - \frac{a_{j1}}{a_{11}}a_{1k}$. This means

$$\begin{aligned} \sum_{j \neq k} |\tilde{\tilde{a}}_{jk}| &= \sum_{j \neq k} \left| a_{jk} - \frac{a_{j1}}{a_{11}}a_{1k} \right| \leq \sum_{j \neq k} |a_{jk}| + \sum_{j \neq k} \left| \frac{a_{j1}}{a_{11}}a_{1k} \right| \\ &= \sum_{j \neq k} |a_{jk}| + \frac{|a_{1k}|}{|a_{11}|} \sum_{j \neq k} |a_{j1}| \quad \left(\text{Note } \sum_{j \neq k} |a_{j1}| + |a_{k1}| = \sum_{j \neq 1} |a_{j1}| < |a_{11}| \right) \\ &< |a_{kk}| - |a_{1k}| + \frac{|a_{1k}|}{|a_{11}|} (|a_{11}| - |a_{k1}|) \\ &= |a_{kk}| - |a_{1k}| + |a_{1k}| - |a_{1k}| \frac{|a_{k1}|}{|a_{11}|} \\ &= |a_{kk}| - |a_{1k}| \frac{|a_{k1}|}{|a_{11}|} = |a_{kk}| - \left| a_{1k} \frac{a_{k1}}{a_{11}} \right| \\ &\leq \left| a_{kk} - a_{1k} \frac{a_{k1}}{a_{11}} \right| = |\tilde{\tilde{a}}_{kk}| \end{aligned}$$

Thus $\tilde{\tilde{A}}$ is also strictly column diagonally dominant after each step of Gaussian elimination is applied, and the top left element $\tilde{\tilde{a}}_{ii}$ satisfies $|\tilde{\tilde{a}}_{ji}| < |\tilde{\tilde{a}}_{ii}|$ for all $j > i$. Therefore no row swaps take place when we apply Gaussian elimination with partial pivoting. \square

2.) a) Proof: Let $A \in \mathbb{C}^{n \times n}$ be invertible, $B \in \mathbb{C}^{n \times m}$ be a rank- k matrix such that $B = XCY^T$ with $X \in \mathbb{C}^{n \times k}$, $Y \in \mathbb{C}^{k \times m}$ and $C \in \mathbb{C}^{k \times k}$ where C is invertible. Show $(A + XCY)^{-1} = A^{-1} - A^{-1}X(C^{-1} + YA^{-1}X)^{-1}YA^{-1}$.
Observe that we can solve the following system for T in two different ways:

① Substitution:

$$\begin{bmatrix} A & X \\ Y & -C^T \end{bmatrix} \begin{bmatrix} T \\ R \end{bmatrix} = \begin{bmatrix} I \\ 0 \end{bmatrix} \Rightarrow \begin{cases} AT + XR = I \\ YT - C^T R = 0 \end{cases}$$

$$YT - C^T R = 0 \Rightarrow YT = C^T R \Rightarrow R = CYT$$

$$\Rightarrow AT + X(CYT) = I \Rightarrow AT + XCYT = I \Rightarrow (A + XCY)T = I$$

$$\Rightarrow T = (A + XCY)^{-1}$$

② Gaussian elimination:

From the lecture 20 notes we see that

$$\begin{bmatrix} A & X \\ Y & -C^T \end{bmatrix} \begin{bmatrix} T \\ R \end{bmatrix} = \begin{bmatrix} I \\ 0 \end{bmatrix} \Rightarrow \begin{bmatrix} I & 0 \\ -YA^{-1} & I \end{bmatrix} \cdot \begin{bmatrix} A & X \\ Y & -C^T \end{bmatrix} \begin{bmatrix} T \\ R \end{bmatrix} = \begin{bmatrix} I & 0 \\ -YA^{-1} & I \end{bmatrix} \begin{bmatrix} I \\ 0 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} A & X \\ 0 & -(C^T + YA^{-1}X) \end{bmatrix} \begin{bmatrix} T \\ R \end{bmatrix} = \begin{bmatrix} I \\ -YA^{-1} \end{bmatrix} \Rightarrow \begin{cases} AT + XR = I \\ -(C^T + YA^{-1}X)R = -YA^{-1} \end{cases}$$

$$-(C^T + YA^{-1}X)R = -YA^{-1} \Rightarrow R = (C^T + YA^{-1}X)^{-1}YA^{-1}$$

$$\Rightarrow AT + X(C^T + YA^{-1}X)^{-1}YA^{-1} = I \Rightarrow AT + X(C^T + YA^{-1}X)^{-1}YA^{-1} = I$$

$$\Rightarrow T + A^{-1}X(C^T + YA^{-1}X)^{-1}YA^{-1} = A^{-1}$$

$$\Rightarrow T = A^{-1} - A^{-1}X(C^T + YA^{-1}X)^{-1}YA^{-1}$$

Thus we obtain $(A + XCY)^{-1} = A^{-1} - A^{-1}X(C^T + YA^{-1}X)^{-1}YA^{-1}$. \square

2.) b) Pseudocode for $(C + ab^T)x = v$.

1.) Calculate the diagonal entries of Δ^{-1} where Δ^{-1} is the inverse of Δ in $FCF^T = \Delta$. This is done by calculating Fc where c is the first column of C , and then taking the reciprocals of Fc .

2.) Calculate $C^{-1}v$ and $C^{-1}a$ by calculating $F^{-1}\Delta^{-1}Fv$ and $F^{-1}\Delta^{-1}Fa$.

3.) Calculate $x = (C + ab^T)^{-1}v$ using the Sherman-Morrison Woodbury formula with $A=C$, $X=a$, $C=1$, and $Y=b^T$.

4.) Return x .

This requires 1 fft for step 1 and 2 ffts and 2 iffts for step 2, so there are 5 total fft calculations. Step 1 is $O(m \log m)$, step 2 is $O(m \log m)$, step 3 is $O(m)$, and step 4 is $O(1)$, so the overall computational complexity is $O(m \log m)$.

3.) a) Let $M = A + B$ where $B = \begin{bmatrix} \uparrow & & & \uparrow \\ b_1 & 0 & \dots & 0 \\ \downarrow & & & \downarrow \\ & & & b_2 \end{bmatrix}$.

We see that if $X = [b_1 \ b_2]$ and $Y = [e_1 \ e_m]$, then

$$B = X I_{2 \times 2} Y.$$

By inspection, B only has two non zero columns, so if b_1 & b_2 are linearly independent then $\text{rank}(B) = 2$. If b_1 & b_2 are linearly dependent, then $\text{rank}(B) = 1$, so the maximum value of $\text{rank}(B)$ is 2.

b.) Pseudocode for $M(t)x(t) = f(t)$.

1.) Compute the LU decomposition of A

2.) Loop through all times t_k

3.) Call stepsolve function in each iteration, where stepsolve does the following:

3.)b) continued.)

- (i.) Compute $A^{-1}f(t_k)$ by using the LU decomp. of A .
- (ii.) Form $X(t_k)$ and $Y(t_k)$ where $B(t_k) = X(t_k)I_{2 \times 2}Y(t_k)^T$.
- (iii.) Compute $A^{-1}X(t_k)$ by using the LU decomp. of A .
- (iv.) Use the Sherman-Morrison Woodbury formula to calculate X_k and return X_k .

The time complexity of step 1 is $\mathcal{O}(m^3)$, step 3 is contained in the loop in step 2, and the loop occurs p times. For each iteration of the loop, step (i) is $\mathcal{O}(m^2)$, step (ii) is $\mathcal{O}(m)$, step (iii) is $\mathcal{O}(m^2)$, and step (iv) is $\mathcal{O}(m^2)$. This means stepsolve is $\mathcal{O}(m^2)$, so since it occurs p times the time complexity of the loop is $\mathcal{O}(m^2p)$. The total time complexity then becomes $\mathcal{O}(m^3 + m^2p)$.