

Homework 4

Math 5424

Numerical Linear Algebra

Alexander Novotny

Zuriah Quinton

November 6, 2023

1. Let \mathbf{L} be a lower triangular matrix and solve $\mathbf{L}\vec{x} = \vec{b}$ by forward substitution. Show that barring overflow or underflow, the computed solution \hat{x} satisfies $(\mathbf{L} + \delta\mathbf{L})\hat{x} = \vec{b}$, where $|\delta l_{ij}| \leq n\varepsilon|l_{ij}|$, where ε is the machine precision. This means that forward substitution is backward stable. Argue that backward substitution for solving upper triangular systems satisfies the same bound.

Proof. We have, for the true solution \vec{x} ,

$$x_i = \left(b_i - \sum_{k=1}^{i-1} l_{ik}x_k \right) / l_{ii}.$$

For the solution \hat{x} computed via backwards substitution,

$$\hat{x}_i = \left(b_i - \sum_{k=1}^{i-1} \left[l_{ik}\hat{x}_k(1 + \varepsilon_k^*) \prod_{j=k}^{i-1} (1 + \varepsilon_j^+) \right] \right) (1 + \varepsilon^-)(1 + \varepsilon') / l_{ii}$$

(where $|\varepsilon_k^*|, |\varepsilon_k^+|, |\varepsilon^-|, |\varepsilon'| < eps$)

$$\begin{aligned} &= \left(b_i - \sum_{k=1}^{i-1} \left[l_{ik}(1 + \varepsilon_k^*) \prod_{j=k}^{i-1} (1 + \varepsilon_j^+) \right] \hat{x}_k \right) / \left(\frac{l_{ii}}{(1 + \varepsilon^-)(1 + \varepsilon')} \right) \\ &= \left(b_i - \sum_{k=1}^{i-1} \left[l_{ik} \left(1 + \varepsilon_k^* + \sum_{j=k}^{i-1} \varepsilon_j^+ + \mathcal{O}(eps^2) \right) \right] \hat{x}_k \right) / (l_{ii}(1 - \varepsilon^- - \varepsilon' + \mathcal{O}(eps^2))). \end{aligned}$$

Then we have that $(\mathbf{L} + \delta\mathbf{L})\hat{x} = \vec{b}$ for

$$\begin{aligned} \delta l_{ij} &= \begin{cases} l_{ij}(-\varepsilon^- - \varepsilon'), & i = j \\ l_{ij}(\varepsilon_j^* + \sum_{k=j}^{i-1} \varepsilon_k^+), & i \neq j \end{cases} \\ &= l_{ij} \begin{cases} -\varepsilon^- - \varepsilon', & i = j \\ \varepsilon_j^* + \sum_{k=j}^{i-1} \varepsilon_k^+, & i \neq j \end{cases} \\ \Rightarrow |\delta l_{ij}| &= |l_{ij}| \begin{cases} |(-\varepsilon^- - \varepsilon')|, & i = j \\ |\varepsilon_j^* + \sum_{k=j}^{i-1} \varepsilon_k^+|, & i \neq j \end{cases} \end{aligned}$$

(note that for $i = 1$, $\varepsilon^- = 0$, since we subtract by 0 - which will never round)

$$\begin{aligned} &\leq |l_{ij}| \begin{cases} |\varepsilon'|, & i = j = 1 \\ |\varepsilon^-| + |\varepsilon'|, & i = j \neq 1 \\ |\varepsilon_j^*| + \sum_{k=j}^{i-1} |\varepsilon_k^+|, & i \neq j \end{cases} \\ &\leq |l_{ij}| \begin{cases} eps, & i = j = 1 \\ 2eps, & i = j \neq 1 \\ (i - j + 1)eps, & i \neq j \end{cases} \\ &\leq n \cdot eps |l_{ij}| \end{aligned}$$



2. Matrix \mathbf{A} is called *strictly column diagonally dominant*, or diagonally dominant for short, if

$$|\alpha_{ii}| > \sum_{j=1, j \neq i}^n |\alpha_{ji}| \quad (1)$$

Show that Gaussian elimination with partial pivoting does not actually permute any rows, i.e., that it is identical to Gaussian elimination without pivoting. Hint: Show that after one step of Gaussian elimination, the trailing $(n-1)$ -by- $(n-1)$ submatrix, the *Schur complement* of α_{11} in \mathbf{A} , is still diagonally dominant.

Proof. For a diagonally dominant matrix \mathbf{A} , the first step of Gaussian elimination with partial pivoting does not permute any rows, since $|\alpha_{ii}| > \sum_{j=1, j \neq i}^n |\alpha_{ji}|$ and thus $|\alpha_{ii}| > |\alpha_{ji}|$ for any $j \neq i$. Then, we have the decomposition for \mathbf{A} given by the first step of Gaussian elimination as

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} \alpha_{11} & \vec{c}^\top \\ \vec{a} & \hat{\mathbf{A}} \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 \\ \vec{l} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \alpha_{11} & \vec{c}^\top \\ 0 & \hat{\mathbf{A}} - \vec{l}\vec{c}^\top \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 \\ \vec{l} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \alpha_{11} & \vec{c}^\top \\ 0 & \mathbf{B} \end{bmatrix} \end{aligned}$$

where $\vec{l} = \frac{\vec{a}}{\alpha_{11}}$. Note that $c_i = \alpha_{1,i+1}$, $l_i = \frac{\alpha_{i+1,1}}{\alpha_{11}}$. Then, denote the submatrix $\mathbf{B} = \hat{\mathbf{A}} - \vec{l}\vec{c}^\top$. Now, it is sufficient to show that \mathbf{B} is diagonally dominant, since then we could continue each step of Gaussian elimination recursively always getting a diagonally dominant submatrix. First note

$$\beta_{ji} = \alpha_{j+1,i+1} - \frac{\alpha_{j+1,1}\alpha_{1,i+1}}{\alpha_{11}} \quad (2)$$

is the expression for the elements in \mathbf{B} by definition. Further, from eq. (1) we can split one of the terms out of the sum, and we have

$$|\alpha_{ii}| - \sum_{\substack{j=1 \\ j \neq i, k}} |\alpha_{ji}| > |\alpha_{ki}|. \quad (3)$$

We have

$$\begin{aligned} |\beta_{ii}| &\stackrel{(2)}{=} \left| \alpha_{i+1,i+1} - \frac{\alpha_{i+1,1}\alpha_{1,i+1}}{\alpha_{11}} \right| \\ &\geq |\alpha_{i+1,i+1}| - \left| \frac{\alpha_{i+1,1}\alpha_{1,i+1}}{\alpha_{11}} \right| \quad \text{by reverse triangle inequality} \\ &\stackrel{(1)}{>} \sum_{\substack{j=1 \\ j \neq i+1}}^n |\alpha_{j,i+1}| - \left| \frac{\alpha_{i+1,1}\alpha_{1,i+1}}{\alpha_{11}} \right| \\ &\stackrel{(3)}{>} \sum_{\substack{j=1 \\ j \neq i+1}}^n |\alpha_{j,i+1}| - \frac{|\alpha_{1,i+1}|}{|\alpha_{11}|} \left(|\alpha_{11}| - \sum_{\substack{j=1 \\ j \neq 1, i+1}}^n |\alpha_{j,1}| \right) \\ &= \sum_{\substack{j=1 \\ j \neq i+1}}^n |a_{j,i+1}| - |\alpha_{1,i+1}| + \frac{|\alpha_{1,i+1}|}{|\alpha_{11}|} \sum_{\substack{j=1 \\ j \neq 1, i+1}}^n |\alpha_{j,1}| \\ &\quad \text{split } j=1 \text{ term out} \\ &= \cancel{|a_{1,i+1}|} + \sum_{\substack{j=2 \\ j \neq i+1}}^n |a_{j,i+1}| - \cancel{|\alpha_{1,i+1}|} + \frac{|\alpha_{1,i+1}|}{|\alpha_{11}|} \sum_{\substack{j=2 \\ j \neq i+1}}^n |\alpha_{j,1}| \end{aligned}$$

$$\begin{aligned}
&= \sum_{\substack{j=2 \\ j \neq i+1}}^n |a_{j,i+1}| + \frac{|\alpha_{1,i+1}|}{|\alpha_{11}|} \sum_{\substack{j=2 \\ j \neq i+1}}^n |\alpha_{j,1}| \\
&= \sum_{\substack{j=2 \\ j \neq i+1}}^n \left(|a_{j,i+1}| + \frac{|\alpha_{1,i+1}| |\alpha_{j,1}|}{|\alpha_{11}|} \right) \\
&\geq \sum_{\substack{j=2 \\ j \neq i+1}}^n \left| a_{j,i+1} - \frac{\alpha_{1,i+1} \alpha_{j,1}}{\alpha_{11}} \right| \\
&\stackrel{(2)}{=} \sum_{\substack{j=1 \\ j \neq i}}^{n-1} |\beta_{ji}|.
\end{aligned}$$

Thus, \mathbf{B} is diagonally dominant. Then, by induction each iteration of Gaussian elimination on a diagonally dominant matrix produces a diagonally dominant submatrix, so no pivoting is required. \odot

3. Let \mathbf{A} , \mathbf{B} , and \mathbf{C} be matrices with dimensions such that the product $\mathbf{A}^\top \mathbf{C} \mathbf{B}^\top$ is well defined. Let \mathcal{X} be the set of matrices \mathbf{X} minimizing $\|\mathbf{A} \mathbf{X} \mathbf{B} - \mathbf{C}\|_F$, and let \mathbf{X}_0 be the unique member of \mathcal{X} minimizing $\|\mathbf{X}\|_F$. Show that $\mathbf{X}_0 = \mathbf{A}^+ \mathbf{C} \mathbf{B}^+$. Hint: Use the SVDs of \mathbf{A} and \mathbf{B} .

Proof. Let \mathbf{B} be of rank r with left singular vectors \vec{u}_i , right singular vectors \vec{v}_i , and non-zero singular values σ_i for $i = 1, \dots, r$. From class, to minimize the given norm, we must have

$$\begin{aligned}
\mathbf{A}^+ \mathbf{C} &= \mathbf{X} \mathbf{B} \\
\iff \mathbf{A}^+ \mathbf{C} \mathbf{B}^+ &= \mathbf{X} \mathbf{B} \mathbf{B}^+ \\
\iff \mathbf{A}^+ \mathbf{C} \mathbf{B}^+ &= \mathbf{X} \left(\sum_{i=1}^r \sigma_i \vec{u}_i \vec{v}_i^\top \right) \left(\sum_{i=1}^r \frac{1}{\sigma_i} \vec{v}_i \vec{u}_i^\top \right) \\
&= \mathbf{X} \sum_{i=1}^r \cancel{\vec{u}_i \vec{v}_i^\top} \vec{v}_i \vec{u}_i^\top \\
&= \mathbf{X} \sum_{i=1}^r \vec{u}_i \vec{u}_i^\top.
\end{aligned}$$

Then our set

$$\mathcal{X} = \left\{ \mathbf{X} \left| \mathbf{X} \sum_{i=1}^r \vec{u}_i \vec{u}_i^\top = \mathbf{A}^+ \mathbf{C} \mathbf{B}^+ \right. \right\}.$$

The matrix $\sum_{i=1}^r \vec{u}_i \vec{u}_i^\top$ is the projection matrix onto $\text{Range}(\mathbf{B})$, so \mathcal{X} is nonempty if and only if $\mathbf{A}^+ \mathbf{C} \mathbf{B}^+ \in \mathcal{X}$, since if the row vectors $(\mathbf{A}^+ \mathbf{C} \mathbf{B}^+)_i^\top \notin \text{Range}(\mathbf{B})$, then this equality could never hold. We ensure that \mathcal{X} is nonempty by verifying

$$\begin{aligned}
\mathbf{A}^+ \mathbf{C} \mathbf{B}^+ \sum_{i=1}^r \vec{u}_i \vec{u}_i^\top &= \mathbf{A}^+ \mathbf{C} \left(\sum_{i=1}^r \frac{1}{\sigma_i} \vec{v}_i \vec{u}_i^\top \right) \sum_{i=1}^r \vec{u}_i \vec{u}_i^\top \\
&= \mathbf{A}^+ \mathbf{C} \sum_{i=1}^r \frac{1}{\sigma_i} \vec{v}_i \cancel{\vec{u}_i^\top} \vec{u}_i \vec{u}_i^\top \\
&= \mathbf{A}^+ \mathbf{C} \mathbf{B}^+.
\end{aligned}$$

As well, we must have $\mathbf{X} = \mathbf{A}^+ \mathbf{C} \mathbf{B}^+ + \tilde{\mathbf{X}}$, where $\tilde{\mathbf{X}} \sum_{i=1}^r \vec{u}_i \vec{u}_i^\top = 0$, i.e. row vectors $\tilde{\mathbf{X}}_i^\top \perp \text{Range}(\mathbf{B})$ for all $\mathbf{X} \in \mathcal{X}$. Then

$$\begin{aligned}
\|\mathbf{X}\|_F^2 &= \|\mathbf{A}^+ \mathbf{C} \mathbf{B}^+ + \tilde{\mathbf{X}}\|_F^2 \\
&= \sum_i \|(\mathbf{A}^+ \mathbf{C} \mathbf{B}^+)_i^\top + \tilde{\mathbf{X}}_i^\top\|_2^2
\end{aligned}$$

(since $(\mathbf{A}^+ \mathbf{C} \mathbf{B}^+)_i^\top \perp \tilde{\mathbf{X}}_i^\top$, by Pythagorean Theorem)

$$\geq \sum_i \|(\mathbf{A}^+ \mathbf{C} \mathbf{B}^+)_i^\top\|_2^2$$

(note that equality only holds for $\tilde{\mathbf{X}} = 0$)

$$= \|\mathbf{A}^+ \mathbf{C} \mathbf{B}^+\|_F^2,$$

so $\mathbf{X}_0 = \mathbf{A}^+ \mathbf{C} \mathbf{B}^+$.



4. Show that the Moore—Penrose pseudoinverse of \mathbf{A} satisfies the following identities:

$$\begin{aligned}\mathbf{A} \mathbf{A}^+ \mathbf{A} &= \mathbf{A}, \\ \mathbf{A}^+ \mathbf{A} \mathbf{A}^+ &= \mathbf{A}^+, \\ \mathbf{A}^+ \mathbf{A} &= (\mathbf{A}^+ \mathbf{A})^\top, \\ \mathbf{A} \mathbf{A}^+ &= (\mathbf{A} \mathbf{A}^+)^\top.\end{aligned}$$

Proof. Assuming we have a full rank matrix \mathbf{A} with more rows than columns, $\mathbf{A}^+ = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top$, and we have

$$\begin{aligned}\mathbf{A} \mathbf{A}^+ \mathbf{A} &= \mathbf{A} (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{A} \\ &= \cancel{\mathbf{A} (\mathbf{A}^\top \mathbf{A})}^{-1} \cancel{(\mathbf{A}^\top \mathbf{A})} \\ &= \mathbf{A}.\end{aligned}$$

Then,

$$\begin{aligned}\mathbf{A}^+ \mathbf{A} \mathbf{A}^+ &= (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{A} \mathbf{A}^+ \\ &= \cancel{(\mathbf{A}^\top \mathbf{A})}^{-1} \cancel{(\mathbf{A}^\top \mathbf{A})} \mathbf{A}^+ \\ &= \mathbf{A}^+.\end{aligned}$$

Further,

$$\begin{aligned}\mathbf{A}^+ \mathbf{A} &= (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{A} \\ &= \cancel{(\mathbf{A}^\top \mathbf{A})}^{-1} \cancel{(\mathbf{A}^\top \mathbf{A})} \\ &= \mathbf{I},\end{aligned}$$

and

$$\begin{aligned}(\mathbf{A}^+ \mathbf{A})^\top &= \mathbf{A}^\top (\mathbf{A}^+)^{\top} \\ &= \mathbf{A}^\top ((\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top)^\top \\ &= \mathbf{A}^\top \mathbf{A} ((\mathbf{A}^\top \mathbf{A})^{-1})^\top \\ &= \mathbf{A}^\top \mathbf{A} ((\mathbf{A}^\top \mathbf{A})^\top)^{-1} \\ &= \cancel{(\mathbf{A}^\top \mathbf{A})} (\cancel{\mathbf{A}^\top \mathbf{A}})^{-1} \\ &= \mathbf{I},\end{aligned}$$

so $\mathbf{A}^+ \mathbf{A} = (\mathbf{A}^+ \mathbf{A})^\top$. Finally,

$$\begin{aligned}(\mathbf{A} \mathbf{A}^+)^\top &= (\mathbf{A}^+)^{\top} \mathbf{A}^\top \\ &= ((\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top)^\top \mathbf{A}^\top \\ &= \mathbf{A} ((\mathbf{A}^\top \mathbf{A})^{-1})^\top \mathbf{A}^\top \\ &= \mathbf{A} ((\mathbf{A}^\top \mathbf{A})^\top)^{-1} \mathbf{A}^\top \\ &= \mathbf{A} (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \\ &= \mathbf{A} \mathbf{A}^+.\end{aligned}$$

Then for a full rank matrix with more columns than rows, $\mathbf{A}^+ = \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1}$, and we have

$$\begin{aligned}\mathbf{A}\mathbf{A}^+ \mathbf{A} &= \mathbf{A}\mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A} \\ &= \cancel{(\mathbf{A}\mathbf{A}^\top)} \cancel{(\mathbf{A}\mathbf{A}^\top)^{-1}} \mathbf{A} \\ &= \mathbf{A}.\end{aligned}$$

Then,

$$\begin{aligned}\mathbf{A}^+ \mathbf{A}\mathbf{A}^+ &= \mathbf{A}^+ \mathbf{A}\mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \\ &= \mathbf{A}^+ \cancel{(\mathbf{A}\mathbf{A}^\top)} \cancel{(\mathbf{A}\mathbf{A}^\top)^{-1}} \\ &= \mathbf{A}^+.\end{aligned}$$

Further,

$$\begin{aligned}(\mathbf{A}^+ \mathbf{A})^\top &= ((\mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1}) \mathbf{A})^\top \\ &= \mathbf{A}^\top (\mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1})^\top \\ &= \mathbf{A}^\top ((\mathbf{A}\mathbf{A}^\top)^{-1})^\top \mathbf{A} \\ &= \mathbf{A}^\top ((\mathbf{A}\mathbf{A}^\top)^\top)^{-1} \mathbf{A} \\ &= \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A} \\ &= \mathbf{A}^+ \mathbf{A}.\end{aligned}$$

Finally,

$$\begin{aligned}(\mathbf{A}\mathbf{A}^+)^\top &= (\mathbf{A}\mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1})^\top \\ &= \mathbf{I},\end{aligned}$$

and

$$\begin{aligned}\mathbf{A}\mathbf{A}^+ &= \mathbf{A}\mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \\ &= \mathbf{I},\end{aligned}$$

so $\mathbf{A}\mathbf{A}^+ = (\mathbf{A}\mathbf{A}^+)^\top$. 😊

5. (a) Describe a variant of Gaussian elimination that introduces zeros into the columns of \mathbf{A} in the order $n : -1 : 2$ and which produces the factorization $\mathbf{A} = \mathbf{U}\mathbf{L}$ where \mathbf{U} is the unit upper triangular and \mathbf{L} is lower triangular.

Answer. Note that for a 1×1 matrix \mathbf{A} , we have $\mathbf{A} = \mathbf{I}\mathbf{A}$, where \mathbf{I} is unit upper triangular, and \mathbf{A} is lower triangular. Then for an $n \times n$ matrix \mathbf{A} , we have

$$\mathbf{A} = \begin{bmatrix} \hat{\mathbf{A}} & \vec{a} \\ \vec{c}^\top & \alpha_{nn} \end{bmatrix} \quad (4)$$

$$= \begin{bmatrix} \mathbf{I} & \vec{l} \\ \vec{0}^\top & 1 \end{bmatrix} \begin{bmatrix} \hat{\mathbf{A}} - \vec{l}\vec{c}^\top & \vec{0} \\ \vec{c}^\top & \alpha_{nn} \end{bmatrix}, \quad (5)$$

for $\vec{l} = \frac{\vec{a}}{\alpha_{11}}$. Assume, by induction, that $\hat{\mathbf{A}} - \vec{l}\vec{c}^\top = \mathbf{U}_1\mathbf{L}_1$ for some $n-1 \times n-1$ unit upper triangular matrix \mathbf{U}_1 and lower triangular matrix \mathbf{L}_1 . Then we have

$$\begin{aligned}\mathbf{A} &= \begin{bmatrix} \mathbf{I} & \vec{l} \\ \vec{0}^\top & 1 \end{bmatrix} \begin{bmatrix} \mathbf{U}_1\mathbf{L}_1 & \vec{0} \\ \vec{c}^\top & \alpha_{nn} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{U}_1 & \vec{l} \\ \vec{0}^\top & 1 \end{bmatrix} \begin{bmatrix} \mathbf{L}_1 & \vec{0} \\ \vec{c}^\top & \alpha_{nn} \end{bmatrix} \\ &= \mathbf{U}\mathbf{L},\end{aligned}$$

where \mathbf{U} is the unit upper triangular and \mathbf{L} is lower triangular. This produces a recursive algorithm, which can be seen in algorithm 1.

Algorithm 1: A recursive algorithm to factorize $\mathbf{A} = \mathbf{UL}$ where \mathbf{U} is unit upper triangular and \mathbf{L} is lower triangular.

Data: \mathbf{A}
Result: \mathbf{U}, \mathbf{L}

```

1  $\vec{l} \leftarrow \vec{a}/\alpha_{11};$ 
2 Factorize  $\hat{\mathbf{A}} - \vec{l}\vec{c}^\top = \mathbf{U}_1\mathbf{L}_1;$ 
3  $\mathbf{U} \leftarrow \begin{bmatrix} \mathbf{U}_1 & \vec{l} \\ \vec{0}^\top & 1 \end{bmatrix};$ 
4  $\mathbf{L} \leftarrow \begin{bmatrix} \mathbf{L}_1 & \vec{0} \\ \vec{c}^\top & \alpha_{nn} \end{bmatrix};$ 
```

Algorithm 2: An iterative algorithm to factorize $\mathbf{A} = \mathbf{UL}$ where \mathbf{U} is unit upper triangular and \mathbf{L} is lower triangular.

Data: \mathbf{A} - an $n \times n$ matrix
Result: \mathbf{UL} - a single matrix which stores the entries of \mathbf{U} above the diagonal and the entries of \mathbf{L} below the diagonal.

```

1  $\mathbf{UL} \leftarrow \mathbf{A};$ 
2 for  $i \leftarrow n$  to 2 do
3    $\mathbf{UL}[1:i-1, i] \leftarrow \mathbf{UL}[1:i-1, i]/\mathbf{UL}_{ii};$ 
   /* Note that  $\mathbf{UL}[1:i-1, 1:i-1]$  is unused for the solution so far - we can use
   this to store  $\hat{\mathbf{A}} - \vec{l}\vec{c}^\top$ . */
4    $\mathbf{UL}[1:i-1, 1:i-1] \leftarrow \mathbf{UL}[1:i-1, 1:i-1] - \mathbf{UL}[1:i-1, i] \cdot \mathbf{UL}[i, 1:i-1];$ 
5 end
```

If we pre-allocate space for the entirety of \mathbf{U}, \mathbf{L} , this algorithm is tail-recursive, and can be de-recursed. As well, only the entries above the diagonal of \mathbf{U} are modified, and only the entries on and below the diagonal of \mathbf{L} are modified - so they can be stored in the same matrix. This leads to algorithm 2.

Note that in the case where it is not necessary to preserve \mathbf{A} , we can use $\mathbf{UL} = \mathbf{A}$

- (b) Based on your algorithm, prove/provide the necessary and sufficient determinant conditions for the existence of the UL decomposition.

Answer.

Theorem 1. The following two statements are equivalent:

- i. There exist a unique unit upper triangular matrix \mathbf{U} and nonsingular lower triangular matrix \mathbf{L} such that $\mathbf{A} = \mathbf{UL}$.
- ii. All trailing principal submatrices of \mathbf{A} are nonsingular.

Proof.

\Rightarrow For a trailing submatrix \mathbf{A}_{ii} , we have

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{1i} \\ \mathbf{A}_{i1} & \mathbf{A}_{ii} \end{bmatrix} = \begin{bmatrix} \mathbf{U}_{11} & \mathbf{U}_{1i} \\ 0 & \mathbf{U}_{ii} \end{bmatrix} \begin{bmatrix} \mathbf{L}_{11} & 0 \\ \mathbf{L}_{i1} & \mathbf{L}_{ii} \end{bmatrix} = \mathbf{UL}.$$

Then $\det(\mathbf{A}_{ii}) = \det(\mathbf{U}_{ii}\mathbf{L}_{ii}) = \det(\mathbf{U}_{ii})\det(\mathbf{L}_{ii}) = 1 \cdot \prod_{k=1}^i (\mathbf{L}_{ii})_{kk} \neq 0$, since \mathbf{U} is unit upper triangular and \mathbf{L} is nonsingular lower triangular.

\Leftarrow Note that this is trivially true for a 1×1 matrix. Then, by induction, assume there exist a unique unit upper triangular matrix \mathbf{U} and nonsingular lower triangular matrix \mathbf{L} such that $\mathbf{A} = \mathbf{UL}$ for all $(n-1) \times (n-1)$ matrices with all trailing principal submatrices nonsingular. For

$$\mathbf{A} = \begin{bmatrix} \hat{\mathbf{A}} & \vec{a} \\ \vec{c}^\top & \alpha_{nn} \end{bmatrix},$$

and

$$\hat{\mathbf{A}} = \mathbf{U}_1\mathbf{L}_1,$$

which exist by assumption, we have

$$\mathbf{A} = \begin{bmatrix} \mathbf{U}_1 & \vec{l} \\ \vec{0}^\top & 1 \end{bmatrix} \begin{bmatrix} \mathbf{L}_1 & \vec{0} \\ \vec{c}^\top & \alpha_{nn} \end{bmatrix},$$

where $\vec{l} = \vec{a}/\alpha_{nn}$. Note, then that $0 \neq \det(\mathbf{A}) = \det(\mathbf{L}_1)\alpha_{nn}$, so $\alpha_{nn} \neq 0$, and the lower triangular matrix is nonsingular.



(c) Write a Matlab code to implement the UL decomposition and apply it to

$$\begin{bmatrix} 1 & 0 & 2 & 1 \\ -4 & 5 & 3 & -1 \\ -1 & 3 & 1 & 1 \\ 0 & 2 & 0 & 1 \end{bmatrix}$$

to verify that your code generates the required decomposition $\mathbf{A} = \mathbf{UL}$.

Answer. The output for this matrix can be found below.

U=

$$\begin{bmatrix} 1 & -1 & 2 & 1 \\ 0 & 1 & 3 & -1 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

L=

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

UL=

$$\begin{bmatrix} 1 & 0 & 2 & 1 \\ -4 & 5 & 3 & -1 \\ -1 & 3 & 1 & 1 \\ 0 & 2 & 0 & 1 \end{bmatrix}$$

The code used to implement the algorithm and verify the given test matrix can be found below.

```

1  use std::{
2      error::Error,
3      fs::{create_dir_all, File},
4      io::Write,
5  };
6
7  use nalgebra::{Dim, Matrix4, OMatrix};
8
9  fn main() -> Result<(), Box<dyn Error>> {
10     // fill A with given values, column major order
11     let a = Matrix4::from_iterator([

```

```

12     1.0, -4.0, -1.0, 0.0, 0.0, 5.0, 3.0, 2.0, 2.0, 3.0, 1.0, 0.0, 1.0, -1.0, 1.0, 1.0,
13 ];
14
15 // compute UL decomposition
16 let ul = ul_decomp(a);
17
18 // extract the upper triangular part
19 let mut upper = ul.upper_triangle();
20 // enforce upper is unit triangular
21 upper.fill_diagonal(1.0);
22 // extract lower triangular part
23 let lower = ul.lower_triangle();
24
25 create_dir_all("./out/")?;
26 let f = File::create("out/output5.txt")?;
27 writeln!(&f, "U={upper}\nL={lower}\nUL={}", upper * lower)?;
28
29 Ok(()) // :)
30 }
31
32 /// Computes UL decomposition in place, note input matrix is overridden
33 fn ul_decomp<R: Dim, C: Dim>(mut ul: OMatrix<f64, R, C>) -> OMatrix<f64, R, C>
34 where
35     nalgebra::DefaultAllocator: nalgebra::allocator::Allocator<f64, R, C>,
36 {
37     // rust is 0 indexed
38     for i in (1..4).rev() {
39         let alphaii = ul[(i, i)];
40         let mut a_vec = ul.view_mut((0, i), (i, 1));
41         a_vec /= alphaii;
42
43         let lct = ul.view((0, i), (i, 1)) * ul.view((i, 0), (1, i));
44         // start (row, col), size (nrows, ncols)
45         let mut a_hat = ul.view_mut((0, 0), (i, i));
46         a_hat -= lct;
47     }
48     ul
49 }

```

6. Even though we rarely need to compute the inverse of a matrix, let us think about it in this problem. Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be an invertible matrix. Describe an algorithm (based on the LU Decomposition/Gaussian Elimination) that computes \mathbf{A}^{-1} with an operation count of $8n^3/3$ flops (ignoring the lower order terms).

Answer. To compute the inverse of a matrix from the LU decomposition, we must solve two triangular systems n times. This is to find the n columns of \mathbf{A}^{-1} from the LU decomposition. Thus, the cost to perform this backward substitution is $2n$ times the cost to solve a single triangular system, which is n^2 . Summing this together, the cost of finding \mathbf{A}^{-1} is $\frac{2}{3}n^3 + 2n(n^2) = \frac{8}{3}n^3 + \mathcal{O}(n^2)$. The algorithm is described in more detail in algorithm 3.

Algorithm 3: An algorithm to compute the inverse of a matrix \mathbf{A} using the LU decomposition.

Data: \mathbf{A} and $n \times n$ matrix

Result: \mathbf{A}^{-1}

1 Factorize $\mathbf{A} = \mathbf{LU}$;	$// \frac{2}{3} n^3$
2 for $i \leftarrow 1$ to n do	
3 Solve $\mathbf{LU}\vec{b}_i = \vec{e}_i$;	$// \sum_{i=1}^n (2 n^2)$
4 end	
5 $\mathbf{A}^{-1} \leftarrow [\vec{b}_1 \ \vec{b}_2 \ \dots \ \vec{b}_n]$;	

7. Inspired by the presentation in [Trefethen and Bau, SIAM Press, 1997], in this problem, we will numerically investigate the growth factor in LU with partial pivoting. In the class (see October 13th notes), we showed that the growth factor ρ_{pp} could be as large as $\rho_{pp} = 2n - 1$. Indeed we had found an example where this upper bound is attained. However, as we mentioned, the algorithm behaves much better in practice. Here, we will try LU with partial pivoting on random matrices with varying dimensions and plot the observations.

Use the command `n = ceil(logspace(1, 3, 1000))` to create (approximately) logarithmically spaced matrix dimensions between 10 and 1000. Some of the dimensions will be repeated. The variable `n` is a vector of size 1000 with entries ranging from 10 to 1000. Then, for every entry `n(i)` of `n`, i.e., for $i = 1, 2, \dots, 1000$, create a random matrix `A` using `A = randn(n(i), n(i))/sqrt(n(i))`. So, we are creating a random matrix of varying dimensions with normally distributed entries having mean zero and standard deviation $\sqrt{n(i)}$. Then, compute the growth factor ρ_{pp} for every `A`. At the end you will have a vector of size 1000 whose entries corresponding to the growth factor for every random `A`. Using the `loglog` command (logarithmic scale both in the horizontal and 2 vertical axes), plot the growth factor vs the matrix dimensions `n`. On the same plot, by using the `hold on` command, plot the growth rate of \sqrt{n} . How is the observed/numerical growth behaving with respect to the theoretical upper bound 2^{n-1} and with respect to \sqrt{n} ? Comment on your observations.

Answer. The growth rate of the calculated ρ_{pp} v.s. \sqrt{n} can be found in fig. 1. Note that \sqrt{n} has been scaled to more easily observe the difference in growth rates. It can be seen that ρ_{pp} grows slightly faster than \sqrt{n} on average, as expected (the average growth rate should be somewhere between \sqrt{n} and $n^{3/4}$).

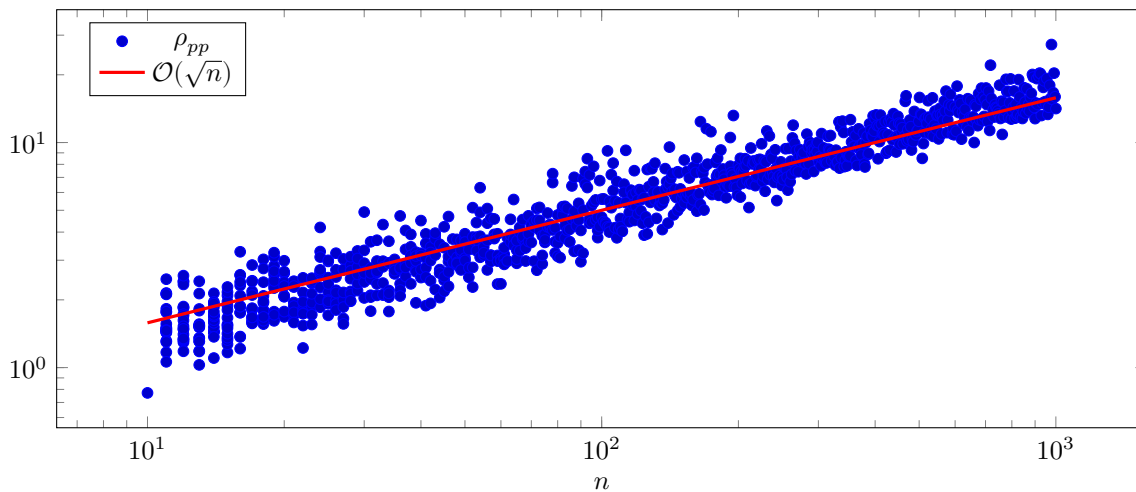


Figure 1: Growth rate of ρ_{pp} vs. \sqrt{n} .

The code used for this problem can be found below.

```

1  use indicatif::ProgressIterator;
2  use nalgebra::{DMatrix, LU};
3  use rand_distr::StandardNormal;
4  use std::{
5      error::Error,
6      fs::{create_dir_all, File},
7      io::Write,
8  };
9
10 fn main() -> Result<(), Box<dyn Error>> {
11     create_dir_all("./out/")?;
12     let f = File::create("out/output7.txt")?;
13     for n in ceil_logspace(1., 3., 1000).progress() {
14         // A = randn(n, n)/sqrt(n)
15         let a = DMatrix::from_distribution(n, n, &StandardNormal, &mut rand::thread_rng())
16             / (n as f64).sqrt();
17
18         // ||A||_max = max_{i,j} |a_{ij}|
19         let a_max = a.amax();
20         let u_max = LU::new(a).u().amax();
21
22         // rho_pp = ||U||_max / ||A||_max
23         let rho_pp = u_max / a_max;
24
25         writeln!(&f, "{n} {rho_pp}")?;
26     }
27 }

```

```
27     Ok(()) // :)
28 }
29
30
31 /// Recreates Matlab ceil(logspace(a, b, n)) generates n points between decades 10^a and 10^b.
32 fn ceil_logspace(a: f64, b: f64, n: usize) -> impl ExactSizeIterator<Item = usize> {
33     let temp = (b - a) / (n as f64 - 1.);
34     (0..n).map(move |i| 10_f64.powf((i as f64) * temp + a).ceil() as usize)
35 }
```