



Citation for published version:

Hook, J & Tisseur, F 2017, 'Incomplete LU Preconditioner Based on Max-Plus Approximation of LU Factorization', *SIAM Journal On Matrix Analysis and Applications (SIMAX)*, vol. 38, no. 4, pp. 1160 - 1189. <https://doi.org/10.1137/16M1094579>

DOI:

[10.1137/16M1094579](https://doi.org/10.1137/16M1094579)

Publication date:

2017

Document Version

Publisher's PDF, also known as Version of record

[Link to publication](#)

Publisher Rights

CC BY

The Version of Record is available via: <http://doi.org/10.1137/16M1094579>

University of Bath

Alternative formats

If you require this document in an alternative format, please contact:
openaccess@bath.ac.uk

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

INCOMPLETE LU PRECONDITIONER BASED ON MAX-PLUS APPROXIMATION OF LU FACTORIZATION*

JAMES HOOK[†] AND FRANÇOISE TISSEUR[‡]

Abstract. We present a new method for the a priori approximation of the orders of magnitude of the entries in the LU factors of a complex or real matrix A . This approximation is used to determine the positions of the largest entries in the LU factors of A , and these positions are used as the sparsity pattern for an incomplete LU factorization preconditioner. Our method uses max-plus algebra and is based solely on the moduli of the entries of A . We also present techniques for predicting which permutation matrices will be chosen by Gaussian elimination with partial pivoting. We exploit the strong connection between the field of Puiseux series and the max-plus semiring to prove properties of the max-plus LU factors. Experiments with a set of test matrices from the University of Florida Sparse Matrix Collection show that our max-plus LU preconditioners outperform traditional level of fill methods and have similar performance to those preconditioners computed with more expensive threshold-based methods.

Key words. max-plus algebra, LU factorization, Hungarian scaling, linear systems of equations, sparse matrices, incomplete LU factorization, preconditioning

AMS subject classifications. 65F08, 65F30, 15A23, 15A80

DOI. 10.1137/16M1094579

1. Introduction. Max-plus algebra is the analogue of linear algebra developed for the binary operations max and plus over the real numbers together with $-\infty$, the latter playing the role of additive identity. Max-plus algebraic techniques have already been used in numerical linear algebra to, for example, approximate the orders of magnitude of the roots of scalar polynomials [19], approximate the moduli of the eigenvalues of matrix polynomials [1, 10, 14], and approximate singular values [9]. These approximations have been used as starting points for iterative schemes and in the design of preprocessing steps to improve the numerical stability of standard algorithms [3, 6, 7, 14, 20]. Our aim is to show how max-plus algebra can be used to approximate the sizes of the entries in the LU factors of a complex or real matrix A and how these approximations can subsequently be used in the construction of an incomplete LU (ILU) factorization preconditioner for A .

In order to be able to apply max-plus techniques to the matrix $A \in \mathbb{C}^{n \times n}$ we must first transform it into a max-plus matrix. We do this using the valuation map

$$(1) \quad \mathcal{V}_c : \mathbb{C} \rightarrow \overline{\mathbb{R}} := \mathbb{R} \cup \{-\infty\}, \quad \mathcal{V}_c(x) = \log |x| \quad (\log 0 = -\infty).$$

The valuation map is applied to matrices componentwise, so that $\mathcal{V}_c(A) \in \overline{\mathbb{R}}^{n \times n}$ is a max-plus matrix. Note that for $x, y \in \mathbb{C}$, $\mathcal{V}_c(xy) = \mathcal{V}_c(x) + \mathcal{V}_c(y)$, and when $|x| \gg |y|$ or $|x| \ll |y|$, then $\mathcal{V}_c(x + y) \approx \max\{\mathcal{V}_c(x), \mathcal{V}_c(y)\}$. This suggests using the operations

*Received by the editors September 19, 2016; accepted for publication (in revised form) by D. B. Szyld August 21, 2017; published electronically October 19, 2017.

<http://www.siam.org/journals/simax/38-4/M109457.html>

Funding: This work was supported by Engineering and Physical Sciences Research Council grant EP/I005293. The work of the second author was also supported by a Royal Society-Wolfson Research Merit Award.

[†]Bath Institute for Mathematical Innovation, University of Bath, Bath BA2 7AY, UK (j.l.hook@bath.ac.uk).

[‡]School of Mathematics, The University of Manchester, Manchester M13 9PL, UK (francoise.tisseur@manchester.ac.uk).

max and plus, which we denote by \oplus and \otimes , respectively, in place of the classical addition and multiplication once we have applied the map \mathcal{V}_c .

The fundamental basis for our approximation of the magnitude of the entries of the LU factors of $A \in \mathbb{C}^{n \times n}$ is

- (a) the fact that the entries in the lower triangle of L and the upper triangle of U can be expressed explicitly in terms of determinants of submatrices S of A , and
- (b) the heuristic that, when the matrix S has large variation in the size of its entries, $\mathcal{V}_c(\det(S)) \approx \text{perm}(\mathcal{V}_c(S))$, where perm is the max-plus permanent.

We use (a) and (b) to define a lower triangular max-plus matrix \mathcal{L} and an upper triangular max-plus matrix \mathcal{U} such that

$$(2) \quad \mathcal{V}_c(L) \approx \mathcal{L}, \quad \mathcal{V}_c(U) \approx \mathcal{U},$$

and we refer to \mathcal{L} and \mathcal{U} as the max-plus LU factors of $\mathcal{A} := \mathcal{V}_c(A) \in \overline{\mathbb{R}}^{n \times n}$. The approximation (2) is a heuristic which only aims to capture the order of magnitude of the entries of L and U . One way to think about the max-plus LU approximation of the LU factors of A is as an intermediate between the true LU factors of A and a symbolic or Boolean factorization which, based purely on the pattern of nonzero entries in A , predicts the nonzero patterns of the LU factors. We show that the matrix-matrix product $\mathcal{L} \otimes \mathcal{U}$ is usually not a factorization of \mathcal{A} but that it “balances” \mathcal{A} , in a sense made precise below.

In order for the max-plus approximation to be useful in practice, it is essential that the cost of computing it be less than the cost of computing the LU factorization exactly. We show that the max-plus LU factors can be computed by solving maximally weighted tree problems. As a result we provide an algorithm for computing the LU approximation of $A \in \mathbb{C}^{n \times n}$ with worst case cost $O(n\tau + n^2 \log n)$, where τ is the number of nonzero entries in A . Note that this cost depends on the number of nonzero entries in A and not on the number of nonzero entries in the LU factors of A . Thus while the approximate LU factors will exhibit fill-in just as in the exact case, the cost of computing the approximation is not affected by fill-in and will therefore be less than that of computing the exact LU factors. If the matrix A is first reordered according to its optimal assignment, so that the product of the moduli of the entries on its diagonal is maximized, then our approximation of the LU factors can be computed in parallel by n separate computations, each of individual cost $O(\tau + n \log n)$. If we seek only the positions and values of the k largest entries in each row of U and column of L , or if we seek only the position and values of the entries that are greater in modulus than some threshold, then this cost can be reduced further.

An approximation of the size of the entries in the LU factors of a sparse matrix A can be used to help construct an ILU preconditioner for solving $Ax = b$, that is, a pair of sparse lower- and upper triangular matrices L, U such that the preconditioned matrix $AU^{-1}L^{-1}$ is more amenable to iterative methods such as GMRES [17]. Two classes of ILU preconditioners are threshold ILU and $\text{ILU}(k)$. In threshold ILU, Gaussian elimination is applied to A , but any computed element with modulus less than some threshold value is set to zero. By storing only the larger entries, threshold ILU is able to compute effective preconditioners.

For $\text{ILU}(k)$ preconditioners, a sparsity pattern for the ILU factors is first computed from a symbolic factorization that determines the level of fill-in of each fill-in entry of A [17, sec. 10.3]. A fill-in entry is dropped when its level of fill is above k , and the corresponding entry in the sparsity pattern matrix is set to zero. The ILU factors are then computed using a variant of Gaussian elimination restricted to the sparsity

pattern, such as that presented in [17, Alg. 10.3]. $\text{ILU}(k)$ preconditioners can be computed using fast parallel algorithms, e.g., [12], but that may not reliably result in effective preconditioners because the sparsity pattern of the factors depends only on the sparsity pattern of A and not on the numerical values of its entries.

Our max-plus LU approximation enables us to take a hybrid approach that offers the best of both of these methods. We use the max-plus LU factors \mathcal{L} and \mathcal{U} to define the sparsity pattern of the ILU preconditioners by allowing only entries with a value over a certain threshold. Provided the entries of \mathcal{L} and \mathcal{U} give good approximations of the size of the true LU entries, our approach results in an ILU pair very close to that obtained through standard threshold ILU. Our method for computing the max-plus sparsity pattern makes n independent calls to Dijkstra's algorithm with an early stopping condition. This is remarkably similar to the basic method for computing an $\text{ILU}(k)$ sparsity pattern, which makes n independent calls to the breadth first search algorithm with an early stopping condition. So it is reasonable to assume that the max-plus sparsity pattern can be computed almost as efficiently as the $\text{ILU}(k)$ sparsity pattern. As a result, the total time for computing the max-plus ILU factors should be competitive with that of the $\text{ILU}(k)$ factors, since once we have the max-plus sparsity pattern, we can use the same efficient algorithm for computing the incomplete factors as $\text{ILU}(k)$ does.

The remainder of this paper is organized as follows. In section 2 we introduce the max-plus permanent and discuss how it can be used to approximate the order of magnitude of the determinant of a complex matrix. This approximation forms the basis of our LU factor approximation. In section 3 we define the max-plus LU factors of a max-plus matrix and argue that they can be used to approximate the orders of magnitude of the entries in the LU factors of a complex matrix. We also show how our max-plus LU factorization can be adapted to include pivoting and examine the special case of Hungarian scaled matrices. In section 4 we examine the connection between max-plus LU factors and the LU decomposition of matrices of Puiseux series, and use this connection to prove several of the theoretical results that are stated earlier in the paper. In section 5 we give a derivation of our different max-plus LU algorithms and describe our max-plus ILU preconditioner. In section 6 we apply our max-plus LU approximation and ILU preconditioning technique to a set of test problems from real life scientific computing problems.

Throughout this paper, complex matrices will be denoted by capital letters, with their entries denoted by the corresponding lowercase letter in the usual way, $A = (a_{ij}) \in \mathbb{C}^{n \times n}$. Matrices of complex Puiseux series will be denoted by capital letters with a tilde, and their entries by the corresponding lowercase letter also with a tilde, $\tilde{A} = (\tilde{a}_{ij}) \in \mathbb{C}\{\{z\}\}^{n \times n}$, where $\mathbb{C}\{\{z\}\}$ denotes the field of Puiseux series. Max-plus matrices will be denoted by calligraphic capital letters, and their entries by the corresponding lowercase calligraphic letter, $\mathcal{A} = (a_{ij}) \in \overline{\mathbb{R}}^{n \times n}$. Since the most important results of this paper are the heuristic max-plus approximations, we will present these results in the style of theorems with a justification following each heuristic in lieu of a proof.

2. Heuristic approximation of the determinant. If we replace the sum by a maximum and the product by a summation in the Leibniz formula for the determinant of $A \in \mathbb{C}^{n \times n}$,

$$\det(A) = \sum_{\pi \in \Pi(n)} \text{sgn}(\pi) \prod_{i=1}^n a_{i, \pi(i)},$$

where $\Pi(n)$ is the set of all permutations on $\{1, \dots, n\}$, and replace the complex scalars $a_{i,\pi(i)}$ by scalars $a_{i,\pi(i)} \in \overline{\mathbb{R}}$, we obtain the formula for the max-plus permanent of $\mathcal{A} = (a_{ij}) \in \overline{\mathbb{R}}^{n \times n}$,

$$(3) \quad \text{perm}(\mathcal{A}) = \max_{\pi \in \Pi(n)} \sum_{i=1}^n a_{i,\pi(i)} = \bigoplus_{\pi \in \Pi(n)} \bigotimes_{i=1}^n a_{i,\pi(i)}.$$

The following heuristic is fundamental to our max-plus LU approximation.

HEURISTIC 2.1. *Let $A \in \mathbb{C}^{n \times n}$ be sparse with nonzero entries that vary widely in magnitude, and let \mathcal{V}_c be as in (1). Then*

$$(4) \quad \mathcal{V}_c(\det(A)) \approx \text{perm}(\mathcal{V}_c(A)).$$

Justification. The determinant of A is a sum of terms the vast majority of which are zero (due to sparsity) and the remainder of which vary widely in order of magnitude (due to the wide variation in entry magnitude). The order of magnitude of the sum of a small number of terms of widely varying magnitude can then be approximated by the order of magnitude of the greatest of those terms, which is precisely $\text{perm}(\mathcal{V}_c(A))$. \square

We show in section 4 that the permanent can also be used to calculate the exact asymptotic growth rate of the determinant of a generic matrix of a Puiseux series, which provides some additional support for Heuristic 2.1. In the meantime let us look at a few examples.

Example 2.2. We use the logarithm in base 10 for \mathcal{V}_c and consider

$$A = \begin{bmatrix} 10 & 0 & 1000 \\ 1 & 10 & 0 \\ 0 & 1 & 1 \end{bmatrix}, \quad \mathcal{A} = \mathcal{V}_c(A) = \begin{bmatrix} 1 & -\infty & 3 \\ 0 & 1 & -\infty \\ -\infty & 0 & 0 \end{bmatrix}.$$

For this example, $\text{perm}(\mathcal{A}) = 3$, which provides an order of magnitude approximation of $\det(A) = -900$ since $\log |\det(A)| \approx 2.95$.

Of course we can easily find counterexamples where the approximation in (4) is very poor. However, we can think of these matrices as occupying a set of small measure, so that the order of magnitude of the determinant of a “typical” complex matrix will be well approximated.

Example 2.3. Consider $A \in \begin{bmatrix} \omega_1 & \omega_2 \\ \omega_3 & \omega_4 \end{bmatrix} \in \mathbb{C}^{2 \times 2}$ with $|\omega_i| = 1$, $i = 1, \dots, 4$, so that $\mathcal{V}_c(A) = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$, where $\mathcal{V}_c(x) := \log_{10} |x|$. Choosing $\omega_i = 1$, $i = 1, \dots, 4$, yields a singular A and $\log |\det(A)| = -\infty$, which is not detected by the max-plus approximation since $\text{perm}(\mathcal{V}_c(A)) = 0$. Likewise whenever $\det(A)$ is close to zero, the max-plus approximation will not be accurate. However, for most choices of ω the approximation will capture the order of magnitude of $\det(A)$. Indeed, if each ω_i is an independent random variable uniformly distributed on the unit circle, then $|\det(A)|$ has expected value $\mathbb{E}(|\det(A)|) = 4/\pi \approx 1$, and for small $\epsilon > 0$ the probability $\mathbb{P}(|\det(A)| \leq \epsilon) \approx \epsilon/\pi$. Thus the choices of ω_i for which the max-plus approximation fails to capture the order of magnitude of $\det(A)$ represent a set of small measure.

3. Max-plus LU factors. An LU decomposition of $A \in \mathbb{C}^{n \times n}$ is a factorization of A into two factors, a unit lower triangular matrix denoted by L and an upper triangular matrix denoted by U such that $A = LU$. The entries of the L and U

factors can be given explicitly in terms of determinants of submatrices of A (see [5, p. 35] or [11, p. 11]) by

$$(5) \quad l_{ik} = \det(A([1 : k - 1, i], 1 : k)) / \det(A(1 : k, 1 : k)), \quad i \geq k,$$

$$(6) \quad u_{kj} = \det(A(1 : k, [1 : k - 1, j])) / \det(A(1 : k - 1, 1 : k - 1)), \quad j \geq k,$$

and $l_{ik} = u_{kj} = 0$ for $i, j < k$, where $A(i : j, k : \ell)$ denotes the submatrix of A formed by the intersection of the rows i to j and columns k to ℓ . If both the numerator and denominator in (5)–(6) are zero, then we use the convention $0/0 = 0$. If the denominator is equal to zero but the numerator is not, then we say that A does not admit an LU decomposition. If all of the denominators in (5)–(6) are nonzero, then $A = LU$ is the unique LU decomposition of A .

Based on these formulae we define the *max-plus LU factors* of $\mathcal{A} \in \overline{\mathbb{R}}^{n \times n}$ to be the unit lower triangular max-plus matrix \mathcal{L} and the upper triangular max-plus matrix \mathcal{U} with entries given by

$$(7) \quad l_{ik} = \text{perm}(\mathcal{A}([1 : k - 1, i], 1 : k)) - \text{perm}(\mathcal{A}(1 : k, 1 : k)), \quad i > k, \quad l_{ii} = 0,$$

$$(8) \quad u_{kj} = \text{perm}(\mathcal{A}(1 : k, [1 : k - 1, j])) - \text{perm}(\mathcal{A}(1 : k - 1, 1 : k - 1)), \quad j \geq k,$$

and $l_{ik} = u_{kj} = -\infty$ if $i, j < k$. If the two terms on the right-hand side of (7) or (8) are $-\infty$, then we use the convention $-\infty - (-\infty) = -\infty$. If the second term is $-\infty$ but the first is not, then we say that \mathcal{A} does not admit max-plus LU factors.

HEURISTIC 3.1. *Let $A \in \mathbb{C}^{n \times n}$, and suppose that $\mathcal{V}_c(A)$ admits max-plus LU factors $\mathcal{L}, \mathcal{U} \in \overline{\mathbb{R}}^{n \times n}$. Then A admits an LU decomposition $A = LU$ with*

$$\mathcal{V}_c(L) \approx \mathcal{L}, \quad \mathcal{V}_c(U) \approx \mathcal{U}.$$

Justification. From Heuristic 2.1, we expect that the determinant of a submatrix of A is zero if and only if the permanent of the corresponding submatrix of $\mathcal{V}_c(A)$ is minus infinity. Therefore if $\mathcal{V}_c(A)$ admits max-plus LU factors, then A admits an LU factorization $A = LU$, where the LU factors are as in (5)–(6). Taking the valuation of these expressions, applying Heuristic 2.1, and comparing to (7)–(8) gives the required result. \square

Example 3.2. The matrix A of Example 2.2 has LU factorization

$$A = \begin{bmatrix} 10 & 0 & 1000 \\ 1 & 10 & 0 \\ 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0.1 & 1 & 0 \\ 0 & 0.1 & 1 \end{bmatrix} \begin{bmatrix} 10 & 0 & 1000 \\ 0 & 10 & -100 \\ 0 & 0 & 11 \end{bmatrix} = LU$$

and max-plus LU factors

$$\mathcal{L} = \begin{bmatrix} 0 & -\infty & -\infty \\ -1 & 0 & -\infty \\ -\infty & -1 & 0 \end{bmatrix}, \quad \mathcal{U} = \begin{bmatrix} 1 & -\infty & 3 \\ -\infty & 1 & 2 \\ -\infty & -\infty & 1 \end{bmatrix},$$

which provide good approximations of the orders of magnitude of the entries in L, U .

Example 3.3. The LU factorization of the matrix A of Example 2.3 with $|\omega_i| = 1$ is given by

$$A = \begin{bmatrix} \omega_1 & \omega_2 \\ \omega_3 & \omega_4 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ \frac{\omega_3}{\omega_1} & 1 \end{bmatrix} \begin{bmatrix} \omega_1 & \omega_2 \\ 0 & \omega_4 - \frac{\omega_2 \omega_3}{\omega_1} \end{bmatrix} = LU,$$

and the max-plus LU factors of $\mathcal{V}_c(A)$ are given by

$$\mathcal{L} = \begin{bmatrix} 0 & -\infty \\ 0 & 0 \end{bmatrix}, \quad \mathcal{U} = \begin{bmatrix} 0 & 0 \\ -\infty & 0 \end{bmatrix}.$$

Since $\mathcal{V}_c(A)$ is independent of the choice of the ω_i since $|\omega_i| = 1$, so are its max-plus LU factors. The (2, 2) entry of U is the only entry where the max-plus approximation is not guaranteed to be perfectly accurate, but for most choices of the ω_i the max-plus approximation captures the order of magnitude of the entries of L and U . There is, however, a small set of parameter values of small measure for which the max-plus approximation is not accurate.

Our definition of the max-plus LU factors of a max-plus matrix was chosen so that that we could use it to approximate the orders of magnitude of the entries in the LU factors of a complex matrix. But what do the max-plus LU factors of a max-plus matrix $\mathcal{A} \in \mathbb{R}^{n \times n}$ tell us about \mathcal{A} ?

THEOREM 3.4. *Suppose that $\mathcal{A} = (a_{ij}) \in \mathbb{R}^{n \times n}$ has max-plus LU factors $\mathcal{L}, \mathcal{U} \in \mathbb{R}^{n \times n}$. Then for each $i, j = 1, \dots, n$ either*

$$(9) \quad (\mathcal{L} \otimes \mathcal{U})_{ij} := \max_{1 \leq k \leq n} (l_{ik} + u_{kj}) > a_{ij},$$

where the maximum is attained by at least two different values of k , or

$$(10) \quad (\mathcal{L} \otimes \mathcal{U})_{ij} = a_{ij}.$$

The proof of Theorem 3.4 is provided in section 4. We say that the max-plus matrix product $\mathcal{L} \otimes \mathcal{U}$ *balances* \mathcal{A} .

3.1. Pivoting. After k steps of Gaussian elimination applied to $A \in \mathbb{C}^{n \times n}$, the matrix A is reduced to

$$(11) \quad M_k \cdots M_1 A = U^{(k)} = \begin{bmatrix} U_{11}^{(k)} & U_{12}^{(k)} \\ 0 & U_{22}^{(k)} \end{bmatrix},$$

where the M_i are Gauss transforms and $U_{11}^{(k)} \in \mathbb{C}^{k \times k}$ is upper triangular. Like the LU factors themselves, the entries of $U^{(k)}$ can be expressed in terms of determinants of submatrices of A , as the next lemma shows.

LEMMA 3.5. *Let $A \in \mathbb{C}^{n \times n}$ have an LU factorization, and let $U^{(k)}$ be as in (11). Then*

$$(12) \quad u_{ij}^{(k)} = \begin{cases} \det(A([1:k, i], [1:k, j])) / \det(A(1:k, 1:k)), & i, j > k, \\ u_{ij} & \text{otherwise,} \end{cases}$$

where $U = U^{(n-1)}$ is the upper triangular factor in the LU factorization of A .

Proof. Suppose that $i, j > k$. Let R_i and C_j be elementary matrices swapping rows $k + 1$ and i , and columns $k + 1$ and j , respectively. Define $A' := R_i A C_j$, and let $U^{(k)}$ be the matrix obtained after performing k steps of Gaussian elimination on A' . Then $U^{(k)} = R_i U^{(k)} C_j$ and in particular $u_{k+1, k+1}^{(k)} = u_{ij}^{(k)}$. The Gauss transform M'_{k+1} at step $k + 1$ has the form $I + m_{k+1} e_{k+1}^T$, where $e_i^T m_{k+1} = 0$ for $i = 1, \dots, k + 1$

so that the $(k + 1, k + 1)$ entries of $U^{(k)}$ and $M'_{k+1}U^{(k)} = U^{(k+1)}$ are the same, that is, $u_{k+1,k+1}^{(k)} = u_{k+1,k+1}^{(k+1)}$. But $u_{k+1,k+1}^{(k+1)} = u'_{k+1,k+1}$, and by (6),

$$\begin{aligned} u'_{k+1,k+1} &= \det(A'(1 : k + 1, [1 : k + 1])) / \det(A'(1 : k, 1 : k)) \\ &= \det(A([1 : k, i], [1 : k, j])) / \det(A(1 : k, 1 : k)). \end{aligned}$$

The next steps of Gaussian elimination leave the (i, j) entries of $U^{(k)}$ with $\min\{i, j\} \leq k$ unchanged so that $u_{ij}^{(k)} = u_{ij}$ for $\min\{i, j\} \leq k$. □

We say that $A \in \mathbb{C}^{n \times n}$ is *partial pivoting free* if

$$(13) \quad |u_{k+1,k+1}^{(k)}| = \max_{k+1 \leq i \leq n} |u_{i,k+1}^{(k)}|, \quad k = 0, \dots, n - 1,$$

where $U^{(0)} = A$. If the matrix A is partial pivoting free, then it is possible to apply Gaussian elimination with partial pivoting to A without the need for any row interchanges.

Let $A = LU$ be the LU decomposition of A , and suppose that we compute an approximate LU pair $\widehat{L}, \widehat{U} \in \mathbb{C}^{n \times n}$ using Gaussian elimination. The *backward error* of these approximate LU factors is equal to the perturbation $\Delta A \in \mathbb{C}^{n \times n}$ such that $A + \Delta A = \widehat{L}\widehat{U}$ and is known to satisfy [8, Lem. 9.6]:

$$\|\Delta A\|_\infty \leq \frac{nu}{1 - nu} \|L\| \|U\|_\infty \leq \frac{nu}{1 - nu} (1 + 2(n^2 - n)\rho_n(A)) \|A\|_\infty,$$

where u is the unit roundoff and $\rho_n(A)$ is the *growth factor* for A defined by

$$(14) \quad \rho_n(A) = \frac{\max_{0 \leq k \leq n-1} (\max_{k \leq i, j \leq n} |u_{ij}^{(k)}|)}{\max_{i, j} |a_{ij}|}.$$

Thus if $\|\Delta A\|_\infty$ is small relative to $\|A\|_\infty$, which certainly happens when $\rho_n(A)$ is small, then the factorization is stable; otherwise it is unstable [8, sec. 9.3].

In analogy to (13), we say that the max-plus matrix $\mathcal{A} \in \overline{\mathbb{R}}^{n \times n}$ is *partial pivoting free* if

$$(15) \quad u_{k+1,k+1}^{(k)} = \max_{k+1 \leq i \leq n} u_{i,k+1}^{(k)}, \quad k = 0, \dots, n - 1,$$

where $u_{ij}^{(0)} := a_{ij}$ and

$$(16) \quad u_{ij}^{(k)} := \begin{cases} \text{perm}(\mathcal{A}([1 : k, i], [1 : k, j])) - \text{perm}(\mathcal{A}(1 : k, 1 : k)), & i, j > k, \\ u_{ij} & \text{otherwise.} \end{cases}$$

Also, in analogy to (14) we define the *max-plus growth factor* of $\mathcal{A} \in \overline{\mathbb{R}}^{n \times n}$ by

$$(17) \quad \varrho_n(\mathcal{A}) = \max_{0 \leq k \leq n-1} \left(\max_{k \leq i, j \leq n} u_{ij}^{(k)} \right) - \max_{1 \leq i, j \leq n} a_{ij}.$$

THEOREM 3.6. *If $\mathcal{A} \in \overline{\mathbb{R}}^{n \times n}$ is partial pivoting free, then $\varrho_n(\mathcal{A}) = 0$.*

The proof of Theorem 3.6 is deferred to the end of section 4.

HEURISTIC 3.7. *For $A \in \mathbb{C}^{n \times n}$ we have $\mathcal{V}_c(\rho_n(A)) \approx \varrho_n(\mathcal{V}_c(A))$.*

Justification. From Lemma 3.5 and Heuristic 2.1 we have $\mathcal{V}_c(U^{(k)}) \approx \mathcal{U}^{(k)}$ for $k = 0, \dots, n - 1$. Then using this and taking the valuation of (14) yields

$$\begin{aligned} \mathcal{V}_c(\rho_n(A)) &= \log \left| \frac{\max_{0 \leq k \leq n-1} \left(\max_{k \leq i, j \leq n} |u_{ij}^{(k)}| \right)}{\max_{i, j} |a_{ij}|} \right| \\ &= \max_{0 \leq k \leq n-1} \left(\max_{k \leq i, j \leq n} \log |u_{ij}^{(k)}| \right) - \max_{1 \leq i, j \leq n} \log |a_{ij}| \\ &\approx \max_{0 \leq k \leq n-1} \left(\max_{k \leq i, j \leq n} u_{ij}^{(k)} \right) - \max_{1 \leq i, j \leq n} a_{ij} = \varrho_n(A). \end{aligned}$$

If $\mathcal{V}_c(A)$ is partial pivoting free, then it follows from Theorem 3.6 that $\varrho_n(\mathcal{V}_c(A)) = 0$ so that, based on Heuristic 3.7, the growth factor $\rho_n(A)$ should be of order one, implying a backward stable LU factorization. As before, this is a heuristic, and it is not difficult to construct counterexample matrices A for which $\mathcal{V}_c(A)$ is partial or full pivoting free but that cannot be factorized in a stable way without further pivoting.

Theorem 3.6 and Heuristic 3.7 suggest applying a permutation P to a given A such that $\mathcal{V}_c(PA)$ is partial pivoting free. We show in section 5.2 how to update our max-plus LU algorithm to include partial pivoting. Another option is to apply Hungarian scaling, which is a two-sided diagonal scaling applied to $A \in \mathbb{C}^{n \times n}$ along with a permutation P that maximizes the product of the moduli of the diagonal entries of the matrix

$$(18) \quad H = PD_1AD_2,$$

where $D_1, D_2 \in \mathbb{R}^{n \times n}$ are nonsingular and diagonal, and such that H 's entries satisfy

$$(19) \quad |h_{ij}| \leq 1, \quad |h_{ii}| = 1, \quad i, j = 1, \dots, n.$$

We refer to any complex matrix satisfying (19) as a *Hungarian matrix*. The max-plus matrix $\mathcal{H} = \mathcal{V}_c(H) \in \overline{\mathbb{R}}^{n \times n}$ is such that $h_{ij} \leq 0$, $h_{ii} = 0$, $i, j = 1, \dots, n$, and is referred to as a *max-plus Hungarian matrix*.

THEOREM 3.8. *Max-plus Hungarian matrices always admit max-plus LU factors.*

Proof. Suppose that $\mathcal{H} \in \mathbb{R}_{\max}^{n \times n}$ is Hungarian. From (7) and (8) we have that

$$\text{perm}(\mathcal{H}(1 : k, 1 : k)) \neq -\infty, \quad k = 1, \dots, n,$$

is a sufficient condition for \mathcal{H} to admit max-plus LU factors. Since $h_{ii} = 0$ for $i = 1, \dots, n$, we have $\text{perm}(\mathcal{H}(1 : k, 1 : k)) \geq 0$ for $k = 1, \dots, n$. \square

THEOREM 3.9. *A max-plus Hungarian matrix is partial pivoting free.*

Proof. It follows from (15) and (16) that \mathcal{H} is partial pivoting free if

$$(20) \quad \text{perm}(\mathcal{H}(1 : k + 1, 1 : k + 1)) \geq \text{perm}(\mathcal{H}([1 : k, i], [1 : k + 1]))$$

for all $i = k + 1, \dots, n$ and for all $k = 0, \dots, n - 1$. But since $h_{ij} \leq 0$ for all i, j , the permanent of any submatrix of \mathcal{H} must be nonpositive. Hence the right-hand side of the inequality in (20) must be less than or equal to zero. Also, since \mathcal{H} has zero diagonal entries, the permanent of any principal leading submatrix of \mathcal{H} is equal to zero. Therefore the inequality in (20) must have left-hand side equal to zero, so that \mathcal{H} is partial pivoting free. \square

Therefore, given $A \in \mathbb{C}^{n \times n}$, we apply Hungarian scaling to obtain $H = PD_1AD_2$, and from Theorems 3.6 and 3.9 and Heuristic 3.7, we expect that it should be possible to factorize H in a stable way without any need for interchange. This preprocessing technique was originally suggested by Olschowka and Neumaier in [15]. They prove that Hungarian scaling removes the need for interchange in Gaussian elimination for some special classes of matrices. While our results do not constitute a definite theorem, they provide some intuitive explanation for the widely observed fact that Hungarian scaling significantly reduces the need for pivoting (see section 6).

Example 3.10. Let $A = \begin{bmatrix} 1 & 10^{-3} \\ 10 & 1 \end{bmatrix}$. We have that $\mathcal{V}_c(A) = \begin{bmatrix} 0 & -3 \\ 1 & 0 \end{bmatrix}$ is not partial pivoting free since $\mathcal{U}^{(0)} = \mathcal{V}_c(A)$ is such that $u_{21}^{(0)} = 1 > u_{11}^{(0)} = 0$. Similarly A is not partial pivoting free. It is easy to check that the matrices PA and $\mathcal{V}_c(PA)$ with $P = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$ are both partial pivoting free. Now a Hungarian scaling for A with $P = I$, $D_1 = \text{diag}(1, 10^{-2})$, and $D_2 = \text{diag}(1, 10^2)$ is given by $H = PD_1AD_2 = \begin{bmatrix} 1 & 10^{-1} \\ 10^{-1} & 1 \end{bmatrix}$ so that $\mathcal{V}_c(H) = \begin{bmatrix} 0 & -1 \\ -1 & 0 \end{bmatrix}$. Theorem 3.9 guarantees that $\mathcal{V}_c(H)$ is partial pivoting free, and it is easy to check that H is also partial pivoting free.

4. Puiseux series. There is a stronger connection between the field of complex Puiseux series and the semiring $\mathbb{R}_{\max} = (\mathbb{R}, \oplus, \otimes)$ than between the field of complex numbers and \mathbb{R}_{\max} , which we now exploit to prove properties of the max-plus LU factors as well as Theorems 3.4 and 3.6 and to provide further justification of Heuristic 2.1. This section is not needed for the derivation of the max-plus LU algorithms presented in section 5.

Complex Puiseux series

$$(21) \quad f(z) = \sum_{i=k}^{\infty} c_i z^{\frac{i}{m}},$$

with $m \in \mathbb{N}$, $k \in \mathbb{Z}$, $c_i \in \mathbb{C}$, $i \geq k$, and $c_k \neq 0$ form an algebraically closed field under addition and multiplication denoted by $\mathbb{C}\{\{z\}\}$. On that field, we define the valuation

$$(22) \quad \mathcal{V}_p : \mathbb{C}\{\{z\}\} \mapsto \overline{\mathbb{R}}, \quad \mathcal{V}_p(f) = -k/m,$$

that is, the valuation of a Puiseux series is minus the degree of its lowest order term. This valuation provides a near homeomorphism between $\mathbb{C}\{\{z\}\}$ and \mathbb{R}_{\max} ,

$$(23) \quad \begin{aligned} \mathcal{V}_p(fg) &= \mathcal{V}_p(f) \otimes \mathcal{V}_p(g) && \text{for all } f, g \in \mathbb{C}\{\{z\}\}, \\ \mathcal{V}_p(f+g) &\leq \mathcal{V}_p(f) \oplus \mathcal{V}_p(g) && \text{for all } f, g \in \mathbb{C}\{\{z\}\}, \\ \mathcal{V}_p(f+g) &= \mathcal{V}_p(f) \oplus \mathcal{V}_p(g) && \text{for almost all } f, g \in \mathbb{C}\{\{z\}\}, \end{aligned}$$

where the third relation holds except for when $\mathcal{V}_p(f) = \mathcal{V}_p(g)$ and where the coefficient of the lowest order term of f is equal to minus that of g . As for complex matrices, the valuation \mathcal{V}_p is applied componentwise to matrices with Puiseux series entries. We decorate matrices in $\mathbb{C}\{\{z\}\}^{n \times n}$ with a tilde to distinguish them from matrices in $\mathbb{C}^{n \times n}$.

Any entry of $\tilde{A} \in \mathbb{C}\{\{z\}\}^{n \times n}$ can be written as

$$\tilde{a}_{ij} = c_{ij} z^{-\mathcal{V}_p(\tilde{a}_{ij})} + \text{higher order terms},$$

where $C = (c_{ij}) =: \mathcal{L}(\tilde{A}) \in \mathbb{C}^{n \times n}$ is the matrix of lowest order term coefficients of \tilde{A} with $\mathcal{L} : \mathbb{C}\{\{z\}\}^{n \times n} \mapsto \mathbb{C}^{n \times n}$. For a set of permutations $\Phi \subset \Pi(n)$, we define the

map $g_\Phi : \mathbb{C}^{n \times n} \mapsto \mathbb{C}$ by

$$(24) \quad g_\Phi(C) = \sum_{\pi \in \Phi} \text{sign}(\pi) \prod_{i=1}^n c_{i\pi(i)}.$$

Note that $g_{\Pi(n)}(C) = \det(C)$. For $\mathcal{A} \in \overline{\mathbb{R}}^{n \times n}$ such that $\text{perm}(\mathcal{A}) \neq -\infty$ we denote by

$$\text{ap}(\mathcal{A}) = \left\{ \pi \in \Pi(n) : \sum_{i=1}^n a_{i\pi(i)} = \text{perm}(\mathcal{A}) \right\}$$

the set of optimal assignments for \mathcal{A} .

The next lemma identifies the set of matrices with Puiseux series entries such that the valuation of the determinant is exactly the permanent of the valuation (see Heuristic 2.1 for matrices with complex entries).

LEMMA 4.1. *Let $\tilde{A} \in \mathbb{C}\{\{z\}\}^{n \times n}$, and suppose that $g_{\text{ap}(\mathcal{V}_p(\tilde{A}))}(\mathcal{L}(\tilde{A})) \neq 0$, where g , ap , and \mathcal{L} are defined above. Then $\mathcal{V}_p(\det(\tilde{A})) = \text{perm}(\mathcal{V}_p(\tilde{A}))$.*

Proof. Let $\mathcal{A} = \mathcal{V}_p(\tilde{A}) \in \overline{\mathbb{R}}^{n \times n}$. First suppose that $\text{perm}(\mathcal{A}) = -\infty$. Then for each permutation $\pi \in \Pi(n)$ there exists i such that $a_{i\pi(i)} = -\infty$ so that $\tilde{a}_{i\pi(i)} = 0$. Thus $\det(\tilde{A}) = \sum_{\pi \in \Pi(n)} \text{sign}(\pi) \prod_{i=1}^n \tilde{a}_{i\pi(i)} = 0$ and $\mathcal{V}_p(\det(\tilde{A})) = \text{perm}(\mathcal{A})$.

Now suppose that $\text{perm}(\mathcal{A}) \neq -\infty$, and let $C = \mathcal{L}(\tilde{A})$. Then

$$\det(\tilde{A}) = \sum_{\pi \in \Pi(n)} \text{sign}(\pi) \prod_{i=1}^n \tilde{a}_{i\pi(i)} = \sum_{\pi \in \Pi(n)} \text{sign}(\pi) \left(z^{-\sum_{i=1}^n a_{i\pi(i)}} \prod_{i=1}^n c_{i\pi(i)} + \text{h.o.t.} \right),$$

where h.o.t. stands for higher order terms. We break the sum into two parts, one over $\text{ap}(\mathcal{A})$ and one over $\Pi(n) \setminus \text{ap}(\mathcal{A})$. We have that

$$\begin{aligned} \sum_{\pi \in \text{ap}(\mathcal{A})} \text{sign}(\pi) \left(z^{-\sum_{i=1}^n a_{i\pi(i)}} \prod_{i=1}^n c_{i\pi(i)} + \text{h.o.t.} \right) &= z^{-\text{perm}(\mathcal{A})} \sum_{\pi \in \text{ap}(\mathcal{A})} \text{sign}(\pi) \prod_{i=1}^n c_{i\pi(i)} \\ &\quad + \text{h.o.t.} \\ &= z^{-\text{perm}(\mathcal{A})} g_{\text{ap}(\mathcal{A})}(C) + \text{h.o.t.}, \end{aligned}$$

where $g_{\text{ap}(\mathcal{A})}(C)$ is defined in (25). Since for $\pi \in \Pi(n) \setminus \text{ap}(\mathcal{A})$, $\sum_{i=1}^n a_{i\pi(i)} < \text{perm}(\mathcal{A})$, $z^{-\sum_{i=1}^n a_{i\pi(i)}}$ is higher order than $z^{-\text{perm}(\mathcal{A})}$ and so is

$$\sum_{\pi \in \Pi(n) \setminus \text{ap}(\mathcal{A})} \text{sign}(\pi) \left(z^{-\sum_{i=1}^n a_{i\pi(i)}} \prod_{i=1}^n c_{i\pi(i)} + \text{h.o.t.} \right).$$

Hence, $\det(\tilde{A}) = z^{-\text{perm}(\mathcal{A})} g_{\text{ap}(\mathcal{A})}(C) + \text{h.o.t.}$, and $\mathcal{V}_p(\det(\tilde{A})) = \text{perm}(\mathcal{A})$ since $g_{\text{ap}(\mathcal{A})}(C) \neq 0$. □

The next lemma will be useful in showing that $\mathcal{V}_p(\det(\tilde{A})) = \text{perm}(\mathcal{V}_p(\tilde{A}))$ holds for generic matrices $\tilde{A} \in \mathbb{C}\{\{z\}\}^{n \times n}$ but also in explaining what we mean by generic in this context.

LEMMA 4.2. *Let g_Φ be as in (24). Then the set*

$$(25) \quad \mathcal{G}_n = \{C \in \mathbb{C}^{n \times n} : g_\Phi(C) \neq 0 \text{ for all nonempty } \Phi \subset \Pi(n)\}$$

is a generic (open and dense) subset of $\mathbb{C}^{n \times n}$.

Proof. For each $\Phi \subset \Pi(n)$, $g_\Phi(C)$ is a polynomial in the coefficients of C . A polynomial is either identically equal to zero or only zero on some low dimensional subset. Therefore

$$V(g_\Phi) = \{C \in \mathbb{C}^{n \times n} : g_\Phi(C) = 0\}$$

is either the whole of $\mathbb{C}^{n \times n}$ or it is a lower dimensional subset of $\mathbb{C}^{n \times n}$. Choose some permutation $\pi \in \Phi$, and define $C_\pi \in \mathbb{C}^{n \times n}$ by $c_{ij} = 1$ if $j = \pi(i)$ and $c_{ij} = 0$ otherwise. By construction we have $g_\Phi(C_\pi) = \text{sign}(\pi) = \pm 1 \neq 0$ and therefore $C_\pi \notin V(g_\Phi)$. Therefore $V(g_\Phi)$ cannot be the whole of $\mathbb{C}^{n \times n}$ and must instead be a lower dimensional subset. Thus $\mathbb{C}^{n \times n} \setminus V(g_\Phi)$ is a generic subset of $\mathbb{C}^{n \times n}$. Finally note that

$$\mathcal{G}_n = \bigcap_{\phi \subset \Pi(n)} \{\mathbb{C}^{n \times n} \setminus V(g_\phi)\}$$

is a finite intersection of generic subsets and is therefore generic. □

Now if $\tilde{A} \in \mathbb{C}\{\{z\}\}^{n \times n}$ is such that $\mathcal{L}(\tilde{A}) \in \mathcal{G}_n$, then $g_{\text{ap}(\mathcal{V}_p(\tilde{A}))}(\mathcal{L}(\tilde{A})) \neq 0$. Lemma 4.2 states that \mathcal{G}_n is a generic set, so that the property $g_{\text{ap}(\mathcal{V}_p(\tilde{A}))}(\mathcal{L}(\tilde{A})) \neq 0$ is a generic property for $\tilde{A} \in \mathbb{C}\{\{z\}\}^{n \times n}$ with respect to the topology induced by the map $\mathcal{L} : \mathbb{C}\{\{z\}\}^{n \times n} \mapsto \mathbb{C}^{n \times n}$. A more intuitive way of understanding this result is that, if we have a matrix \tilde{A} where the leading order coefficients $\mathcal{L}(\tilde{A})$ have been chosen at random, according to a continuous distribution, then with probability one $g_{\text{ap}(\mathcal{V}_p(\tilde{A}))}(\mathcal{L}(\tilde{A})) \neq 0$ will hold. We then say that $\mathcal{V}_p(\det(\tilde{A})) = \text{perm}(\mathcal{V}_p(\tilde{A}))$ holds for “almost all” $\tilde{A} \in \mathbb{C}\{\{z\}\}^{n \times n}$.

Example 4.3. Consider

$$\tilde{A} = \begin{bmatrix} z^{-1} & 0 & z^{-3} \\ 1 & z^{-1} & 0 \\ 0 & 1 & 1 \end{bmatrix}, \quad \mathcal{V}_p(\tilde{A}) = \begin{bmatrix} 1 & -\infty & 3 \\ 0 & 1 & -\infty \\ -\infty & 0 & 0 \end{bmatrix}, \quad \mathcal{L}(\tilde{A}) = \begin{bmatrix} 1 & 0 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix}.$$

It is easy to check that $\mathcal{L}(\tilde{A}) \in \mathcal{G}_3$, $\det(\tilde{A}) = -z^{-3} + z^{-2}$, and that $\mathcal{V}_p(\det(\tilde{A})) = \text{perm}(\mathcal{V}_p(\tilde{A})) = 3$ as expected from Lemma 4.1.

We now show how to use Puiseux series to further justify Heuristic 2.1.

Justification of Heuristic 2.1. For $f(z) = cz^{-\mathcal{V}_p(f)} + \text{h.o.t.} \in \mathbb{C}\{\{z\}\}$ with $c \neq 0$, we have that $\log_{|z|} |f(z)| \rightarrow -\mathcal{V}_p(f)$ as $|z| \rightarrow 0$. For a desired level of approximation accuracy, we define the *domain of the asymptotic regime of f* to be the neighborhood of zero, denoted by $R(f) \subset \mathbb{C}$, such that $\log_{|z_0|} |f(z_0)| \approx -\mathcal{V}_p(f)$ whenever $z_0 \in R(f)$. Now suppose that $x_0 \in \mathbb{C}$ is some value of interest and that we know $\mathcal{V}_p(f)$ but not f , where $f \in \mathbb{C}\{\{z\}\}$ is a Puiseux series with $f(\tilde{z}_0) = x_0$ for some $\tilde{z}_0 \in \mathbb{C}$. Then, assuming that \tilde{z}_0 is in the domain of the asymptotic regime of f , we have

$$(26) \quad \mathcal{V}_c(x_0) = \log |x_0| = \log |f(\tilde{z}_0)| = \log |\tilde{z}_0| \log_{|\tilde{z}_0|} |f(\tilde{z}_0)| \approx -\log |\tilde{z}_0| \mathcal{V}_p(f).$$

This approximation falls short of being a theorem because, given only \tilde{z}_0 and $\mathcal{V}_p(f)$, we have no way of guaranteeing that \tilde{z}_0 is in the domain of the asymptotic regime of f . In other words, there is no uniform scale for determining what constitutes a small value of $z \in \mathbb{C}$.

We can apply the same idea to approximate the determinant of $A \in \mathbb{C}^{n \times n}$. Suppose that we know $\mathcal{V}_p(\tilde{A})$ with $\tilde{A} \in \mathbb{C}\{\{z\}\}^{n \times n}$ such that $\tilde{A}(\tilde{z}_0) = A$ for some $\tilde{z}_0 \in \mathbb{C}$

and $\mathcal{L}(\tilde{A}) \in \mathcal{G}_n$. Then, assuming \tilde{z}_0 is in the domain of the asymptotic regime of \tilde{A} , it follows from (26) that

$$(27) \quad \mathcal{V}_c(A) \approx -\log |\tilde{z}_0| \mathcal{V}_p(\tilde{A}).$$

Since $\tilde{A}(\tilde{z}_0) = A$ we have $\det(\tilde{A}(\tilde{z}_0)) = \det(A)$. Assuming that \tilde{z}_0 is in the domain of the asymptotic regime of $\det(\tilde{A}) \in \mathbb{C}\{\{z\}\}$ and applying (26), we have

$$\mathcal{V}_c(\det(A)) \approx -\log |\tilde{z}_0| \mathcal{V}_p(\det(\tilde{A})).$$

Using Lemma 4.1 and (27), we obtain that $\mathcal{V}_c(\det(A)) \approx \text{perm}(\mathcal{V}_c(A))$, which provides another justification for Heuristic 2.1. \square

We will need the next lemma.

LEMMA 4.4. *Let $C \in \mathbb{C}^{n \times n} \setminus \mathcal{G}_n$ with \mathcal{G}_n as in (25). Then any $k \times k$ submatrix of C is in $\mathbb{C}^{k \times k} \setminus \mathcal{G}_k$.*

Proof. If $P, Q \in \mathbb{C}^{n \times n}$ are permutation matrices, then $C \in \mathbb{C}^{n \times n} \setminus \mathcal{G}_n$ if and only if $PCQ \in \mathbb{C}^{n \times n} \setminus \mathcal{G}_n$, so it suffices to prove the result for the principal submatrix C of order k , which we denote by S . For any $\Psi \subset \Pi(k)$ we construct $\Phi \subset \Pi(n)$ by setting $\Phi = \{\pi[\varphi] \in \Pi(n) : \varphi \in \Psi\}$, where

$$\pi[\varphi](i) = \begin{cases} \varphi(i), & 1 \leq i \leq k, \\ i, & k+1 \leq i \leq n. \end{cases}$$

Then

$$g_\Phi(C) = \sum_{\pi \in \Phi} \text{sign}(\pi) \prod_{i=1}^n c_{i\pi(i)} = \sum_{\varpi \in \Psi} \text{sign}(\varpi) \prod_{i=1}^k s_{i\varpi(i)} \prod_{i=k+1}^n c_{ii} = g_\Psi(S) \prod_{i=k+1}^n c_{ii}$$

so that $g_\Psi(S) = 0$ if and only if $g_\Phi(C) = 0$. Hence $C \in \mathbb{C}^{n \times n} \setminus \mathcal{G}_n$ if and only if $S \in \mathbb{C}^{k \times k} \setminus \mathcal{G}_k$. \square

As for complex matrices, an LU factorization of $\tilde{A} \in \mathbb{C}\{\{z\}\}^{n \times n}$ is a factorization of \tilde{A} into a lower triangular matrix $\tilde{L} \in \mathbb{C}\{\{z\}\}^{n \times n}$ with ones on the diagonal and an upper triangular matrix $\tilde{U} \in \mathbb{C}\{\{z\}\}^{n \times n}$ such that $\tilde{A} = \tilde{L}\tilde{U}$. When the factorization exists, the nonzero entries of \tilde{L} and \tilde{U} can be defined as in (5)–(6) with \tilde{A} in place of A . The next result should be compared to Heuristic 3.1.

THEOREM 4.5. *Let $\tilde{A} \in \mathbb{C}\{\{z\}\}^{n \times n}$ be such that $\mathcal{L}(\tilde{A}) \in \mathcal{G}_n$, and suppose that $\mathcal{V}_p(\tilde{A}) \in \mathbb{R}^{n \times n}$ admits max-plus LU factors $\mathcal{L}, \mathcal{U} \in \mathbb{R}^{n \times n}$. Then \tilde{A} admits an LU factorization $\tilde{A} = \tilde{L}\tilde{U}$, where*

$$\mathcal{V}_p(\tilde{L}) = \mathcal{L}, \quad \mathcal{V}_p(\tilde{U}) = \mathcal{U}.$$

If for $\tilde{A} \in \mathbb{C}\{\{z\}\}^{n \times n}$, $\mathcal{V}_p(\tilde{A})$ does not admit max-plus LU factors, then \tilde{A} does not admit an LU factorization.

Proof. Let $\mathcal{A} = \mathcal{V}_p(\tilde{A})$. From Lemmas 4.1 and 4.4 we have

$$\mathcal{V}_p(\det(\tilde{A}([i_1, \dots, i_k], [j_1, \dots, j_k]))) = \text{perm}(\mathcal{A}([i_1, \dots, i_k], [j_1, \dots, j_k]))$$

for all submatrices of \tilde{A} since $\mathcal{L}(\tilde{A}) \in \mathcal{G}_n$. Therefore a submatrix of \tilde{A} has zero determinant if and only if the corresponding submatrix of $\mathcal{V}_p(\tilde{A})$ has permanent equal

to minus infinity. Thus if \mathcal{A} admits max-plus LU factors, then an LU factorization of \tilde{A} exists with entries given by (5)–(6).

If $\mathcal{V}_p(\tilde{A})$ does not have max-plus LU factors, then this means that for some i, j, k the first term on the right-hand side of (7) or (8) is equal to $-\infty$ but the second term is not. As a result, the denominator on the right-hand side of (5) or (6) is equal to 0 but the numerator is not, so \tilde{A} does not have an LU factorization. \square

Recall from section 3 that for $\mathcal{A}, \mathcal{B}, \mathcal{C} \in \overline{\mathbb{R}}^{n \times n}$ the product $\mathcal{A} \otimes \mathcal{B}$ balances \mathcal{C} if for every $i, j = 1, \dots, n$ either $(\mathcal{A} \otimes \mathcal{B})_{ij} = c_{ij}$ or $(\mathcal{A} \otimes \mathcal{B})_{ij} = \max_{1 \leq k \leq n} (a_{ik} + b_{kj}) > c_{ij}$, where the maximum must be attained by at least two different values of k .

LEMMA 4.6. *Let $\tilde{A}, \tilde{B} \in \mathbb{C}\{\{z\}\}^{n \times n}$. Then the product $\mathcal{V}_p(\tilde{A}) \otimes \mathcal{V}_p(\tilde{B})$ balances $\mathcal{V}_p(\tilde{A}\tilde{B})$.*

Proof. We have $(\tilde{A}\tilde{B})_{ij} = cz^{-\max_{1 \leq k \leq n} (\mathcal{V}_p(\tilde{a}_{ik}) + \mathcal{V}_p(\tilde{b}_{kj}))} + \text{h.o.t.}$, where $c \in \mathbb{C}$ is the coefficient of the lowest order term in the sum. Therefore

$$(28) \quad \mathcal{V}_p(\tilde{A}\tilde{B})_{ij} = \max_{1 \leq k \leq n} (\mathcal{V}_p(\tilde{a}_{ik}) + \mathcal{V}_p(\tilde{b}_{kj})) = (\mathcal{V}_p(\tilde{A}) \otimes \mathcal{V}_p(\tilde{B}))_{ij},$$

unless $c = 0$, which is only possible if the maximum in (28) is attained more than once, in which case $\mathcal{V}_p(\tilde{A}\tilde{B})_{ij} < \max_{1 \leq k \leq n} (\mathcal{V}_p(\tilde{a}_{ik}) + \mathcal{V}_p(\tilde{b}_{kj})) = (\mathcal{V}_p(\tilde{A}) \otimes \mathcal{V}_p(\tilde{B}))_{ij}$. \square

We are now ready to prove Theorems 3.4 and 3.6.

Proof of Theorem 3.4. Suppose that $\mathcal{A} = \mathcal{V}_p(\tilde{A}) \in \overline{\mathbb{R}}^{n \times n}$ admits max-plus LU factors $\mathcal{L}, \mathcal{U} \in \overline{\mathbb{R}}^{n \times n}$, where $\tilde{A} \in \mathbb{C}\{\{z\}\}^{n \times n}$ is such that $\mathcal{L}(\tilde{A}) \in \mathcal{G}_n$. Then by Theorem 4.5, \tilde{A} has LU factorization $\tilde{A} = \tilde{L}\tilde{U}$ with $\mathcal{V}_p(\tilde{L}) = \mathcal{L}$ and $\mathcal{V}_p(\tilde{U}) = \mathcal{U}$, and by Lemma 4.6 the product $\mathcal{L} \otimes \mathcal{U}$ balances $\mathcal{A} = \mathcal{V}_p(\tilde{A})$. \square

Proof of Theorem 3.6. Let $\tilde{A} \in \mathbb{C}\{\{z\}\}^{n \times n}$ satisfy the conditions in the statement of Theorem 4.5. Now let $\tilde{U}^{(k)} = \tilde{M}_k \cdots \tilde{M}_1 \tilde{A} = \tilde{M}_k \tilde{U}^{(k-1)}$ be the matrix obtained after k steps of Gaussian elimination applied to $\tilde{A} = \tilde{U}^{(0)}$, where $\tilde{M}_k = I - \tilde{m}_k e_k^T$, with e_k the k th unit vector, $\tilde{m}_k = [0, \dots, 0, \tilde{m}_{k+1,k}, \dots, \tilde{m}_{n,k}]^T$, and $\tilde{m}_{ik} = \tilde{u}_{ik}^{(k-1)} / \tilde{u}_{kk}^{(k-1)}$. By Lemma 4.6, the product $\mathcal{V}_p(\tilde{M}_k) \otimes \mathcal{V}_p(\tilde{U}^{(k-1)})$ balances $\mathcal{V}_p(\tilde{U}^{(k)})$, which yields

$$(29) \quad \mathcal{V}_p(\tilde{u}_{ij}^{(k)}) =: u_{ij}^{(k)} \leq \max\{u_{ij}^{(k-1)}, u_{ik}^{(k-1)} + u_{kj}^{(k-1)} - u_{kk}^{(k-1)}\}, \quad i, j \geq k.$$

Next, we show by induction on k that $u_{ij}^{(k)} \leq \max_{p,q} a_{p,q}$ for all i, j , and k . Since $\tilde{U}^{(0)} = \tilde{A}$, $\mathcal{V}_p(\tilde{U}^{(0)}) = \mathcal{V}_p(\tilde{A})$ so that $u_{ij}^{(0)} \leq \max_{p,q} a_{p,q}$ for all i, j . Assume that $u_{ij}^{(\ell)} \leq \max_{p,q} a_{p,q}$ for all i, j and $\ell \leq k - 1$. Since $\mathcal{A} = \mathcal{V}_p(\tilde{A})$ is partial pivoting free, $u_{kk}^{(k-1)} \geq u_{ik}^{(k-1)}$, which combined with (29) and the induction hypothesis gives

$$u_{ij}^{(k)} \leq \max\{u_{ij}^{(k-1)}, u_{kj}^{(k-1)}\} \leq \max_{p,q} a_{p,q}$$

for all i, j . Hence $\varrho_n(\mathcal{A}) = \max_{0 \leq k \leq n-1} (\max_{k \leq i, j \leq n} u_{ij}^{(k)}) - \max_{1 \leq i, j \leq n} a_{ij} \leq 0$. But, by definition, $\varrho_n(\mathcal{A}) \geq 0$, so $\varrho_n(\mathcal{A}) = 0$. \square

5. Max-plus LU algorithm. Computing the max-plus LU factors directly from the formulae (7)–(8) is computationally expensive, as each entry in either the lower part of \mathcal{L} or the upper part of \mathcal{U} requires the computation of two max-plus permanents or, in other words, the solution of two optimal assignment problems. The best known

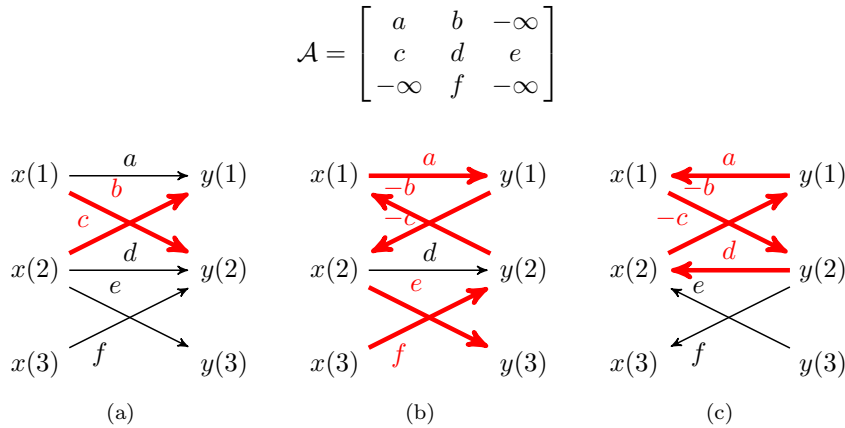


FIG. 1. (a) Bipartite graph G of A with matching $M = \{e_{12}, e_{21}\}$ highlighted with thicker lines. (b) Residual graph $R_G(M)$ with directed path $\sigma = \{e_{32}, e_{12}, e_{11}, e_{21}, e_{23}\}$ highlighted with thicker lines. (c) Transpose residual graph $R_G^T(M)$ with cycle $c = \{e_{12}, e_{22}, e_{21}, e_{11}\}$ highlighted with thicker lines.

algorithms for computing an optimal assignment of $A \in \mathbb{R}_{\max}^{n \times n}$ have worst case cost $O(n\tau + n^2 \log n)$, where τ is the number of nonzeros in A . So the computation of the LU factors using (7)–(8) can cost as much as $O(n^2\tau + n^3 \log n)$ operations. We now describe a more efficient approach, which consists of simultaneously computing all the permanents needed for all the entries in a row of \mathcal{U} or a column of \mathcal{L} , while sharing some of the computation along the way. The bipartite graph setup underpinning our method will be familiar to readers who already have some knowledge of primal dual optimal assignment solvers such as the Hungarian algorithm.

To $A = (a_{ij}) \in \overline{\mathbb{R}}^{n \times n}$ we associate a bipartite graph $G = (X, Y; E)$ with left vertices $X = \{x(1), \dots, x(n)\}$, right vertices $Y = \{y(1), \dots, y(n)\}$, and edge set E . We include an edge e_{ij} in E from $x(i)$ to $y(j)$ with weight $w(e_{ij}) = a_{ij}$ whenever $a_{ij} \neq -\infty$. Thus the edges out of a left vertex $x(i)$ represent the finite entries in the i th row of A (see Figure 1(a)).

A *matching* M is a subset of E with the property that no vertex in M is incident to more than one edge. Vertices which are incident to edges in M are said to be *matched*. The *weight* of a matching $w(M)$ is the sum of its edge weights. Given a matching M , we define the *residual graph* $R_G(M)$ to be the bipartite graph obtained from G by reversing the direction of all of the edges in M and changing the sign of the edges' weights (see Figure 1(b)). Note that we maintain the labelling of edges even when they are reversed. Thus e_{ij} labels either the forward edge from $x(i)$ to $y(j)$ or the backward edge from $y(j)$ to $x(i)$, depending on whether or not it has been reversed. We do not switch to labelling this edge as e_{ji} . We define the weight $w(\sigma)$ of a directed path or cycle σ through $R_G(M)$ to be equal to the sum of the weights of its constitute edges in $R_G(M)$. For the directed path $\sigma = \{e_{32}, e_{12}, e_{11}, e_{21}, e_{23}\}$ through $R_G(M)$ in Figure 1(b), $w(\sigma) = f - b + a - c + e$. We denote by $R_G^T(M)$ the *transpose residual graph* obtained from $R_G(M)$ by reversing the direction of all edges (see Figure 1(c)).

Given a subset S of the edges of $R_G(M)$, we *augment* M according to S , written

as $M\Delta S$, by taking all the edges that appear in either M or S but not both; that is,

$$M\Delta S := \{M \cup S\} \setminus \{M \cap S\}.$$

When we augment with respect to a path/cycle through $R_G(M)$, we treat the path/cycle as a set of edges. For the path $\sigma = \{e_{32}, e_{12}, e_{11}, e_{21}, e_{23}\}$ in Figure 1(b), we have that $M\Delta\sigma = \{e_{11}, e_{23}, e_{32}\}$, which is a matching between the left vertices $\{x(1), x(2), x(3)\}$ and the right vertices $\{y(1), y(2), y(3)\}$ with weight $w(M\Delta\sigma) = a + e + f = w(M) + w(\sigma)$.

Our max-plus LU algorithm relies on the following result.

PROPOSITION 5.1. *Let $\mathcal{A} \in \overline{\mathbb{R}}^{n \times n}$ have max-plus LU factors $\mathcal{L} = (l_{ij})$ and $\mathcal{U} = (u_{ij})$, and for $k = 1, \dots, n$ let $\mathcal{U}^{(k)} = (u_{ij}^{(k)})$ be the intermediate upper factors defined in (16). Let G be the bipartite graph associated with \mathcal{A} with left vertices $X = \{x(1), \dots, x(n)\}$ and right vertices $Y = \{y(1), \dots, y(n)\}$. Then there exists a sequence of matchings M_ℓ between the left vertices $\{x(1), \dots, x(\ell)\}$ and right vertices $\{y(1), \dots, y(\ell)\}$, $\ell = 1, \dots, n$, such that with $M_0 = \emptyset$,*

- (i) $M_k = M_{k-1} \Delta \sigma$, where σ is the maximally weighted path through $R_G(M_{k-1})$ from $x(k)$ to $y(k)$;
- (ii) u_{ki} is either the weight of the maximally weighted path through $R_G(M_{k-1})$ from $x(k)$ to $y(i)$ for $i \geq k$, or $-\infty$ if there is no such path;
- (iii) l_{ik} is either the weight of the maximally weighted path through $R_G^T(M_k)$ from $x(k)$ to $x(i)$ for $i > k$, or $-\infty$ if there is no such path;
- (iv) $u_{ik}^{(k-1)}$ is either the weight of the maximally weighted path through $R_G(M_{k-1})$ from $x(i)$ to $y(k)$ for $i \geq k$, or $-\infty$ if there is no such path.

The proof of Proposition 5.1 is technical and is left to Appendix A. According to Proposition 5.1(i), the sequence of maximally weighted matchings M_1, \dots, M_n can be obtained iteratively starting with $M_0 = \emptyset$. At step $k > 0$, M_{k-1} is augmented with respect to a maximally weighted path through $R_G(M_{k-1})$ from $x(k)$ to $y(k)$ to form M_k . From Proposition 5.1(ii), the entries in the k th row of \mathcal{U} are given by the weights of $n - k$ different maximally weighted paths through $R_G(M_{k-1})$. Likewise, from Proposition 5.1(iii), the entries in the k th columns of \mathcal{L} are given by the weights of $n - k$ different maximally weighted paths through $R_G^T(M_k)$. For a given k , the weights of these maximally weighted paths can be obtained by solving two maximally weighted spanning tree problems rooted at $x(k)$, one through $R_G(M_{k-1})$, the other through $R_G^T(M_k)$. A maximally weighted spanning tree T through $R_G(M)$ rooted at $x(k)$ consists of the maximally weighted paths from $x(k)$ to every reachable left and right vertex. The *depth* of a reachable vertex is the weight of the corresponding maximally weighted path in T . If T does not reach a vertex, then this vertex has depth $-\infty$. All these facts lead to the following algorithm.

ALGORITHM 5.2 (max-plus LU). *Given $\mathcal{A} \in \overline{\mathbb{R}}^{n \times n}$, this algorithm returns a unit lower triangular matrix \mathcal{L} and an upper triangular matrix \mathcal{U} such that \mathcal{L}, \mathcal{U} are max-plus LU factors for \mathcal{A} .*

- % G denotes the bipartite graph associated with \mathcal{A} with left vertices*
% $X = \{x(1), \dots, x(n)\}$ and right vertices $Y = \{y(1), \dots, y(n)\}$.
- 1 Set the lower part of \mathcal{U} and strictly upper part of \mathcal{L} to $-\infty$, and the diagonal entries of \mathcal{L} to 0.
 - 2 Set $M_0 = \emptyset$.
 - 3 for $k = 1 : n$
 - 4 Compute the maximally weighted spanning tree T through $R_G(M_{k-1})$


```

rooted at  $x(k)$ .
5   for  $j = k : n$ 
6        $u_{kj} = \text{depth of } y(j) \text{ in } T$ .
7   end
8    $M_k = M_{k-1} \triangle \sigma$ , where  $\sigma$  is the maximally weighted path through
    $R_G(M_{k-1})$  from  $x(k)$  to  $y(k)$ .
9   Compute the maximally weighted spanning tree  $T'$  through  $R_G^T(M_k)$ 
   rooted at  $x(k)$ .
10  for  $i = k + 1 : n$ 
11       $l_{ik} = \text{depth of } x(k) \text{ in } T'$ .
12  end
13  end
    
```

If \mathcal{A} does not admit max-plus LU factors, then at some step k of Algorithm 5.2, $y(k)$ will have depth $-\infty$ and there will be no path from $x(k)$ to $y(k)$, so the algorithm will not be able to augment the matching M_{k-1} in line 8.

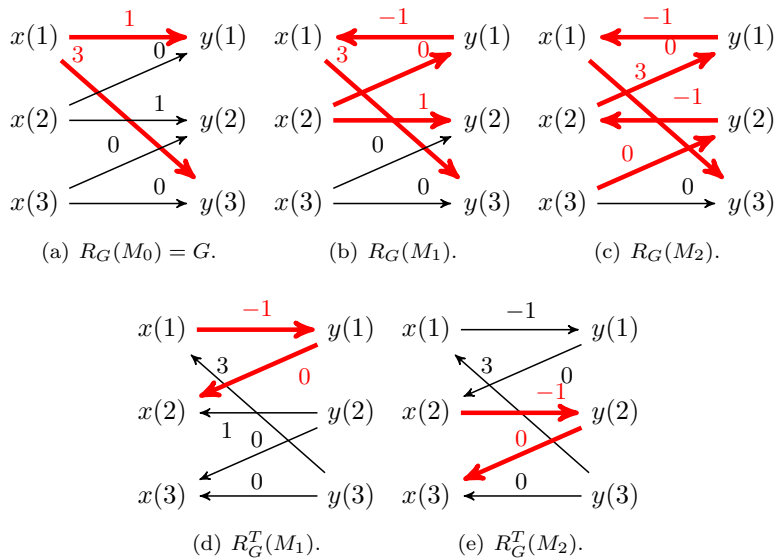


FIG. 2. Bipartite graphs produced by Algorithm 5.2 for the matrix \mathcal{A} of Example 5.3. Maximally weighted spanning trees and paths are highlighted in thicker red lines.

Example 5.3. We apply Algorithm 5.2 to compute the max-plus LU factors for

$$\mathcal{A} = \begin{bmatrix} 1 & -\infty & 3 \\ 0 & 1 & -\infty \\ -\infty & 0 & 0 \end{bmatrix}.$$

$k = 1$. The maximally weighted spanning tree T through $R_G(M_0) = G$ rooted at $x(1)$ is highlighted with thicker red lines in Figure 2(a). The depths of the Y vertices (i.e., 1 for $y(1)$, $-\infty$ for $y(2)$, as it is not reached by the spanning tree, and 3 for $y(3)$) give the entries for the first row of \mathcal{U} .

The maximally weighted path σ through $R_G(M_0)$ from $x(1)$ to $y(1)$ consists of a single edge $\sigma = \{e_{11}\}$ so that $M_1 = M_0 \triangle \sigma = \{e_{11}\}$, yielding the residual graph $R_G(M_1)$ displayed in Figure 2(b).

The maximally weighted spanning tree T' through $R_G^T(M_1)$ rooted at $x(1)$ is highlighted with thicker red lines in Figure 2(d). The depths of the X vertices give the entries for the first column of \mathcal{L} .

$k = 2$. Figure 2(b) highlights the maximally weighted spanning tree T through the residual graph $R_G(M(1))$ rooted at $x(2)$. The depths of the Y vertices give the entries for the second row of \mathcal{U} .

The maximally weighted path σ through $R_G(M_1)$ from $x(2)$ to $y(2)$ consists of a single edge $\sigma = \{e_{22}\}$ (Figure 2(b)) so that $M(2) = M(1) \Delta \sigma = \{e_{11}, e_{22}\}$. The residual graph $R_G(M(2))$ is shown in Figure 2(c). The maximally weighted spanning tree T' through $R_G^T(M(2))$ rooted at $x(2)$ is highlighted in Figure 2(e). The depths of the X vertices give the entries for the second column of \mathcal{L} below the diagonal.

$k = 3$. The maximally weighted spanning tree T through $R_G(M_2)$ rooted at $x(3)$ is highlighted in Figure 2(c). The depth of the $y(3)$ vertex gives the entry for the third row of \mathcal{U} .

The algorithm returns the max-plus LU factors

$$\mathcal{L} = \begin{bmatrix} 0 & -\infty & -\infty \\ -1 & 0 & -\infty \\ -\infty & -1 & 0 \end{bmatrix}, \quad \mathcal{U} = \begin{bmatrix} 1 & -\infty & 3 \\ -\infty & 1 & 2 \\ -\infty & -\infty & 1 \end{bmatrix}.$$

Algorithm 5.2 requires the solution of maximally weighted spanning tree problems in steps 4 and 9. Note that the spanning tree T in step 4 provides the maximally weighted path σ needed in step 8. To efficiently solve the maximally weighted spanning tree problems for a given bipartite graph $G = (X, Y; E)$, we follow an approach taken by Orlin and Lee for the optimal assignment problem [16]. Their approach consists of adjusting the edge weights of G by defining a potential $\phi : X, Y \mapsto \mathbb{R}$ so that, for each edge $e \in E$ from a vertex a to a vertex b , the new edge weight is given by $w'(e) = w(e) - \phi(a) + \phi(b)$ with the property that $w'(e) \leq 0$. This leads to adjusted path weights $w'(\sigma) = w(\sigma) - \phi(a) + \phi(b)$ for a path σ from vertex a to vertex b . Hence if σ is a maximally weighted path from a to b for the original bipartite graph, then it stays maximally weighted for the bipartite graph with adjusted weights, and vice versa. Now since all the adjusted edge weights are nonpositive, Dijkstra's algorithm can then be used to compute the maximally weighted spanning trees and the depth $w'(\sigma)$ of each of its maximally weighted paths σ . Then the depth of σ for the original graph G is given by $w(\sigma) = w'(\sigma) + \phi(a) - \phi(b)$.

The computational cost of adjusting the weights using the technique in [16] is $O(\tau)$ operations, where τ is the number of edges in G . Dijkstra's algorithm solves the maximally weighted spanning tree problem with worst case cost $O(\tau + n \log n)$ for a graph with n vertices and τ edges. Since we need to compute $2n$ of such spanning trees, our max-plus LU algorithm applied to $\mathcal{A} \in \overline{\mathbb{R}}^{n \times n}$ will have worst case cost $O(n\tau + n^2 \log n)$, where τ is the number of finite entries in \mathcal{A} .

5.1. Max-plus LU algorithm for Hungarian matrices. Algorithm 5.2 simplifies if we first apply a Hungarian scaling and an optimal assignment to \mathcal{A} to produce a Hungarian scaled and reordered max-plus matrix \mathcal{H} . In particular, if $\mathcal{A} = \mathcal{V}_c(A)$ with $A \in \mathbb{C}^{n \times n}$, then $\mathcal{H} = \mathcal{V}_c(H)$ with H as in (18) is Hungarian. The next lemma shows that there is no need to compute the sequence of maximally weighted matchings M_1, \dots, M_n (see step 8 of Algorithm 5.2), as these come for free for Hungarian matrices.

LEMMA 5.4. Let $\mathcal{H} \in \overline{\mathbb{R}}^{n \times n}$ be a Hungarian matrix (i.e., $h_{ij} \leq 0$, $h_{ii} = 0$, $1 \leq i, j \leq n$), and let $G = (X, Y; E)$ be the corresponding bipartite graph. Then the sets of edges $M_k = \{e_{11}, \dots, e_{kk}\}$, $k = 1, \dots, n$, are maximally weighted matchings between the left vertices $\{x(1), \dots, x(k)\}$ and the right vertices $\{y(1), \dots, y(k)\}$ for $k = 1, \dots, n$.

Proof. We note that every principal submatrix \mathcal{H}_k of order k of a Hungarian matrix \mathcal{H} is Hungarian. Since \mathcal{H}_k has nonpositive entries, $\sum_{i=1}^k h_{i, \pi(i)} \leq 0$ for any $\pi \in \Pi(k)$. Now for $\pi = \text{id}$, the identity permutation, $\sum_{i=1}^k h_{ii} = 0$ so that $\pi = \text{id}$ is an optimal assignment for \mathcal{H}_k . But an optimal assignment corresponds to a permutation representing a maximally weighted perfect matching between the left and right vertices of the bipartite graph associated with \mathcal{H}_k ; in other words, $M_k = \{e_{11}, \dots, e_{kk}\}$ is a maximally weighted matching between $\{x(1), \dots, x(k)\}$ and $\{y(1), \dots, y(k)\}$. \square

Knowing this sequence of maximally weighted matchings a priori enables us to parallelize Algorithm 5.2. We no longer need to compute the maximally weighted paths through $R_G(M_k)$ before we can form M_{k+1} and compute the maximally weighted paths through $R_G(M_{k+1})$. Instead we can treat each $R_G(M_k)$ separately, computing the k th row of \mathcal{U} and $(k - 1)$ th column of \mathcal{L} in parallel. This approach yields the following algorithm.

ALGORITHM 5.5 (max-plus LU for a Hungarian matrix). *Given a Hungarian matrix $\mathcal{H} \in \overline{\mathbb{R}}^{n \times n}$, this algorithm returns a unit lower triangular matrix \mathcal{L} and an upper triangular matrix \mathcal{U} such that \mathcal{L}, \mathcal{U} are max-plus LU factors for \mathcal{H} .*

```

% G denotes the bipartite graph associated with H with left vertices
% X = {x(1), ..., x(n)} and right vertices Y = {y(1), ..., y(n)}.
1 Set the lower part of U and strictly upper part of L to -∞, and the diagonal
  entries of L to 0.
2 Set M_0 = ∅.
3 for k = 1 : n
4   Set M_k = {e_{11}, ..., e_{kk}}.
5   Compute the maximally weighted spanning tree T through R_G(M_{k-1})
     rooted at x(k).
6   for j = k : n
7     u_{kj} = depth of y(j) in T.
8   end
9   Compute the maximally weighted spanning tree T' through R_G^T(M_k)
     rooted at x(k).
10  for i = k + 1 : n
11    l_{ik} = depth of x(k) in T'.
12  end
13 end
    
```

Because the entries of \mathcal{H} are nonpositive, we can use Dijkstra’s algorithm to compute the maximally weighted spanning trees in steps 5 and 9. Thus our max-plus LU algorithm applied to $\mathcal{H} \in \mathbb{R}_{\max}^{n \times n}$ will have worst case cost $O(n\tau + n^2 \log n)$, where τ is the number of finite entries in \mathcal{H} .

Dijkstra’s algorithm permanently labels vertices in decreasing order of their depth in the tree. This means that when we run Dijkstra’s algorithm to compute the k th row of \mathcal{U} or k th column of \mathcal{L} , we are given the position and value of the entries one at a time in decreasing order of their value. If we are only interested in entries that are greater than some threshold, or if we only want to know the m largest entries in each

row, then we can stop Dijkstra's algorithm earlier and reduce the cost considerably. The exact cost of this implementation will depend heavily on the particular details of the problem matrix. This approach will not work for a non-Hungarian scaled matrix as, while Dijkstra's algorithm will always label vertices in decreasing order, it will be in the order of their adjusted depths. So we cannot be sure that we have found the m largest entries until we have computed all of the adjusted depths and then converted them back into their true depths using the potential.

5.2. Max-plus LU algorithm with partial pivoting. At each step k of Algorithm 5.6, the unmatched left vertices are permuted to maximize the weight of the augmenting path. For this we find a maximally weighted path through $R_G(M_{k-1})$ from $\{x(k), \dots, x(n)\}$ to $y(k)$. If this maximally weighted path begins at $x(i)$, then we interchange $x(i)$ with $x(k)$ and perform step k of Algorithm 5.2. This is analogous to interchanging rows in Gaussian elimination with partial pivoting. Note that we can compute the maximally weighted path through $R_G(M_{k-1})$ from $\{x(k), \dots, x(n)\}$ to $y(k)$ by solving a single maximally weighted spanning tree problem. This is done by adding a root vertex r and connecting it to each unmatched left vertex $x(j)$, $j = k, \dots, n$, with an edge of weight zero (see Figure 3(a)). We then compute the maximally weighted spanning tree through $R_G(M_{k-1}) \cup \{r\}$ rooted at r . Then, rather than choosing the starting vertex of the maximally weighted path from $\{x(k), \dots, x(n)\}$ to $y(k)$, we take the second vertex on the maximally weighted path from r to $y(k)$. This ensures the same choice of vertex, but the algorithm requires us to compute fewer maximally weighted spanning trees. This approach yields the following algorithm.

ALGORITHM 5.6 (max-plus LU with partial pivoting). *Given $\mathcal{A} \in \overline{\mathbb{R}}^{n \times n}$, this algorithm returns a permutation π , a unit lower triangular matrix \mathcal{L} , and an upper triangular matrix \mathcal{U} such that \mathcal{L}, \mathcal{U} are max-plus LU factors for the partial pivoting free matrix $\mathcal{P}_\pi \otimes \mathcal{A}$, where $(\mathcal{P}_\pi)_{ij} = 0$ for $j = \pi(i)$ and $(\mathcal{P}_\pi)_{ij} = -\infty$ otherwise.*

- % G_π denotes the bipartite graph associated with \mathcal{A} with left vertices
- % $X_\pi = \{x(\pi(1)), \dots, x(\pi(n))\}$ and right vertices $Y = \{y(1), \dots, y(n)\}$
- % for some permutation π .
- 1 Set the lower part of \mathcal{U} and strictly upper part of \mathcal{L} to $-\infty$, and the diagonal entries of \mathcal{L} to 0.
- 2 Set $M_0 = \emptyset$ and $\pi = [1, 2, \dots, n]$.
- 3 for $k = 1 : n$
- 4 Add a root vertex r and connect it to each left vertex $x(\pi(j))$,
 $j = k, \dots, n$, by an edge of weight zero.
- 5 Compute the maximally weighted spanning tree T through $R_{G_\pi}(M_{k-1})$ rooted at r .
- 6 Swap $\pi(k)$ with $\pi(i)$, where $x(\pi(i))$ is the 2nd vertex on the maximally weighted path from r to $y(k)$.
- 7 for $j = k, \dots, n$
- 8 $u_{kj} = \text{depth of } y(j) \text{ in } T$.
- 9 end
- 10 $M_k = M_{k-1} \triangle \sigma$, where σ is the maximally weighted path through $R_{G_\pi}(M_{k-1})$ from $x(\pi(k))$ to $y(k)$.
- 11 Compute the maximally weighted spanning tree T' through $R_{G_\pi}^T(M_k)$ rooted at $x(\pi(k))$.
- 12 for $i = k + 1 : n$
- 13 $l_{ik} = \text{depth of } x(\pi(k)) \text{ in } T'$.
- 14 end
- 15 end

Line 6 in Algorithm 5.6 is the row interchanging step. Since we added a root vertex r before the left vertices $X_\pi = \{x(\pi(1)), \dots, x(\pi(n))\}$, we need to look for the second vertex, say $x(\pi(i))$, on the maximally weighted path from r to $y(k)$, the first vertex being r . We then swap $\pi(k)$ with $\pi(i)$, which corresponds to interchanging $x(\pi(k))$ with $x(\pi(i))$.

PROPOSITION 5.7. For $\mathcal{A} \in \overline{\mathbb{R}}^{n \times n}$, Algorithm 5.6 returns a permutation π such that the permuted matrix $\mathcal{P}_\pi \otimes \mathcal{A}$ is partial pivoting free.

Proof. Applying the result of Proposition 5.1(iv) to the matrix $\mathcal{P}_\pi \otimes \mathcal{A}$, we have that $u_{i,k+1}^{(k)}$ is the weight of the maximally weighted path through $R_{G_\pi}(M_k)$ from $x(i)$ to $y(k+1)$ for $i \geq k+1$. But from line 6 in Algorithm 5.6 we are guaranteed because of the interchanging that the maximally weighted path from $\{x(k+1), \dots, x(n)\}$ to $y(k+1)$ is the path from $x(k+1)$ to $y(k+1)$ so that $u_{k+1,k+1}^{(k)} = \max_{k+1 \leq i \leq n} u_{i,k+1}^{(k)}$, $k = 0, \dots, n-1$, which is the definition of max-plus partial pivoting free as given in (15). \square

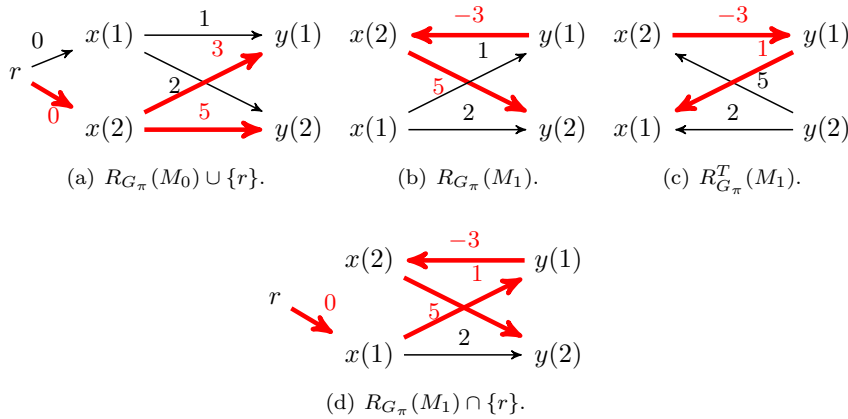


FIG. 3. Bipartite graphs produced by Algorithm 5.6 applied to $\mathcal{A} = \begin{bmatrix} 1 & 2 \\ 3 & 5 \end{bmatrix}$.

Example 5.8. Let us apply Algorithm 5.6 to $\mathcal{A} = \begin{bmatrix} 1 & 2 \\ 3 & 5 \end{bmatrix}$. Let $\pi = [1, 2]$.

$k = 1$. We adjoin a root vertex r to the bipartite graph $R_{G_\pi}(M_0) = G_\pi$ and connect it to the left vertices of $x(\pi(1))$, $x(\pi(2))$ by an edge of weight zero. The maximally weighted spanning tree T through $R_{G_\pi}(M_0) \cup \{r\}$ rooted at r is shown with thicker lines in Figure 3(a). The depths of the Y vertices (i.e., 3 for $y(1)$ and 5 for $y(2)$) define the entries of the first row of \mathcal{U} . Since the maximally weighted path is the one from $x(\pi(2))$, we swap $\pi(1)$ with $\pi(2)$ so that $\pi = [2, 1]$. Then the maximally weighted path σ through $R_{G_\pi}(M_0)$ from $x(\pi(1))$ to $y(1)$ consists of a single edge $\{e_{\pi(1),1}\} = \{e_{21}\}$, so that $M_1 = M_0 \triangle \{e_{21}\} = \{e_{21}\}$, yielding the residual graph $R_{G_\pi}(M_1)$ displayed in Figure 3(b). The maximally weighted spanning tree T' through $R_{G_\pi}^T(M_1)$ rooted at $x(\pi(1)) = x(2)$ is highlighted with thicker red lines in Figure 3(c). The depths of the path from $x(\pi(1))$ to $x(\pi(2))$, i.e., -2 , defines ℓ_{21} .

$k = 2$. We adjoin a root vertex r to the bipartite graph $R_{G_\pi}(M_1)$. Figure 3(d) highlights the maximally weighted spanning tree T through $R_{G_\pi}(M_1) \cup \{r\}$ rooted at r . The depths of the Y vertices define the last row of \mathcal{U} .

Hence we obtain

$$P_\pi \otimes \mathcal{A} = \begin{bmatrix} -\infty & 0 \\ 0 & -\infty \end{bmatrix} \otimes \begin{bmatrix} 1 & 2 \\ 3 & 5 \end{bmatrix} = \begin{bmatrix} 3 & 5 \\ 1 & 2 \end{bmatrix}, \quad \mathcal{L} = \begin{bmatrix} 0 & -\infty \\ -2 & 0 \end{bmatrix}, \quad \mathcal{U} = \begin{bmatrix} 3 & 5 \\ -\infty & 3 \end{bmatrix}.$$

5.3. Max-plus ILU preconditioner. Given the max-plus LU factors \mathcal{L} and \mathcal{U} of $\mathcal{V}_c(A)$, we define the max-plus ILU preconditioner as follows. For a threshold t we store

$$(30) \quad S_{ij} = \begin{cases} 1 & \text{if } l_{ij} \geq \log t + \max_k \log |a_{ik}| \text{ or } u_{ij} \geq \log t + \max_k \log |a_{ik}|, \\ 0 & \text{otherwise.} \end{cases}$$

We then compute the ILU factors for A restricted to positions where S is nonzero using, for example, the general static pattern ILU algorithm described in [17, Alg. 10.1] or the more practical variant [17, Alg. 10.3].

6. Numerical experiments. For our numerical experiments, we select all real nonsymmetric matrices in the University of Florida sparse matrix collection [4] of size $100 \leq n \leq 5000$ that have numeric value symmetry no greater than 0.9, and that are structurally nonsingular. When two matrices from the same group have size n and number of nonzero entries within 1% of each other, then we consider these matrices as duplicate and remove one of them. This leaves us with a total of 260 matrices. We used MATLAB version R2015b to perform the computations. Note that our max-plus LU algorithms are implemented as research codes rather than efficient implementations and, for this reason, we only work with matrices of moderate sizes. For the valuation \mathcal{V}_c we use the logarithm to base 10.

6.1. Stability of Gaussian elimination with no pivoting. The aim of our first set of experiments is to compare the numerical stability of Gaussian elimination with no pivoting applied to Hungarian scaled and reordered matrices H in (18), and to reordered matrices $P_\pi A$, where π is the permutation returned by Algorithm 5.6. For each matrix A in the test set we construct H using the HSL code MC64 [18] and $P_\pi A$ using our MATLAB implementation of Algorithm 5.6. Theorem 3.6 together with Heuristic 3.7 suggests that the growth factors for both H and $P_\pi A$ are of order one since both $\mathcal{V}_c(H)$ and $\mathcal{V}_c(P_\pi A)$ are partial pivoting free. Although this is just a heuristic, we expect the LU factorization with no pivoting of H and $P_\pi A$ to have better numerical stability than for A . To examine the stability of Gaussian elimination on these two classes of matrices we compute the relative backward errors

$$\eta_X = \frac{\|X - \widehat{L}\widehat{U}\|_F}{\|X\|_F}$$

for $X = H$ and $X = P_\pi A$, where \widehat{L} and \widehat{U} are the computed LU factors of the LU factorization of X . We also use Gaussian elimination with no pivoting to compute the LU factorizations of the original matrices A . For each class of matrices, i.e., $X = A, H$ and $P_\pi A$, we plot in Figure 4 the proportion of problems for which we are able to compute LU factors without breakdown and with $\eta_X \leq \alpha$ against α . If the factorization breaks down or if $\eta_X \geq 10^{-1}$, we record a fail. Without pivoting or scaling, the LU factorization fails for almost half of the test matrices A , and $\eta_A \leq 10^{-10}$ for 53% of the test matrices. After applying the max-plus LU permutation P_π to A , the number of failed LU factorizations falls from 120 to 60, and $\eta_{P_\pi A} \leq 10^{-10}$ for 64% of the test matrices. The number of failed LU factorizations is lower for Hungarian scaled matrices H (only 23 fails), and $\eta_H \leq 10^{-10}$ for 86% of the test matrices. Since our aim is to build a new class of ILU preconditioners and because

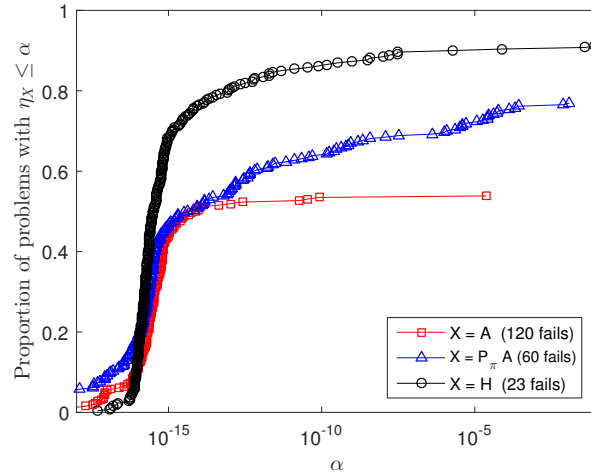


FIG. 4. Proportion of problems with relative backward error $\eta_X \leq \alpha$ for $X = A, H, P_\pi A$.

- for the vast majority of matrices in the test set a reasonably stable LU factorization with no pivoting is possible if A is Hungarian scaled and reordered into H ,
- Hungarian scaling has been shown experimentally to be beneficial for iterative methods [2],
- the max-plus LU algorithm is easier to implement for Hungarian scaled matrices,

from here on we will work with the subset of test problems for which the Hungarian scaled and reordered matrix H can be factorized with no pivoting and with $\eta_H < 0.1$. This results in a test subset of 233 matrices.

6.2. Max-plus LU approximation. The max-plus LU approximation can assist in the computation of an ILU preconditioner by providing a prediction of the positions of larger entries in the LU factors of H . One way of measuring the quality of this prediction is to treat the max-plus LU approximation as a binary classifier. For the LU factorization $H = LU$ of a matrix H from the test set and its max-plus LU factors \mathcal{L}, \mathcal{U} , we predict that $|l_{ij}| \geq 10^{-t}$ for $i > j$ if and only if $l_{ij} \geq -t$, and likewise for the entries of U . The entries of L and U are then labeled as true positive or true negative according to the scheme displayed in Figure 5(a). (For example, for a given t , the (i, j) entry of \mathcal{L} is true positive if $l_{ij} \geq -t$ and $\log_{10} |l_{ij}| \geq -t$.) The accuracy of the classifier is defined by

$$\text{accuracy} = \frac{\text{number of true positives and true negatives}}{\text{number of nonzeros in } L \text{ and } U},$$

and the precision is defined by

$$\text{precision} = \frac{\text{number of true positives}}{\text{number of entries in } L \text{ and } U \text{ such that } |l_{ij}| \geq 10^{-t} \text{ and } |u_{ij}| \geq 10^{-t}}.$$

We also define the soft accuracy and soft precision, which are calculated in the same way using the number of entries in L and U that are soft true positive and soft true negative, where the labelling “soft true positive” and “soft true negative” is done

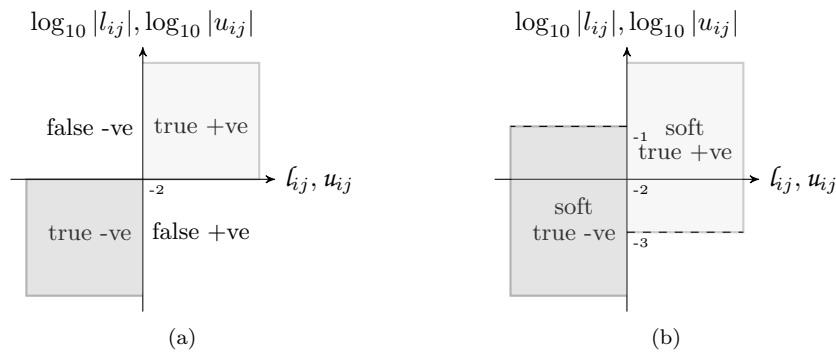


FIG. 5. Labelling of the nonzero entries in the L and U factors for calculating (soft) accuracy and (soft) precision for $t = 2$. The axes cross at $(-2, -2)$.

TABLE 1

Proportion of test problems with with (soft) accuracy and (soft) precision greater than p for $t = 2$.

p	0.8	0.85	0.9	0.95
$A(p)$	85%	83%	79%	71%
$SA(p)$	97%	93%	90%	85%
$P(p)$	86%	83%	80%	59%
$SP(p)$	93%	91%	89%	72%

according to the scheme displayed in Figure 5(b). We record in Table 1 the following proportions:

$A(p)$ = proportion of test examples with accuracy $\geq p$,

$SA(p)$ = proportion of test examples with soft accuracy $\geq p$,

$P(p)$ = proportion of test examples with precision $\geq p$,

$SP(p)$ = proportion of test examples with soft precision $\geq p$,

for $p = 0.80, 0.85, 0.9, 0.95$ and $t = 2$. The scores in Table 1 are quite high and show that, for our test set, the max-plus LU factors provide a good prediction of the position of the larger entries in L and U .

6.3. Behavior of max-plus ILU preconditioning. In this section, we examine the behavior of our max-plus ILU preconditioner on our test set of 233 Hungarian matrices H . We apply GMRES with a right ILU preconditioner to the systems $Hx = b$. For the ILU preconditioner, we use one of the following:

1. Threshold ILU as implemented in the `ilu` function of MATLAB with the option `setup.type = 'crout'`, which performs the Crout version of the threshold ILU factorization (see [17, sec. 10.4.6]) and drop tolerance equal to 10^{-2} , which means that the nonzero and nondiagonal entries of U and L satisfy

$$|u_{ij}| \geq 10^{-2} \|H(i, :)\|_2, \quad |l_{ij}| \geq 10^{-2} \|H(:, j)\|_2 / |u_{jj}|.$$

Since the Hungarian scaled matrices we are working with typically have all row/col norms $\mathcal{O}(1)$, this choice of dropping strategy is roughly equivalent to the max-plus method we describe below.

TABLE 2
 Number of test matrices N for which $ILU(k)$ is used with level k .

Level k	0	1	2	3	4	5	≥ 6
# of test matrices	104	45	28	10	8	4	34

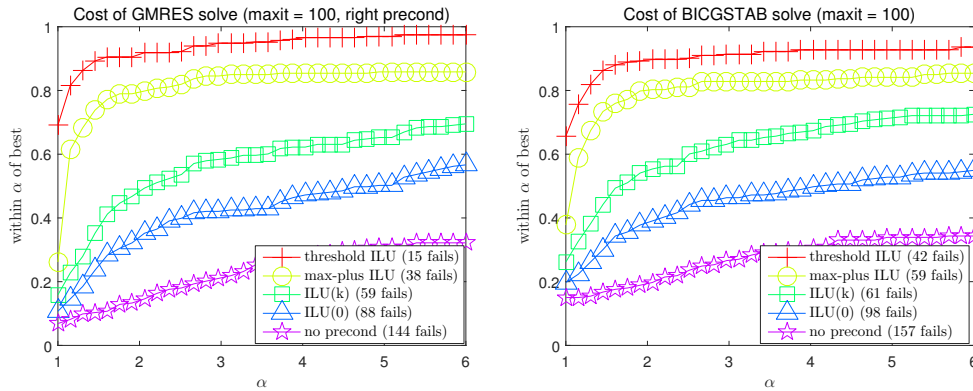


FIG. 6. Performance profile comparing relative costs of GMRES (left plot) and BiCGSTAB (right plot) solves for different ILU preconditioning strategies.

2. $ILU(k)$ through the MATLAB function `iluk` from [13], which is an optimized implementation of the incomplete LU factorization with fill level k described in [17, Chap. 10].
3. $ILU(0)$ with zero level of fill-in.
4. Max-plus ILU by forming the pattern matrix S in (30) with $t = 10^{-2}$ and by calling `iluk` with input parameters set up to bypass the symbolic factorization.

The level k for $ILU(k)$ is chosen as the smallest integer such that the resulting ILU factors are denser than those obtained by the max-plus method. We justify this choice at the end of this section. The distribution of the levels k is shown in Table 2.

We use unrestarted GMRES and stop the iterations when either the approximate solution x satisfies $\|Hx - b\|_2 \leq 10^{-5}\|b\|_2$, in which case we say that GMRES has converged, or the maximum number of iterations has been reached (`maxit=100`), in which case the test matrix is marked as a fail. For each preconditioned system, we record the number of GMRES iterations required for convergence multiplied by the sum of the number of nonzero entries in H and the number of nonzero entries in the ILU factors. For the unpreconditioned systems we record the number of GMRES iterations multiplied by the number of nonzero entries in H . These measures give an approximation of the cost of the GMRES solves but do not include the cost of constructing the preconditioner. The left plot in Figure 6 is a performance profile that compares this cost measure over the different ILU strategies. For the unpreconditioned systems, GMRES fails to converge in less than `maxit = 100` iterations for 144 out of 233 problems and has clearly the worse performance. Systems preconditioned by $ILU(k)$ have a cost more than double that of the best method for about 50% of the problems. The figure shows that the performance of the max-plus ILU preconditioned systems is close to that of the standard threshold ILU preconditioned systems. For about 80% of problems the cost of the max-plus method is within a factor of 2 of

the best method. We obtain an almost identical performance profile plot for GMRES with left ILU preconditioners.

We performed the same set of tests but with BiCGSTAB in place of GMRES for the iterative solver. The corresponding performance profile (see Figure 6(right)) shows that the different ILU preconditioners exhibit the same behavior as for GMRES.

The cost measure that we have used for these performance profiles tends to decrease with the number of nonzero positions included in the ILU factors for all three different ILU techniques. As a result, our choice for the level k is slightly generous to the ILU(k) method. Of course the total cost of solving the linear system should also include the cost of computing the ILU factors, which generally increases with the number of nonzero positions allowed, so that in practice choosing denser ILU factors is not always an advantage. We use this cost measure here for its simplicity. In future work we hope to develop an efficient implementation of our max-plus ILU algorithm that can be timed against state-of-the-art ILU(k) and threshold ILU methods.

7. Conclusion. We have presented a new method for approximating the order of magnitude of the entries in the LU factors of a matrix $A \in \mathbb{C}^{n \times n}$. This method uses max-plus algebra and is based solely on the moduli of the entries in A . If the matrix A is first Hungarian scaled and reordered, then this LU approximation can be computed in parallel with n independent computations, of cost $O(\tau + \log n)$. If we seek only the positions and values of the largest entries in the LU factors, then this cost can be reduced further.

We have shown that this approximation can be used to help compute an ILU preconditioner for A . First we reorder and rescale A to obtain a Hungarian matrix H , then we compute the positions of the largest entries in the LU factors of H , and finally we use these positions as the sparsity pattern for an ILU preconditioner. The resulting preconditioner tends to outperform the comparable ILU(k) preconditioner and have performance very close to a comparable threshold ILU preconditioner when applied as preconditioner to GMRES or BiCGSTAB.

Like the traditional threshold ILU method, the idea behind our max-plus preconditioner is to try to capture the largest entries in the exact LU factors of A . By attempting to find the locations of these large entries a priori using our heuristic max-plus approximation, we trade some accuracy for potential computational advantage. Inaccuracy in the approximation results in slightly less effective preconditioners when compared to the threshold method. But since the columns of the max-plus sparsity pattern can be computed independently of each other and since we can then use the same techniques as in ILU(k) to compute the numerical factors, there is potential to develop a scalable parallel algorithm for max-plus ILU preconditioning.

The numerical examples presented in this paper represent a proof of principal that the max-plus ILU technique can be advantageous in the solution of sparse linear systems.

Appendix A. This appendix presents several technical results needed to prove Proposition 5.1. We refer to section 5 for notation and definitions.

LEMMA A.1. *Let G be a bipartite graph with left vertices $X = \{x(1), \dots, x(n)\}$ and right vertices $Y = \{y(1), \dots, y(n)\}$, and let M be a matching between $\{x(1), \dots, x(k)\}$ and $\{y(1), \dots, y(k)\}$. The following statements hold:*

- (i) *If σ is a direct path through the residual graph $R_G(M)$ from the unmatched left vertex $x(k+1)$ to the unmatched right vertex $y(k+1)$, then $M \Delta \sigma$ is a matching between the left vertices $\{x(1), \dots, x(k), x(k+1)\}$ and the right*

vertices $\{y(1), \dots, y(k), y(k+1)\}$, with weight $w(M\Delta\sigma) = w(M) + w(\sigma)$.

- (ii) If C is a cycle through the residual graph $R_G(M)$, then $M\Delta C$ is a matching between the same vertices as M with weight $w(M\Delta C) = w(M) + w(C)$.
- (iii) If σ is a direct path through the residual graph $R_G(M)$ from the unmatched left vertex $x(k+1)$ to the matched left vertex $x(k)$, then $M\Delta\sigma$ is a matching between the left vertices $\{x(1), \dots, x(k+1)\} \setminus \{x(k)\}$ and the right vertices $\{y(1), \dots, y(k)\}$ with weight $w(M\Delta\sigma) = w(M) + w(\sigma)$.
- (iv) If S_1 and S_2 are disjoint subsets of edges, each either a path or a cycle in $R_G(M)$, then $w(M\Delta(S_1 \cup S_2)) = w(M) + w(S_1) + w(S_2)$.

Proof. (i) There are no edges into any unmatched left vertices or out of any unmatched right vertices, so that σ can only visit its origin vertex $x(k+1)$, destination vertex $y(k+1)$, as well as the vertices matched by M . Since σ is a path from $x(k+1)$ it must include exactly one edge out of this vertex, and since this vertex is unmatched in M there can be no edge out of it in M . Thus $M\Delta\sigma$ contains exactly one edge incident to the origin vertex $x(k+1)$. Likewise for the destination vertex, $M\Delta\sigma$ contains exactly one edge incident to $y(k+1)$.

Let u be a matched left vertex visited by σ . Then σ must include an edge into u , which, being a right-to-left edge, must be an edge in M ; σ must also include an edge out of u , which, being a left-to-right edge, must not be in M . Thus $M\Delta\sigma$ contains exactly one edge incident to u . Likewise if v is a matched right vertex visited by σ , then $M\Delta\sigma$ contains exactly one edge incident to v .

Matched vertices that are not visited by σ are unaffected; likewise unmatched vertices are unaffected. Thus $M\Delta\sigma$ is a subset of E , with exactly one edge incident to each of the left vertices $\{x(1), \dots, x(k), x(k+1)\}$ and each of the right vertices $\{y(1), \dots, y(k), y(k+1)\}$, and no edge incident to any other vertices.

The weight of $M\Delta\sigma$ is equal to $w(M)$ plus the weight of any edges in σ but not in M minus the weight of any edges in M and σ . Since any edge in $M \cap \sigma$ is a backward edge with an extra minus sign and any edge in σ/M is a forwards edge without an extra minus sign, we have $w(M\Delta\sigma) = w(M) + w(\sigma)$.

(ii) The proof of (ii) is the same as for (i) but without the origin or destination vertices.

(iii) The argument is the same as (i) except for the destination vertex $x(k)$. This vertex is incident to exactly one edge in M , which is a right-to-left vertex; since σ ends at this vertex it must also contain this edge, and therefore $M\Delta\sigma$ does not contain an edge incident to $x(k)$.

(iv) Since $S_1 \cup S_2 = \emptyset$ we have $M\Delta(S_1 \cup S_2) = (M\Delta S_1)\Delta S_2$. Hence,

$$w((M\Delta S_1)\Delta S_2) = w(M\Delta S_1) + \hat{w}(S_2),$$

where $\hat{w}(S_2)$ is the weight of the edge set S_2 in the residual graph $R_G(M\Delta S_1)$. However, since S_1 and S_2 are disjoint we have $\hat{w}(S_2) = w(S_2)$, where $w(S_2)$ is the weight of the edge set S_2 in the residual graph $R_G(M)$. (This follows because none of the edges affected by augmenting with respect to S_1 are in the set S_2 .) Therefore we have

$$\begin{aligned} w(M\Delta(S_1 \cup S_2)) &= w((M\Delta S_1)\Delta S_2) \\ &= w(M\Delta S_1) + w(S_2) \\ &= w(M) + w(S_1) + w(S_2). \end{aligned} \quad \square$$

For the residual graph $R_G(M)$ in Figure 1(b) and

- (i) the path $\sigma = \{e_{32}, e_{12}, e_{11}, e_{21}, e_{23}\}$, we have that $M\Delta\sigma = \{e_{11}, e_{23}, e_{32}\}$, which is a matching between the left vertices $\{x(1), x(2), x(3)\}$ and the right vertices $\{y(1), y(2), y(3)\}$ with weight $w(M\Delta\sigma) = a + e + f = w(M) + w(\sigma)$;
- (ii) the cycle $c = \{e_{22}, e_{12}, e_{11}, e_{21}\}$ through $R_G(M)$ with weight $w(c) = d - b + a - c$, we have that $M\Delta c = \{e_{11}, e_{22}\}$, which is a matching between the left vertices $\{x(1), x(2)\}$ and the right vertices $\{y(1), y(2)\}$ with weight $w(M\Delta c) = a + d = w(M) + w(c)$;
- (iii) the path $\sigma = \{e_{32}, e_{12}, e_{11}, e_{21}\}$ through $R_G(M)$ with weight $w(\sigma) = f - b + a - c$, we have that $M\Delta\sigma = \{e_{11}, e_{32}\}$, which is a matching between the left vertices $\{x(1), x(3)\}$ and the right vertices $\{y(1), y(2)\}$, with weight $w(M\Delta\sigma) = a + f = w(M) + w(\sigma)$.

LEMMA A.2. Let G be the bipartite graph associated with $\mathcal{A} \in \overline{\mathbb{R}}^{n \times n}$ with left vertices $\{x(1), \dots, x(n)\}$ and right vertices $\{y(1), \dots, y(n)\}$. Then $\text{perm}(\mathcal{A}([i_1, \dots, i_k], [j_1, \dots, j_k]))$ is the weight of the maximally weighted matching between the left vertices $\{x(i_1), \dots, x(i_k)\}$ and the right vertices $\{y(j_1), \dots, y(j_k)\}$.

Proof. Any matching M between the left vertices $\{x(1), \dots, x(k)\}$ and the right vertices $\{y(1), \dots, y(k)\}$ can be represented by a unique permutation $\pi \in \Pi(k)$ so that M_π is the matching that matches $x(i)$ to $y(\pi(i))$. Now consider $w(M_\pi) = \sum_{i=1}^k a_{x(i)y(\pi(i))}$ so that the weight of the maximally weighted matching is given by

$$\max_{\pi \in \Pi(k)} w(M_\pi) = \max_{\pi \in \Pi(k)} \sum_{i=1}^k a_{x(i)y(\pi(i))} = \text{perm}(\mathcal{A}([i_1, \dots, i_k], [j_1, \dots, j_k])). \quad \square$$

LEMMA A.3. Let $G = (X, Y; E)$ be a bipartite graph, and let $M_\pi \subset E$, $M_\omega \subset E$ be matchings defined by the permutations $\pi \in \Pi(k)$ and $\omega \in \Pi(k+1)$, respectively. Then there exists a path σ through $R_G(M)$ from $x(k+1)$ to $y(k+1)$ as well as disjoint cycles C_1, \dots, C_m through $R_G(M)$ such that $M_\omega = M_\pi \Delta (\sigma \cup C_1 \cup \dots \cup C_m)$.

Proof. We will prove the lemma by constructing the path and cycles as follows. Set $\sigma(1) = x(k+1)$, then set

$$\begin{aligned} \sigma(2) &= y(\omega(k+1)), & \sigma(3) &= x(\pi^{-1}\omega(k+1)), \\ \sigma(4) &= y(\omega\pi^{-1}\omega(k+1)), & \sigma(5) &= x(\pi^{-1}\omega\pi^{-1}\omega(k+1)), \dots \end{aligned}$$

There is no right vertex $y(j)$ with $\pi^{-1}(j) = k+1$ since the domain of π is $\{1, \dots, k\}$. Also all subsequent vertices visited by the constructed path can have only one predecessor, as π and ω are permutations. Therefore σ cannot contain any cycle and must terminate since there are only finitely many vertices that it can visit without repetition. The only way that the path can terminate is by reaching a vertex where the next step is not well defined, and the only such vertex is $y(k+1)$. The constructed path therefore terminates at $\sigma(2\ell) = y(k+1)$.

Next we pick any left vertex x matched by M_π , which is not visited by σ . We construct a new path starting from x with the same rule that we used for constructing σ . Since there are no possible termination points, where π^{-1} or ω are not defined, this new path must be cyclic. If the constructed cycle is of length 2, then we discount the cycle but still record the constituent vertices as having been visited. We construct further cycles C_1, \dots, C_m until all matched vertices have been visited.

By construction each vertex matched by M_π either has the same matching under M_ω or is incident to two edges in $\sigma \cup C_1 \cup \dots \cup C_m$, one edge from M_π and one from

M_ω . When we augment by taking the symmetric difference, the edge from M_π is replaced by the one from M_ω . Likewise the origin and destination vertices are each incident to an edge which is in M_ω but not in M_π , so these edges are also included when we augment. \square

The following theorem shows us how we can obtain a sequence of maximally weighted matchings by augmenting with respect to maximally weighted paths through the residual graph. This is the mechanism by which we will compute all of the entries in the max-plus LU factors.

THEOREM A.4. *Let $G = (X, Y; E)$ be a bipartite graph, and let M be a maximally weighted matching between the left vertices $\{x(1), \dots, x(k)\}$ and the right vertices $\{y(1), \dots, y(k)\}$. The following hold:*

- (i) *If σ is a maximally weighted path through $R_G(M)$ from the unmatched left vertex $x(k+1)$ to the unmatched right vertex $y(k+1)$, then $M \Delta \sigma$ is a maximally weighted matching between the left vertices $\{x(1), \dots, x(k+1)\}$ and the right vertices $\{y(1), \dots, y(k+1)\}$.*
- (ii) *If σ be a maximally weighted path through $R_G(M)$ from the unmatched left vertex $x(k+1)$ to the matched left vertex $x(k)$, then $M \Delta \sigma$ is a maximally weighted matching between the left vertices $\{x(1), \dots, x(k-1), x(k+1)\}$ and the right vertices $\{y(k), \dots, y(k)\}$.*

Proof. (i) Let σ be a path through $R_G(M)$, and let C_1, \dots, C_m be disjoint cycles in $R_G(M)$. Then by Lemma A.1(iv), $w(M \Delta (\sigma \cup C_1 \cup \dots \cup C_m)) = w(M) + w(\sigma) + w(C_1) + \dots + w(C_m)$. It follows from Lemma A.1(ii) that if C is single cycle, then $M \Delta C$ is a matching between the same vertices as M , and since M is the maximally weighted matching on its matched vertices, $w(M \Delta C) = w(M) + w(C) \leq w(M)$, and thus $w(C) \leq 0$, showing that any cycle in $R_G(M)$ must have nonpositive weight.

Now consider $\max_{M'} w(M')$, where the maximum is taken over all matchings from $\{x(1), \dots, x(k+1)\}$ to $\{y(k), \dots, y(k+1)\}$. Lemma A.3 tells us that every matching M' from $\{x(1), \dots, x(k+1)\}$ to $\{y(k), \dots, y(k+1)\}$ can be written as the augmentation of the matching M with respect to some path and cycles. Therefore

$$\max_{M'} w(M') = w(M) + \max_{\sigma, C_1, \dots, C_m} w(\sigma) + w(C_1) + \dots + w(C_m),$$

where the maximum is taken over all paths through $R_G(M)$ from $x(k+1)$ to $y(k+1)$ and disjoint cycles C_1, \dots, C_m in $R_G(M)$. Since the cycle weights are all nonpositive we have

$$\max_{M'} w(M') \leq w(M) + \max_{\sigma} w(\sigma),$$

and since the upper bound is attained by the matching $M' = M \Delta \arg \max_{\sigma} w(\sigma)$, where $\arg \max_{\sigma} w(\sigma)$ is a maximally weighted path through $R_G(M)$ from $x(k+1)$ to $y(k+1)$, we have

$$\max_{M'} w(M') = w \left(M \Delta \arg \max_{\sigma} w(\sigma) \right).$$

Hence $M \Delta \arg \max_{\sigma} w(\sigma)$ is a maximally weighted matching between the left vertices $\{x(1), \dots, x(k+1)\}$ and the right vertices $\{y(1), \dots, y(k+1)\}$.

- (ii) This follows from Lemmas A.1 and A.3 in analogy to Theorem A.4. \square

We are now ready to prove Proposition 5.1.

Proof of Proposition 5.1. Statement (i) is simply Theorem A.4(i).

Statement (ii): Let M_{k-1} be a maximally weighted matching between the left vertices $\{x(1), \dots, x(k-1)\}$ and the right vertices $\{y(1), \dots, y(k-1)\}$, and let σ

be a maximally weighted path through $R_G(M_{k-1})$ from the unmatched left vertex $x(k)$ to the unmatched right vertex $y(j)$ for $j \geq k$. Then from Lemma A.2 we have $\text{perm}(\mathcal{A}(1 : k - 1, 1 : k - 1)) = w(M_{k-1})$. From Theorem A.4 we have that $M_{k-1} \Delta \sigma$ is the maximally weighted matching between the left vertices $\{x(1), \dots, x(k)\}$ and the right vertices $\{y(1), \dots, y(k-1), y(j)\}$. Thus from Lemma A.2 we have $\text{perm}(\mathcal{A}(1 : k, [1 : k - 1, j])) = w(M_{k-1} \Delta \sigma)$. Finally, using Lemma A.1, we have $w(M_{k-1} \Delta \sigma) = w(M_{k-1}) + w(\sigma)$, and using the expression for u_{kj} for $j \geq k$ in (7), we have that

$$\begin{aligned} u_{kj} &= \text{perm}(\mathcal{A}(1 : k, [1 : k - 1, j])) - \text{perm}(\mathcal{A}(1 : k - 1, 1 : k - 1)) \\ &= w(M_{k-1}) + w(\sigma) - w(M_{k-1}) = w(\sigma). \end{aligned}$$

Statement (iii): Let M_k be a maximally weighted matching between the left vertices $\{x(1), \dots, x(k)\}$ and the right vertices $\{y(1), \dots, y(k)\}$, and let σ be a maximally weighted path through $R_G(M_k)$ from the unmatched left vertex $x(i)$ with $i > k$ to the matched left vertex $x(k)$. Then from Lemma A.2 we have $\text{perm}(\mathcal{A}(1 : k, 1 : k)) = w(M_k)$. From Theorem A.4 we have that $M_k \Delta \sigma$ is the maximally weighted matching between the left vertices $\{x(1), \dots, x(k-1), x(i)\}$ and the right vertices $\{y(1), \dots, y(k)\}$. Thus from Lemma A.2 we have $\text{perm}(\mathcal{A}([1 : k - 1, i], 1 : k)) = w(M_k \Delta \sigma)$. Finally, using Lemma A.1, we have $w(M_k \Delta \sigma) = w(M_k) + w(\sigma)$, and using the expression for l_{ik} for $i > k$ in (7), we have that

$$\begin{aligned} l_{ik} &= \text{perm}(\mathcal{A}([1 : k - 1, i], 1 : k)) - \text{perm}(\mathcal{A}(1 : k, 1 : k)) \\ &= w(M_k) + w(\sigma) - w(M_k) = w(\sigma). \end{aligned}$$

Statement (iv): From (16) we have

$$(31) \quad u_{ik}^{(k-1)} = \text{perm}(\mathcal{A}([1 : k - 1, i], [1 : k])) - \text{perm}(\mathcal{A}(1 : k - 1, 1 : k - 1)), \quad i \geq k.$$

Let M_{k-1} be a maximally weighted matching between the left vertices $\{x(1), \dots, x(k-1)\}$ and the right vertices $\{y(1), \dots, y(k-1)\}$, and let σ be a maximally weighted path through $R_G(M_{k-1})$ from the unmatched left vertex $x(i)$ with $i \leq k$ to the unmatched right vertex $y(k)$. Then from Lemma A.2 we have

$$\text{perm}(\mathcal{A}(1 : k - 1, 1 : k - 1)) = w(M_{k-1}).$$

From Theorem A.4 we have that $M_{k-1} \Delta \sigma$ is the maximally weighted matching between the left vertices $\{x(1), \dots, x(k-1), x(i)\}$ and the right vertices $\{y(1), \dots, y(k-1), y(k)\}$. Thus from Lemma A.2 we have

$$\text{perm}(\mathcal{A}([1 : k - 1, i], [1 : k])) = w(M_{k-1} \Delta \sigma).$$

Finally, using Lemma A.1, we have $w(M_{k-1} \Delta \sigma) = w(M_{k-1}) + w(\sigma)$, and substituting back into (31), we have that

$$\begin{aligned} u_{ij}^{(k-1)} &= \text{perm}(\mathcal{A}([1 : k - 1, i], [1 : k])) - \text{perm}(\mathcal{A}(1 : k - 1, 1 : k - 1)) \\ &= w(M_{k-1}) + w(\sigma) - w(M_{k-1}) = w(\sigma). \quad \square \end{aligned}$$

REFERENCES

- [1] M. AKIAN, S. GAUBERT, AND A. MARCHESINI, *Tropical bounds for eigenvalues of matrices*, Linear Algebra Appl., 446 (2014), pp. 281–303, <https://doi.org/10.1016/j.laa.2013.12.021>.

- [2] M. BENZI, J. C. HAWS, AND M. TUMA, *Preconditioning highly indefinite and nonsymmetric matrices*, SIAM J. Sci. Comput., 22 (2000), pp. 1333–1353, <https://doi.org/10.1137/S1064827599361308>.
- [3] D. A. BINI AND V. NOFERINI, *Solving polynomial eigenvalue problems by means of the Ehrlich-Aberth method*, Linear Algebra Appl., 439 (2013), pp. 1130–1149, <https://doi.org/10.1016/j.laa.2013.02.024>.
- [4] T. A. DAVIS AND Y. HU, *The University of Florida sparse matrix collection*, ACM Trans. Math. Software, 38 (2011), 1, <https://doi.org/10.1145/2049662.2049663>.
- [5] F. R. GANTMACHER, *The Theory of Matrices*, Vol. 1, Chelsea, New York, 1959; ISBN 0-8284-0131-4.
- [6] S. GAUBERT AND M. SHARIFY, *Tropical scaling of polynomial matrices*, in Positive Systems, Lecture Notes in Control Inform. Sci. 389, Springer-Verlag, Berlin, 2009, pp. 291–303.
- [7] S. HAMMARLING, C. J. MUNRO, AND F. TISSEUR, *An algorithm for the complete solution of quadratic eigenvalue problems*, ACM Trans. Math. Software, 39 (2013), 18, <https://doi.org/10.1145/2450153.2450156>.
- [8] N. J. HIGHAM, *Accuracy and Stability of Numerical Algorithms*, 2nd ed., SIAM, Philadelphia, 2002, <https://doi.org/10.1137/1.9780898718027>.
- [9] J. HOOK, *Max-plus singular values*, Linear Algebra Appl., 486 (2015), pp. 419–442, [https://doi.org/10.1016/0024-3795\(95\)00455-6](https://doi.org/10.1016/0024-3795(95)00455-6).
- [10] J. HOOK AND F. TISSEUR, *Max-plus eigenvalues and singular values: A useful tool in numerical linear algebra*, manuscript, 2016.
- [11] A. S. HOUSEHOLDER, *The Theory of Matrices in Numerical Analysis*, Blaisdell, New York, 1964; ISBN 0-486-61781-5.
- [12] D. HYSOM AND A. POTHEN, *A scalable parallel algorithm for incomplete factor preconditioning*, SIAM J. Sci. Comput., 22 (2001), pp. 2194–2215, <https://doi.org/10.1137/S1064827500376193>.
- [13] K. MILLER, *ILU(k) Preconditioner*, <https://uk.mathworks.com/matlabcentral/fileexchange/48320-ilu-k-preconditioner>.
- [14] V. NOFERINI, M. SHARIFY, AND F. TISSEUR, *Tropical roots as approximations to eigenvalues of matrix polynomials*, SIAM J. Matrix Anal. Appl., 36 (2015), pp. 138–157, <https://doi.org/10.1137/14096637X>.
- [15] M. OLSCHOWKA AND A. NEUMAIER, *A new pivoting strategy for Gaussian elimination*, Linear Algebra Appl., 240 (1996), pp. 131–151, [https://doi.org/10.1016/0024-3795\(94\)00192-8](https://doi.org/10.1016/0024-3795(94)00192-8).
- [16] J. B. ORLIN AND Y. LEE, *Quickmatch—A Very Fast Algorithm for the Assignment Problem*, Working papers 3547-93, Sloan School of Management, Massachusetts Institute of Technology (MIT), Cambridge, MA, 1993.
- [17] Y. SAAD, *Iterative Methods for Sparse Linear Systems*, 2nd ed., SIAM, Philadelphia, 2003, <https://doi.org/10.1137/1.9780898718003>.
- [18] SCIENCE AND TECHNOLOGY FACILITIES COUNCIL, *The HSL Mathematical Software Library*, Fortran codes for large scale scientific computation, <http://www.hsl.rl.ac.uk/>.
- [19] M. SHARIFY, *Scaling Algorithms and Tropical Methods in Numerical Matrix Analysis: Application to the Optimal Assignment Problem and to the Accurate Computation of Eigenvalues*, Ph.D. thesis, Ecole Polytechnique, Palaiseau, France, 2011; available from <http://hal.archives-ouvertes.fr/docs/00/64/38/36/PDF/thesis.pdf>.
- [20] M. VAN BAREL AND F. TISSEUR, *Polynomial eigenvalue solver based on tropically scaled Lagrange linearization*, Linear Algebra Appl., 2017, to appear, <https://doi.org/10.1016/j.laa.2017.04.025>.