

# CSE 6643 Homework 3

Karl Hiner

## 1 One-upping [25 pts]

Let  $\mathbf{A} \in \mathbb{R}^{m \times m}$  have full rank. Assume that we have already computed the QR decomposition of  $\mathbf{A}$ . For  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^m$ , we call the matrix  $\mathbf{B} = \mathbf{A} + \mathbf{u}\mathbf{v}^T$  a rank-1 update of  $\mathbf{A}$ .

### (a) [5 pts]

Prove that if  $\mathbf{v}^T \mathbf{A}^{-1} \mathbf{u} \neq -1$ , then  $\mathbf{B}$  is invertible.

Let  $\mathbf{x} \in \mathbb{R}^m$  be a non-zero vector. Then,

$$\begin{aligned} \mathbf{v}^T \mathbf{A}^{-1} \mathbf{u} &\neq -1 \\ 1 &\neq -\mathbf{v}^T \mathbf{A}^{-1} \mathbf{u} \\ \mathbf{v}^T \mathbf{x} &\neq -(\mathbf{v}^T \mathbf{A}^{-1} \mathbf{u})(\mathbf{v}^T \mathbf{x}) \\ \mathbf{x} &\neq -\mathbf{A}^{-1} \mathbf{u}(\mathbf{v}^T \mathbf{x}) \\ \mathbf{A} \mathbf{x} &\neq -\mathbf{u} \mathbf{v}^T \mathbf{x} \\ (\mathbf{A} + \mathbf{u} \mathbf{v}^T) \mathbf{x} &\neq \mathbf{0}. \\ \mathbf{B} \mathbf{x} &\neq \mathbf{0} \end{aligned}$$

Thus, if  $\mathbf{v}^T \mathbf{A}^{-1} \mathbf{u} \neq -1$ , there are no nontrivial solutions for  $\mathbf{B} \mathbf{x} = \mathbf{0}$ , and so  $\mathbf{B}$  is invertible.

### (b) [10 pts]

Design an algorithm that provably solves the system of equations  $\mathbf{B} \mathbf{x} = \mathbf{b}$  in  $O(m^2)$  operations.

$$\begin{aligned} \mathbf{x} &= \mathbf{B}^{-1} \mathbf{b} && \text{(assume } \mathbf{B} \text{ is invertible)} \\ &= (\mathbf{A} + \mathbf{u} \mathbf{v}^T)^{-1} \mathbf{b} && \text{(since } \mathbf{B} = \mathbf{A} + \mathbf{u} \mathbf{v}^T) \\ &= \left( \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1} \mathbf{u} \mathbf{v}^T \mathbf{A}^{-1}}{1 + \mathbf{v}^T \mathbf{A}^{-1} \mathbf{u}} \right) \mathbf{b} && \text{(Sherman-Morrison formula [1])} \\ &= \mathbf{A}^{-1} \mathbf{b} - \frac{\mathbf{A}^{-1} \mathbf{u} \mathbf{v}^T}{1 + \mathbf{v}^T \mathbf{A}^{-1} \mathbf{u}} \mathbf{A}^{-1} \mathbf{b} && \text{(note: (1a) above), } \mathbf{B} \text{ invertible} \implies \text{den} \neq 0 \\ &= \tilde{\mathbf{x}} - \frac{\tilde{\mathbf{u}} \mathbf{v}^T}{1 + \mathbf{v}^T \tilde{\mathbf{u}}} \tilde{\mathbf{x}} && \text{(let } \tilde{\mathbf{x}} \equiv \mathbf{A}^{-1} \mathbf{b}, \tilde{\mathbf{u}} \equiv \mathbf{A}^{-1} \mathbf{u}) \\ &= \left[ \mathbf{I} + \left( \frac{-1}{1 + \mathbf{v}^T \tilde{\mathbf{u}}} \right) \tilde{\mathbf{u}} \mathbf{v}^T \right] \tilde{\mathbf{x}} && \text{(rearrange)} \\ &= (\mathbf{I} + \alpha \tilde{\mathbf{u}} \mathbf{v}^T) \tilde{\mathbf{x}} && \text{(let } \alpha \equiv \frac{-1}{1 + \mathbf{v}^T \tilde{\mathbf{u}}}) \end{aligned}$$

Thus, we have derived an expression for  $\mathbf{x}$  in terms of a scalar  $\alpha$  and vectors  $\tilde{\mathbf{u}}, \tilde{\mathbf{x}}$ .

Vectors  $\tilde{\mathbf{u}}$  and  $\tilde{\mathbf{x}}$  are both defined in terms of  $\mathbf{A}^{-1} = \mathbf{R}^{-1}\mathbf{Q}^T$ , where we assume  $\mathbf{R}$  and  $\mathbf{Q}$  have already been computed. Thus, both vectors can be computed using back substitution in  $O(m^2)$ .  $\alpha$  can then be computed in  $O(m)$ , and the final expression for  $\mathbf{x}$  involves a vector-scalar product, a vector outer product, an identity addition, and a matrix-vector multiplication, for a total of

$$O(m) + O(m^2) + O(m) + O(m^2) = O(2m) + O(m^2) = O(m^2)$$

operations.

Here is the algorithm:

1. Solve  $\mathbf{R}\tilde{\mathbf{x}} = \mathbf{Q}^T\mathbf{b}$  for  $\tilde{\mathbf{x}}$  using back substitution. ( $O(m^2)$ )
2. Solve  $\mathbf{R}\tilde{\mathbf{u}} = \mathbf{Q}^T\mathbf{u}$  for  $\tilde{\mathbf{u}}$  using back substitution. ( $O(m^2)$ )
3. Compute the scalar  $\alpha = \frac{-1}{1 + \mathbf{v}^T\tilde{\mathbf{u}}}$ . ( $O(m)$ )
4. Finally, compute the solution  $\mathbf{x} = (\mathbf{I} + ((\alpha\tilde{\mathbf{u}})\mathbf{v}^T))\tilde{\mathbf{x}}$ . ( $O(m^2)$ )

Since the highest-order term across all steps is  $O(m^2)$ , the total number of operations is  $O(m^2)$ .

### (c) r-upping [10 pts]

Extend the algorithm from the previous exercise to the case of  $\mathbf{B} = \mathbf{A} + \mathbf{U}\mathbf{V}^T$ , for  $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{m \times r}$  and  $r \ll m$ . Calculate the asymptotic complexity of the resulting algorithm.

The Sherman-Morrison formula was the key step above. Here, we will use the generalization of that formula, the Woodbury matrix identity [2]. Given the definitions above, the Woodbury matrix identity states that if  $(\mathbf{I} + \mathbf{V}^T\mathbf{A}^{-1}\mathbf{U})$  is invertible,

$$\mathbf{B}^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{I} + \mathbf{V}^T\mathbf{A}^{-1}\mathbf{U})^{-1}\mathbf{V}^T\mathbf{A}^{-1}.$$

Thus, we can express the solution  $\mathbf{x}$  as

$$\begin{aligned} \mathbf{x} &= \left( \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{I} + \mathbf{V}^T\mathbf{A}^{-1}\mathbf{U})^{-1}\mathbf{V}^T\mathbf{A}^{-1} \right) \mathbf{b} \\ &= \mathbf{A}^{-1}\mathbf{b} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{I} + \mathbf{V}^T\mathbf{A}^{-1}\mathbf{U})^{-1}\mathbf{V}^T\mathbf{A}^{-1}\mathbf{b} \\ &= \tilde{\mathbf{x}} - \tilde{\mathbf{U}}(\mathbf{I} + \mathbf{V}^T\tilde{\mathbf{U}})^{-1}\mathbf{V}^T\tilde{\mathbf{x}} && (\text{let } \tilde{\mathbf{x}} \equiv \mathbf{A}^{-1}\mathbf{b}, \tilde{\mathbf{U}} \equiv \mathbf{A}^{-1}\mathbf{U}) \\ &= \tilde{\mathbf{x}} - \tilde{\mathbf{U}}\mathbf{C}^{-1}\mathbf{V}^T\tilde{\mathbf{x}} && (\text{let } \mathbf{C} = \mathbf{I} + \mathbf{V}^T\tilde{\mathbf{U}}) \end{aligned}$$

Note that, since  $\tilde{\mathbf{U}} \in \mathbb{R}^{m \times r}$ , and  $\mathbf{V}^T \in \mathbb{R}^{r \times m}$ , we can compute the matrix product  $\tilde{\mathbf{U}}\mathbf{V}^T$  using  $O(m^2r)$  operations.

Here is the algorithm:

1. Solve  $\mathbf{R}\tilde{\mathbf{x}} = \mathbf{Q}^T\mathbf{b}$  for  $\tilde{\mathbf{x}}$  using back substitution (same as step 1 in (b)). ( $O(m^2)$ )
2. Compute  $\tilde{\mathbf{U}} = \mathbf{A}^{-1}\mathbf{U}$  by solving each of the linear systems  $\mathbf{R}\tilde{\mathbf{u}}_i = \mathbf{Q}^T\mathbf{u}_i$  for  $\tilde{\mathbf{u}}_i$ , where  $\tilde{\mathbf{u}}_i$  is the  $i$ th column of  $\tilde{\mathbf{U}}$  and  $\mathbf{u}_i$  is the  $i$ th column of  $\mathbf{U}$ . This can be done using back substitution (as in step 2 in (b)) for each of the  $r$  columns. ( $O(m^2r)$ )
3. *I was unable to complete this.* I feel that this is very close. Finding  $\mathbf{C}^{-1}$  should be possible using another linear system, similarly in  $O(m^2r)$  time, and I think it can be done in a way that actually avoids an expensive final computation step by the correct formulation of this linear system (ideally using only  $O(mr)$  matrix-vector/vector-vector multiplications).

This solution should be computable in  $O(m^2r)$  time.

## 2 You Factor [25 pts]

In class we have seen that if  $\tilde{\mathbf{x}} \in \mathbb{R}^m$  is the solution to the system  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , as computed by unpivoted LU factorization, we have

$$(\mathbf{A} + \mathbf{E})\tilde{\mathbf{x}} = \mathbf{b}. \quad (1)$$

For  $u$  the unit roundoff error and  $\tilde{\mathbf{L}}, \tilde{\mathbf{U}}$  the LU factors computed in finite precision, we have

$$|\mathbf{E}| \leq mu(2|\mathbf{A}| + 4|\tilde{\mathbf{L}}||\tilde{\mathbf{U}}|) + O(u^2). \quad (2)$$

Here,  $|\cdot|$  signifies the element-wise absolute values and  $\leq$  is interpreted element-wise, as well. In this problem, we investigate the conclusions from this bound in the case of row-pivoted LU factorization.

### (a) [7.5 pts]

Deduce that under row-pivoted LU factorization and taking  $\|\cdot\|_\infty$  to signify the vector-infinity norm, we have

$$\|\mathbf{E}\|_\infty \leq mu(2\|\mathbf{A}\|_\infty + 4m\|\tilde{\mathbf{U}}\|_\infty) + O(u^2). \quad (3)$$

This prompts us to investigate the growth factor  $\rho := \frac{\|\mathbf{U}\|_\infty}{\|\mathbf{A}\|_\infty}$  of row-pivoted LU factorization.

During row-pivoted LU factorization, we pick the maximum value over a column during pivot selection. Thus, every element of  $\tilde{\mathbf{L}}$  has magnitude  $\leq 1$ , by construction. This leads to the following inequality:

$$(|\tilde{\mathbf{L}}||\tilde{\mathbf{U}}|)_{ij} = \sum_{k=1}^m |\tilde{L}_{ik}||\tilde{U}_{kj}| \leq m|\tilde{U}|_{ij} \quad (4)$$

Thus, we have

$$\begin{aligned} |\mathbf{E}| &\leq mu(2|\mathbf{A}| + 4|\tilde{\mathbf{L}}||\tilde{\mathbf{U}}|) + O(u^2) \\ &\leq mu(2|\mathbf{A}| + 4m|\tilde{\mathbf{U}}|) + O(u^2). \end{aligned}$$

Since all operators here apply element-wise, the desired result then follows immediately, by the definition of the vector-infinity norm applied to matrices, which applies an element-wise maximum.

### (b) [5 pts]

Verify that the rows  $\mathbf{u}_i^T, \mathbf{a}_i^T$  of  $\mathbf{U}, \mathbf{A}$  satisfy

$$\mathbf{u}_i^T = \mathbf{a}_i^T - \sum_{j=1}^{i-1} \mathbf{L}_{ij} \mathbf{u}_j^T. \quad (5)$$

Expressing the LU factorization of  $\mathbf{A}$  in terms of rows and columns:

$$\mathbf{A} = \mathbf{LU}$$

$$\begin{bmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \dots \\ \mathbf{a}_m^T \end{bmatrix} = \mathbf{L} \begin{bmatrix} \mathbf{u}_1^T \\ \mathbf{u}_2^T \\ \dots \\ \mathbf{u}_m^T \end{bmatrix}$$

By the definition of matrix multiplication, with  $\mathbf{X}_i$  denoting the  $i$ th row of matrix  $\mathbf{X}$ ,

$$\mathbf{A}_{ij} = \mathbf{L}_i \cdot (\mathbf{U}^T)_j = \sum_{k=1}^m \mathbf{L}_{ik} \mathbf{U}_{kj}.$$

Since  $\mathbf{L}$  is lower-triangular and  $\mathbf{U}$  is upper-triangular,  $\mathbf{L}_{ij} = \mathbf{U}_{ji} = 0$  for  $i < j$ . Thus, we have

$$\begin{aligned} \mathbf{A}_{ij} &= \sum_{k=1}^i \mathbf{L}_{ik} \mathbf{U}_{kj} \\ \mathbf{A}_{ij} &= \sum_{k=1}^{i-1} \mathbf{L}_{ik} \mathbf{U}_{kj} + \mathbf{L}_{ii} \mathbf{U}_{ij} \\ &= \sum_{k=1}^{i-1} \mathbf{L}_{ik} \mathbf{U}_{kj} + \mathbf{U}_{ij} && \text{(since } \mathbf{L}_{ii} = 1) \\ \mathbf{U}_{ij} &= \mathbf{A}_{ij} - \sum_{k=1}^{i-1} \mathbf{L}_{ik} \mathbf{U}_{kj} && \text{(rearranging terms)} \\ \mathbf{u}_i^T &= \mathbf{a}_i^T - \sum_{k=1}^{i-1} \mathbf{L}_{ik} \mathbf{u}_k^T && \text{(restate in terms of rows. QED)} \end{aligned}$$

**(c) [5 pts]**

Use part (b) to show that  $\|\mathbf{U}\|_\infty \leq 2^{m-1} \|\mathbf{A}\|_\infty$ .

Note: I got some guidance from fellow student Samuel Talkington about this problem.

First, we can use the result from (b) to find the max value of the first row of  $\mathbf{U}$ :

$$\begin{aligned} \|\mathbf{u}_1^T\|_\infty &= \left\| \mathbf{a}_1^T - \sum_{j=1}^0 \mathbf{L}_{1j} \mathbf{u}_j^T \right\|_\infty \\ &= \|\mathbf{a}_1^T\|_\infty. \end{aligned}$$

We can use this result to find a bound on the second row:

$$\begin{aligned} \|\mathbf{u}_2^T\|_\infty &= \left\| \mathbf{a}_2^T - \sum_{j=1}^1 \mathbf{L}_{2j} \mathbf{u}_j^T \right\|_\infty && \text{(from (b))} \\ &\leq \|\mathbf{a}_2^T\|_\infty + \left\| \sum_{j=1}^1 \mathbf{L}_{2j} \mathbf{u}_j^T \right\|_\infty && \text{(triangle inequality, def. of } \infty\text{-norm)} \\ &= \|\mathbf{a}_2^T\|_\infty + \|\mathbf{L}_{21} \mathbf{u}_1^T\|_\infty && \text{(only one element in sum)} \\ &\leq \|\mathbf{a}_2^T\|_\infty + \|\mathbf{u}_1^T\|_\infty && \text{(since } \max |\mathbf{L}_{ij}| \leq 1) \\ &= \|\mathbf{a}_2^T\|_\infty + \|\mathbf{a}_1^T\|_\infty && \text{(from above result for row 1)} \end{aligned}$$

Continuing to the third row:

$$\begin{aligned}
\|\mathbf{u}_3^T\|_\infty &\leq \|\mathbf{a}_3^T\|_\infty + \left\| \sum_{j=1}^2 \mathbf{L}_{3j} \mathbf{u}_j^T \right\|_\infty \\
&\leq \|\mathbf{a}_3^T\|_\infty + \|\mathbf{L}_{32} \mathbf{u}_2^T\|_\infty + \|\mathbf{L}_{31} \mathbf{u}_1^T\|_\infty && \text{(expand sum, triangle inequality)} \\
&\leq \|\mathbf{a}_3^T\|_\infty + \|\mathbf{u}_2^T\|_\infty + \|\mathbf{u}_1^T\|_\infty && \text{(since } \max |\mathbf{L}_{ij}| \leq 1) \\
&= \|\mathbf{a}_3^T\|_\infty + \|\mathbf{a}_2^T\|_\infty + 2\|\mathbf{a}_1^T\|_\infty && \text{(from above results for rows 1/2)}
\end{aligned}$$

and the fourth:

$$\begin{aligned}
\|\mathbf{u}_4^T\|_\infty &\leq \|\mathbf{a}_4^T\|_\infty + \|\mathbf{u}_3^T\|_\infty + \|\mathbf{u}_2^T\|_\infty + \|\mathbf{u}_1^T\|_\infty \\
&= \|\mathbf{a}_4^T\|_\infty + \|\mathbf{a}_3^T\|_\infty + 2\|\mathbf{a}_2^T\|_\infty + 4\|\mathbf{a}_1^T\|_\infty
\end{aligned}$$

We can see the pattern emerging:

$$\begin{aligned}
\|\mathbf{u}_i^T\|_\infty &\leq \|\mathbf{a}_i^T\|_\infty + \sum_{j=1}^{i-1} \|\mathbf{u}_j^T\|_\infty \\
&= \|\mathbf{a}_i^T\|_\infty + \sum_{j=1}^{i-1} 2^{(j-1)} \|\mathbf{a}_{i-j}^T\|_\infty
\end{aligned}$$

Thus, the bound for the  $m$ th row is

$$\begin{aligned}
\|\mathbf{u}_m^T\|_\infty &\leq \|\mathbf{a}_m^T\|_\infty + \sum_{j=1}^{m-1} 2^{(j-1)} \|\mathbf{a}_{m-j}^T\|_\infty \\
&\leq 2^{m-1} \left( \max_i \|\mathbf{a}_i^T\|_\infty \right) \\
&= 2^{m-1} \|\mathbf{A}\|_\infty,
\end{aligned}$$

where the last equality is from the definition of the vector  $\infty$ -norm applied to matrices. This last row,  $\mathbf{u}_m^T$ , has the largest  $\infty$ -norm of all rows, and thus will contain the largest element in the matrix, and so we have shown that  $\|\mathbf{U}\|_\infty \leq 2^{m-1} \|\mathbf{A}\|_\infty$ .

**(d) [7.5 pts]**

Consider matrices of the form

$$\mathbf{A} = \begin{pmatrix} 1 & & & & 1 \\ -1 & 1 & & & 1 \\ -1 & -1 & 1 & & 1 \\ -1 & -1 & -1 & 1 & 1 \\ -1 & -1 & -1 & -1 & 1 \end{pmatrix}. \tag{6}$$

Derive the growth factor in this case as a function of  $m$ . How does this relate to part (c)?

The  $\mathbf{PA} = \mathbf{LU}$  factorization for  $\mathbf{A}$  is:

$$\begin{pmatrix} 1 & & & & 1 \\ -1 & 1 & & & 1 \\ -1 & -1 & 1 & & 1 \\ -1 & -1 & -1 & 1 & 1 \\ -1 & -1 & -1 & -1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & & & & \\ -1 & 1 & & & \\ -1 & -1 & 1 & & \\ -1 & -1 & -1 & 1 & \\ -1 & -1 & -1 & -1 & 1 \end{pmatrix} \begin{pmatrix} 1 & & & & 1 \\ & 1 & & & 2 \\ & & 1 & & 4 \\ & & & 1 & 8 \\ & & & & 16 \end{pmatrix} \quad (7)$$

The growth factor is then

$$\rho = \frac{\|\mathbf{U}\|_\infty}{\|\mathbf{A}\|_\infty} = \frac{16}{1} = 16 = 2^4 = 2^{m-1}.$$

In part (c), we showed that an upper bound for the maximum value of  $\mathbf{U}$  can be related to the matrix dimension and the maximum value of a matrix  $\mathbf{A}$ , with

$$\|\mathbf{U}\|_\infty \leq 2^{m-1} \|\mathbf{A}\|_\infty.$$

Rearranging terms, we can see that this also provides an upper bound for the growth factor of an  $m \times m$  matrix  $\mathbf{A}$ :

$$\frac{\|\mathbf{U}\|_\infty}{\|\mathbf{A}\|_\infty} \leq 2^{m-1}.$$

Since  $\rho = 2^{m-1}$  for  $\mathbf{A}$ , we can conclude that  $\mathbf{A}$  exemplifies a worst-case growth factor for a square matrix of its size.

### 3 [25 pts]

Suppose that  $\mathbf{A} \in \mathbb{K}^{m \times m}$  is strictly column diagonally dominant, meaning that for all  $1 \leq k \leq m$ ,

$$|\mathbf{A}_{kk}| > \sum_{j \neq k} |\mathbf{A}_{jk}|. \quad (8)$$

Show that if LU factorization with row pivoting is applied to  $\mathbf{A}$ , no row interchange takes place.

To show that no row interchange takes place, we can show the following:

1. Show that  $\mathbf{A}$ 's strict column diagonal property implies no row interchange takes place during the first factorization step. Then the first permutation matrix  $\mathbf{P}^1 = \mathbf{I}$ , and we have  $\mathbf{L}^1 \mathbf{P}^1 \mathbf{A} = \mathbf{L}^1 \mathbf{A} = \mathbf{U}^1$ , with  $\mathbf{U}^1$  being the first intermediate matrix on the way to  $\mathbf{U}$  after completing factorization step  $k = 1$ . Express  $\mathbf{U}^1$  as

$$\mathbf{U}^1 = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{1k>1} \\ \mathbf{0} & \mathbf{A}^1 \end{bmatrix},$$

with  $\mathbf{A}^1 = \mathbf{A}_{2:m,2:m} \in \mathbb{K}^{m-1 \times m-1}$ .

2. Given (1), show that  $\mathbf{A}^1$  is also strictly diagonally dominant. Then, apply LU factorization with row pivoting to  $\mathbf{A}^1$ , with  $\hat{\mathbf{P}} \mathbf{A}^1 = \hat{\mathbf{L}} \hat{\mathbf{U}}$ . We conclude from (1) that  $\hat{\mathbf{P}}^1 = \mathbf{I}$ , and we have  $\hat{\mathbf{L}}^1 \mathbf{A}^1 = \hat{\mathbf{U}}^1$ , where

$$\hat{\mathbf{U}}^1 = \begin{bmatrix} \mathbf{A}_{11}^1 & \mathbf{A}_{1,2:m}^1 \\ \mathbf{0} & \mathbf{A}^2 \end{bmatrix},$$

with  $\mathbf{A}^2 = \mathbf{A}_{2:m,2:m}^1 \in \mathbb{K}^{m-2 \times m-2}$ .

If we show (1) and (2) hold, then we could apply step (1) again to matrix  $\mathbf{A}^2$  to show that no row interchange takes place at step  $k = 2$ , and again apply step (2) to show  $\mathbf{A}^3$  is also strictly diagonally dominant, etc. Thus, if we show that (1) and (2) hold, we can inductively infer that no row interchanges take place during all factorization steps  $1 \leq k \leq m - 1$ .

1. Show that  $|\mathbf{A}_{kk}| > \sum_{j \neq k} |\mathbf{A}_{jk}| \implies \mathbf{P}^1 = \mathbf{I}$ : During the first factorization step, we choose a pivot row  $j_1 = \operatorname{argmax}_j |\mathbf{A}_{j1}|$ .

$$\begin{aligned} |\mathbf{A}_{kk}| > \sum_{j \neq k} |\mathbf{A}_{jk}| &\geq \max_{j \neq k} |\mathbf{A}_{jk}| \implies \\ |\mathbf{A}_{kk}| &= \max_j |\mathbf{A}_{jk}| \implies \\ |\mathbf{A}_{11}| &= \max_j |\mathbf{A}_{j1}| \implies \\ 1 &= \operatorname{argmax}_j |\mathbf{A}_{j1}|, \end{aligned}$$

and so we choose  $j_1 = 1$ . In other words, no rows are swapped during this first step, and  $\mathbf{P}^1 = \mathbf{I}$ .

2. Show that  $\mathbf{A}^1$  is also strictly diagonally dominant.

For brevity, let  $\hat{j} \equiv j + 1$  and  $\hat{k} = k + 1$  for  $j, k \in \mathbb{N}$ , and let  $\alpha_k \equiv \frac{\mathbf{A}_{1k}}{\mathbf{A}_{11}}$ . Then, after the first factorization step, the  $k$ th column of  $\mathbf{A}^1$  can be defined in terms of  $\mathbf{A}$  as:

$$\begin{bmatrix} \mathbf{A}_{1k}^1 \\ \mathbf{A}_{2k}^1 \\ \dots \\ \mathbf{A}_{kk}^1 \\ \dots \\ \mathbf{A}_{(m-1),k}^1 \end{bmatrix} = \begin{bmatrix} \mathbf{A}_{2\hat{k}} - \frac{\mathbf{A}_{21}}{\mathbf{A}_{11}} \mathbf{A}_{1\hat{k}} \\ \mathbf{A}_{3\hat{k}} - \frac{\mathbf{A}_{31}}{\mathbf{A}_{11}} \mathbf{A}_{1\hat{k}} \\ \dots \\ \mathbf{A}_{\hat{k}\hat{k}} - \frac{\mathbf{A}_{\hat{k}1}}{\mathbf{A}_{11}} \mathbf{A}_{1\hat{k}} \\ \dots \\ \mathbf{A}_{m\hat{k}} - \frac{\mathbf{A}_{m1}}{\mathbf{A}_{11}} \mathbf{A}_{1\hat{k}} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_{2\hat{k}} - \alpha_{\hat{k}} \mathbf{A}_{21} \\ \mathbf{A}_{3\hat{k}} - \alpha_{\hat{k}} \mathbf{A}_{31} \\ \dots \\ \mathbf{A}_{\hat{k}\hat{k}} - \alpha_{\hat{k}} \mathbf{A}_{\hat{k}1} \\ \dots \\ \mathbf{A}_{m\hat{k}} - \alpha_{\hat{k}} \mathbf{A}_{m1} \end{bmatrix}$$

So we can express any element of  $\mathbf{A}^1$  in terms of  $\mathbf{A}$  as

$$\mathbf{A}_{jk}^1 = \mathbf{A}_{(j+1),(k+1)} + \frac{\mathbf{A}_{1,(k+1)}}{\mathbf{A}_{11}} \mathbf{A}_{(j+1),1} \equiv \mathbf{A}_{j\hat{k}} - \alpha_{\hat{k}} \mathbf{A}_{j1}.$$

Thus, for  $\mathbf{A}^1$  to be strictly diagonally dominant, we must have

$$\begin{aligned} |\mathbf{A}_{kk}^1| &> \sum_{j \neq k}^{m-1} |\mathbf{A}_{jk}^1| \\ |\mathbf{A}_{\hat{k}\hat{k}} - \alpha_{\hat{k}} \mathbf{A}_{\hat{k}1}| &> \sum_{j \neq k}^{m-1} |\mathbf{A}_{\hat{j}\hat{k}} - \alpha_{\hat{k}} \mathbf{A}_{j1}| \\ |\mathbf{A}_{kk} - \alpha_k \mathbf{A}_{k1}| &> \sum_{j > 1, j \neq k}^m |\mathbf{A}_{jk} - \alpha_k \mathbf{A}_{j1}|, k > 1. \end{aligned}$$

Starting with the LHS, with  $k > 1$ :

$$\begin{aligned} |\mathbf{A}_{kk} - \alpha_k \mathbf{A}_{k1}| &\geq |\mathbf{A}_{kk}| - |\alpha_k \mathbf{A}_{k1}| && (\mathbf{A} \text{ SDD} \rightarrow \mathbf{A}_{kk} \mathbf{A}_{11} \geq \mathbf{A}_{k1} \mathbf{A}_{1k}) \\ &> \sum_{j \neq k} |\mathbf{A}_{jk}| - |\alpha_k \mathbf{A}_{k1}| && (\mathbf{A} \text{ is SDD}) \\ &= \sum_{j > 1, j \neq k} |\mathbf{A}_{jk}| + |\mathbf{A}_{1k}| - |\alpha_k \mathbf{A}_{k1}| && (\text{move } j = 1 \text{ out of sum}) \\ &= \sum_{j > 1, j \neq k} |\mathbf{A}_{jk}| + |\alpha_k \mathbf{A}_{11}| - |\alpha_k \mathbf{A}_{k1}| && (\text{sub for } \mathbf{A}_{1k} \text{ using def. of } \alpha) \\ &> \sum_{j > 1, j \neq k} |\mathbf{A}_{jk}| + \sum_{j > 1} |\alpha_k \mathbf{A}_{j1}| - |\alpha_k \mathbf{A}_{k1}| && (\mathbf{A} \text{ is SDD}) \\ &= \sum_{j > 1, j \neq k} |\mathbf{A}_{jk}| + \sum_{j > 1, j \neq k} |\alpha_k \mathbf{A}_{j1}| && (\text{combine right two terms}) \\ &= \sum_{j > 1, j \neq k} (|\mathbf{A}_{jk}| + |\alpha_k \mathbf{A}_{j1}|) && (\text{combine sums}) \\ &\geq \sum_{j > 1, j \neq k} |\mathbf{A}_{jk} - \alpha_k \mathbf{A}_{j1}| && (\text{triangle inequality. QED}) \end{aligned}$$

## 4 Pivoting [25 pts]

### (a) [5 pts]

Go to section (a) of the file `HW3_your_code.jl` and implement a function that takes in a matrix  $\mathbf{LU} \in \mathbb{K}^{m \times m}$  containing the upper triangular part of  $\mathbf{U}$  as well as the strict lower triangular part of  $\mathbf{L}$ , as well as an array  $\mathbf{P} \in \{1, \dots, m\}^m$  that encodes the permutation matrix  $\mathbf{P}$  by  $\mathbf{P}[j] = i \Leftrightarrow \mathbf{P}_{ij} = 1$ . Your function should not allocate any memory.

### (b) [5 pts]

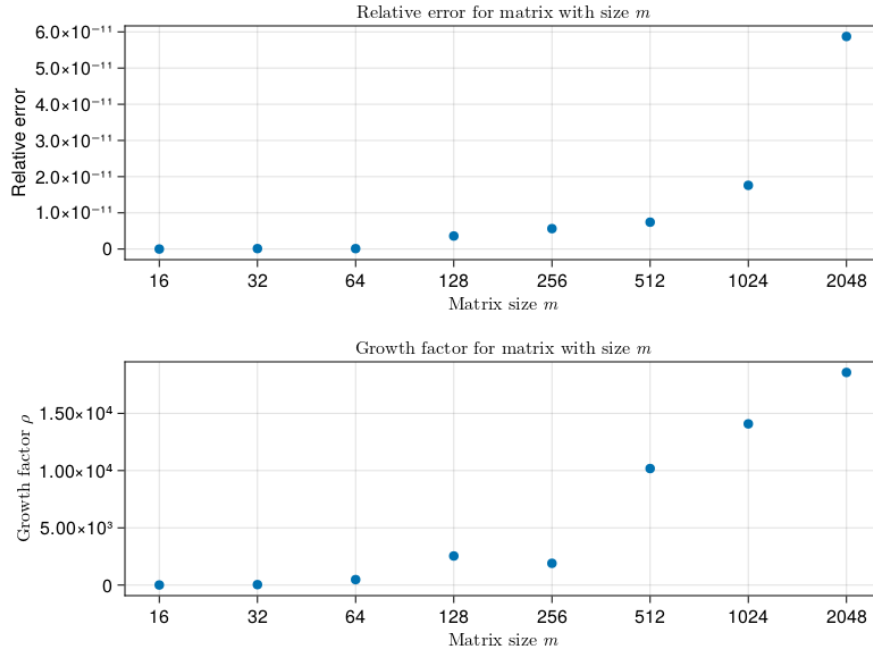
Go to section (b) of the file `HW3_your_code.jl` and implement the unpivoted LU factorization. Check your code by ensuring that the assertions in section a + b of `HW3_driver.jl` do not produce any errors.

### (c) [5 pts]

Generate families of random  $m \times m$  matrices and vectors of length  $m$ . Plot as a function of the size  $m$ , the relative error of the solution obtained from your code in parts (a,b) and the growth



factor introduced in problem 2. Report the floating point type used by your program. You can use the code provided in the second homework as a starting point for creating and saving plots.



**Figure 1** Relative errors and growth factors, without pivoting

See Fig. 1. The floating point type used by my program is `Float64`. All tested matrices are normally distributed with a mean of 0 and a standard deviation of 1, with a unit diagonal offset (to help unpivoted factorization not get stuck), and no row permutations. Relative error was computed as

$$\frac{\|\tilde{\mathbf{b}} - \mathbf{b}\|_2}{\|\mathbf{b}\|_2}.$$

**(d) [5 pts]**

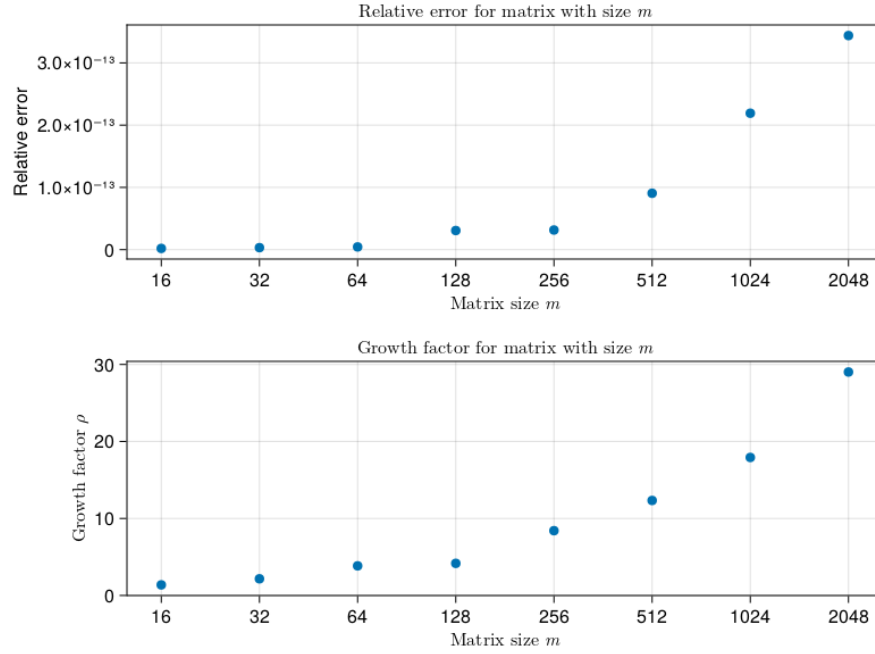
Go to section (d) of the file `HW3_your_code.jl` and implement the unpivoted LU factorization. Your code should pass the assertions in section (c) of `HW3_driver.jl`. Your function should take the matrix  $A$  as an input to modify in place, and return an integer array  $P$  according to the specifications of (a). Repeat the experiment of (c) using the pivoted LU factorization.

See Fig. 2. Note that, as expected, row pivoting results in a dramatic reduction in the growth factor, bringing it down from a max of  $\approx 2 \cdot 10^4$  to a max of  $\approx 30$ ! Also, the relative errors are consistently reduced by 2-4 decimal places (varying across experiments).

*Note that the two experiments were run with different random matrices, so we can't make an exact 1:1 comparison.*

**(e) [5 pts]**

Go to section (e) of the file `HW3_driver.jl` and implement a function that takes an integer  $m$  as an input and returns an  $m \times m$  matrix as introduced in problem 2 (d). Plot the error of the



**Figure 2** Relative errors and growth factors, with pivoting

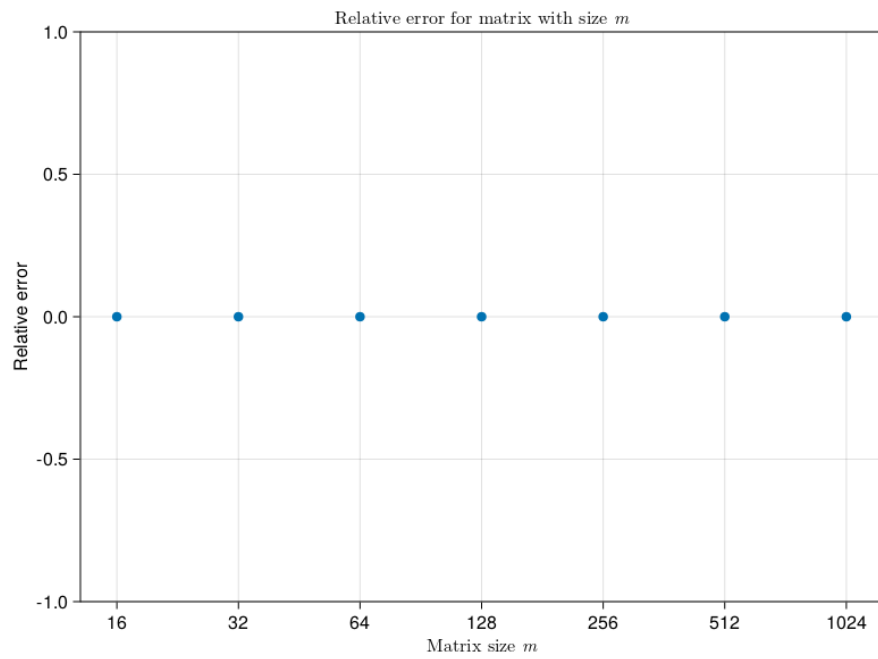
solution when solving equations in this matrix as a function of  $m$ . Compare the error to the built-in solution (the `\` operator). Draw your conclusions from this comparison.

See Fig. 3. For this special kind of maximum-growth matrix, we can see that there is *no difference* between our solution and the built-in `\` operator.

When I print out the resulting  $LU$  factorization, the results look exactly accurate, so I have convinced myself this is not an error. I suppose this results in part from all numbers in the input matrix and the  $LU$  matrices being represented exactly by floating point representations of natural numbers. Beyond this, I am not sure, and I look forward to learning more about this phenomenon! (I have a feeling there's a simple explanation I'm missing...)

## References

- [1] Jack Sherman and Winifried J. Morrison. Adjustment of an inverse matrix corresponding to a change in one element of a given matrix. *Annals of Mathematical Statistics*, 21:124–127, 1950.
- [2] Max A. Woodbury. Inverting modified matrices. *Statistical Research Group, Memo.*, Rep. no. 42, 1950.



**Figure 3** Relative errors and growth factors, with pivoting  
Using growth matrix