

LOW-RANK APPROXIMATION IN THE FROBENIUS NORM BY COLUMN AND ROW SUBSET SELECTION*

ALICE CORTINOVIS[†] AND DANIEL KRESSNER[†]

Abstract. A CUR approximation of a matrix A is a particular type of low-rank approximation $A \approx CUR$, where C and R consist of columns and rows of A , respectively. One way to obtain such an approximation is to apply column subset selection to A and A^T . In this work, we describe a numerically robust and much faster variant of the column subset selection algorithm proposed by Deshpande and Rademacher, which guarantees an error close to the best approximation error in the Frobenius norm. For cross approximation, in which U is required to be the inverse of a submatrix of A described by the intersection of C and R , we obtain a new algorithm with an error bound that stays within a factor $k + 1$ of the best rank- k approximation error in the Frobenius norm. To the best of our knowledge, this is the first deterministic polynomial-time algorithm for which this factor is bounded by a polynomial in k . Our derivation and analysis of the algorithm is based on derandomizing a recent existence result by Zamarashkin and Osinsky. To illustrate the versatility of our new column subset selection algorithm, an extension to low multilinear rank approximations of tensors is provided as well.

Key words. column subset selection, cross approximation, low-rank approximation, tensors

AMS subject classifications. 15A23, 65F30, 65F99

DOI. 10.1137/19M1281848

1. Introduction. Given an $m \times n$ matrix A and an integer k , typically much smaller than m and n , the *column subset selection problem* aims at determining an index set $I \subset \{1, \dots, n\}$ of cardinality k such that the corresponding k columns $A(:, I)^1$ represent a good approximation of the range of A . This problem has broad applications in a diversity of disciplines, such as scientific computing, model reduction, and statistical data analysis. Compared to other approximations, such as the one obtained from applying a singular value decomposition (SVD) to A , choosing a subset of columns may lend itself to enhanced interpretability and structure preservation. Examples include unsupervised feature selection [5, 30], sensor selection [27], and data mining [19].

While column subset selection is a classical problem in numerical linear algebra, closely connected to rank-revealing QR factorizations [9, 10, 24], new significant theoretical and algorithmic developments have been achieved during the last two decades within the model reduction and theory of algorithms communities. In particular, this concerns the interplay between column subset selection and interpolation [3, 11] as well as the development and analysis of greedy and randomized algorithms; see [6, 17, 18, 20, 31, 40] for a few references representing this research direction.

This paper is concerned with algorithmic improvements and extensions of the seminal work by Deshpande, Rademacher, and coauthors [14, 15] on column subset

*Received by the editors August 19, 2019; accepted for publication (in revised form) August 5, 2020; published electronically November 9, 2020.

<https://doi.org/10.1137/19M1281848>

Funding: The work of the first author was supported by the Swiss National Science Foundation research project “Fast Algorithms from Low-Rank Updates,” grant 200020_178806.

[†]Institute of Mathematics, EPF Lausanne, 1015 Lausanne, Switzerland (alice.cortinovic@epfl.ch, daniel.kressner@epfl.ch).

¹We use MATLAB’s colon notation $A(:, I)$ to indicate the submatrix of A corresponding to all row indices and column indices in I .

selection. In [15] the existence of an index set I such that

$$(1) \quad \|A - A(:, I)A(:, I)^\dagger A\|_F^2 \leq (k+1) \left(\sigma_{k+1}^2(A) + \cdots + \sigma_{\min\{m, n\}}^2(A) \right)$$

has been established. Here, $\|\cdot\|_F$ and $(\cdot)^\dagger$ denote the Frobenius norm and the Moore–Penrose inverse of a matrix, respectively. We let $\sigma_1(A) \geq \sigma_2(A) \geq \cdots$ denote the singular values of A . Note that $A(:, I)A(:, I)^\dagger$ is an orthogonal projector and the bound (1) measures how well all the columns of A are approximated by the subset of columns contained in I . The bound (1) is remarkable because the SVD of A implies that the best approximation error $\|A - QQ^\dagger A\|_F^2$ attained by an *arbitrary* $m \times k$ matrix Q is given by $\sigma_{k+1}^2(A) + \cdots + \sigma_{\min\{m, n\}}^2(A)$. The bound (1) is larger by a factor that is only linear in k . More generally, we will call any quasi-optimal bound with a factor that is at most polynomial in k (and independent of m, n , or A) a *polynomial bound*. The proof of (1) proceeds by defining a suitable discrete probability distribution on index tuples such that the expected value of the error with respect to this distribution satisfies the bound. This then implies the existence of at least one index set satisfying the bound as well. We remark that the factor $k+1$ in (1) cannot be improved [15, Proposition 3.3]. In [14], a deterministic algorithm has been developed by derandomizing this approach using the method of conditional expectations. These conditional expectations are given in terms of coefficients of certain characteristic polynomials, and the algorithm from [14] attains efficiency by cheaply updating these coefficients. However, it is well known that working with characteristic polynomials in finite precision arithmetic is prone to massive numerical cancellation [34] and, as we will see, the algorithm from [14] is also affected by numerical instability. Our first contribution, presented in section 2, consists of deriving a formulation of the algorithm that updates singular values instead of coefficients of characteristic polynomials. While our new variant enjoys the same favorable complexity, numerical experiments with matrices of different singular value decay indicate that it is numerically robust, achieving (1) even when the right-hand side is at the level of unit roundoff. Based on a minor extension of the theory from [14, 15], we will also present a modification of the column selection strategy that results in significant speedups of the algorithm.

In section 3, we extend the developments from [14] to the problem of determining a rank- k approximation of the form

$$A \approx CUR,$$

where $C = A(:, J)$ and $R = A(I, :)$ contain k selected columns and rows of A , respectively. There is a simple and well-established strategy for deriving such an approximation; see, e.g., [17, 36]: One first applies column subset selection to A and A^T in order to determine C and R , respectively. Given C and R , the choice

$$(2) \quad U = C^\dagger A R^\dagger$$

then minimizes the Frobenius norm error. We will show that this strategy combined with (1) results in an error that is at most a factor $\sqrt{2k+2}$ larger than the best rank- k approximation error. While this is clearly a favorable bound, the choice (2) comes with a disadvantage. Constructing an approximation of the form (2) involves matrix-vector products with A , but this may be too costly when the full matrix A is not readily available, and cheap, possibly heuristic strategies are available for determining I, J . One example for such a situation is the Chebfun2 construction for approximating bivariate functions [38, section 2.1], which uses a coarse discretization

to cheaply determine I, J and then evaluates the full matrix A only along the cross containing the rows and columns determined by I and J , respectively. The choice $U = A(I, J)^{-1}$ then leads to a rank- k approximation of the form

$$A \approx A(:, J)A(I, J)^{-1}A(I, :),$$

which is often called cross approximation. Choosing I, J via column subset selection is not advisable in this setting; it may lead to (nearly) singular $A(I, J)$ and result in an unfavorable approximation error. On the other hand, Goreinov and Tyrtshnikov [22] have established a polynomial bound for cross approximation in the maximum norm when choosing I, J such that the volume of $A(I, J)$ is maximal. Recently, Zamarashkin and Osinsky [42] derived a polynomial bound in the Frobenius norm by extending the techniques from [15]. However, as far as we know, there is no polynomial-time deterministic algorithm that guarantees a polynomial bound (in any norm); popular greedy algorithms lead to exponential bounds [12, 25] at best. One major contribution of this work is to derive such an algorithm via an extension of [14]; our algorithm guarantees a Frobenius norm error that is at most a factor $k + 1$ larger than the best approximation error.

Section 4 contains an extension to the Tucker decomposition of tensors, which is suitable for tensors of low order. In particular, we derive a deterministic algorithm that obtains a multilinear low-rank approximation that is constructed from the fibers of the tensor and satisfies a polynomial bound. Although our approach is a relatively straightforward extension of (2) and related approaches have been proposed in the literature [16, 23, 35], we are not aware that such an algorithm has been explicitly spelled out and analyzed.

2. Column subset selection. We start by providing more details on the approach from [14, 15] for the column subset selection problem. In the following we consider a matrix $A \in \mathbb{R}^{m \times n}$ with $m \leq n$ and rank at least k . We let a_i denote the i th column of A , and $\pi_{i_1, \dots, i_k} A$ the orthogonal projection of A on the subspace spanned by the columns a_{i_1}, \dots, a_{i_k} , that is,

$$\pi_{i_1, \dots, i_k} A := A(:, I) \cdot A(:, I)^\dagger \cdot A = QQ^T A,$$

where $I = (i_1, \dots, i_k) \in \{1, \dots, n\}^k$ and Q denotes an orthonormal basis of $A(:, I)$. Let us emphasize that I is now a *tuple*. Although order is not important and we are ultimately interested in an index *set*, working with tuples simplifies the subsequent definition and manipulation of probability distributions. The *volume* of a rectangular matrix $B \in \mathbb{R}^{m \times k}$ with $k \leq m$ is defined as $\text{Vol}(B) := \prod_{i=1}^k \sigma_i(B)$. Note that $\text{Vol}^2(B) = \det(B^T B)$.

We now define a discrete probability distribution on integer tuples of the form $I \in \{1, \dots, n\}^k$ corresponding to a selection of k columns from A . For this purpose, let $X = (X_1, \dots, X_k)$ be a k -tuple of random variables with values in $\{1, \dots, n\}$ such that

$$(3) \quad \mathbb{P}(X = I) := \frac{\text{Vol}^2(A(:, I))}{\sum_{J \in \{1, \dots, n\}^k} \text{Vol}^2(A(:, J))}.$$

It follows from the definition that $\text{Vol}(A(:, I)) = 0$ whenever i_1, \dots, i_k contain repeated indices. Then [15, Theorem 1.3] shows that

$$(4) \quad \mathbb{E}[\|A - \pi_{X_1, \dots, X_k} A\|_F^2] \leq (k + 1) (\sigma_{k+1}^2 + \dots + \sigma_m^2).$$

In particular, this implies the existence of I satisfying this bound.

In view of (3) and the prominent role played by maximum volume submatrices in low-rank approximation [22], it is tempting to expect that the k columns of maximum volume satisfy (1). However, choosing such columns is not only an NP-hard problem [8], but these might also fail to satisfy (1). For instance, for $k = 1$ consider the $2 \times n$ matrix

$$A = \begin{bmatrix} a(1+\varepsilon) & b & b & \dots & b \\ -b(1+\varepsilon) & a & a & \dots & a \end{bmatrix},$$

with $a^2 + b^2 = 1$ and $\varepsilon > 0$. The column of maximum volume (that is, of maximum Euclidean norm) is the first one. The approximation error obtained by this choice is given by $\|A - \pi_1 A\|_F^2 = n - 1$, which is much larger than $2\sigma_2^2 = 2(1+\varepsilon)^2$ for ε sufficiently small. Note that choosing any of the other columns yields the best approximation error $(1+\varepsilon)^2 = \sigma_2^2$.

2.1. Algorithm by Deshpande and Rademacher. Deshpande and Rademacher [14] derived a deterministic algorithm for column subset selection by derandomizing (4) using the method of conditional expectations.

More specifically, the first step of the algorithm chooses an index i_1 such that

$$\mathbb{E} [\|A - \pi_{X_1, \dots, X_k} A\|_F^2 \mid X_1 = i_1]$$

is minimized. By construction, this quantity still satisfies the bound (4). More generally, having $t-1$ indices i_1, \dots, i_{t-1} selected, step t chooses an index i_t such that

$$(5) \quad \mathbb{E} [\|A - \pi_{X_1, \dots, X_k} A\|_F^2 \mid X_1 = i_1, \dots, X_{t-1} = i_{t-1}, X_t = i_t]$$

is minimized. After k steps we arrive at an index set I of cardinality k such that the desired bound (1) holds.

For the algorithm to be practical, it is crucial to compute the conditional expectations (5) efficiently. Lemma 21 in [14] shows that

$$\mathbb{E} [\|A - \pi_{X_1, \dots, X_k} A\|_F^2 \mid X_1 = i_1, \dots, X_t = i_t] = (k-t+1) \frac{c_{m-k+t-1}(BB^T)}{c_{m-k+t}(BB^T)},$$

where the right-hand side involves the matrix $B = A - \pi_{i_1, \dots, i_t} A$ and coefficients $c_j \equiv c_j(BB^T)$ of the characteristic polynomial

$$(6) \quad (-\lambda)^m + c_{m-1}(-\lambda)^{m-1} + \dots + c_1(-\lambda) + c_0 := \det(BB^T - \lambda I).$$

It is therefore required to compute in every step for all values of i , the ratios

$$(7) \quad \frac{c_{m-k+t-1}(B_i B_i^T)}{c_{m-k+t}(B_i B_i^T)},$$

where $B_i = A - \pi_{i_1, \dots, i_{t-1}, i} A$.

In the following, we discuss the computation of (7) and show how the minimization problem (5) can be relaxed in order to accelerate the search for suitable indices.

2.2. Computation of characteristic polynomial coefficients. Assuming that the first $t-1$ indices have been selected, we set $B := A - \pi_{i_1, \dots, i_{t-1}} A$. Then

$$B_i = B - \pi_i B = \left(I - \frac{b_i b_i^T}{\|b_i\|_2^2} \right) B$$

is a rank-1 modification of B . Deshpande and Rademacher [14] propose two methods to compute (7) for $i = 1, \dots, n$. In the following, we summarize them briefly.

1. Algorithm 2 in [14] computes BB^T explicitly and then computes $B_i B_i^T$ as a rank-2 update of BB^T for every $i = 1, \dots, n$. The characteristic polynomial of $B_i B_i^T$ is computed by establishing a similarity transformation to a matrix in Frobenius normal form [7, section 16.6]. Fast matrix-matrix multiplication and inversion can be exploited so that the cost of this approach is $O(nm^\omega \log m)$, where $\omega \leq 2.373$ is the best exponent of matrix-matrix multiplication complexity.
2. Algorithm 3 in [14] first computes the thin SVD

$$(8) \quad B = U \Sigma V^T, \quad U \in \mathbb{R}^{m \times m}, \quad \Sigma \in \mathbb{R}^{m \times m}, \quad V \in \mathbb{R}^{n \times m},$$

where U and V have orthonormal columns and Σ is a diagonal matrix. Then it computes the characteristic polynomial of BB^T from the squared singular values of B , and the auxiliary polynomials $g_j(x) = \prod_{\ell \neq j} (x - \sigma_\ell^2(B))$ for $j = 1, \dots, m$. For $h = m - k + t$ and $h = m - k + t - 1$, the coefficient $c_h(B_i B_i^T)$ can then be computed as the coefficient of x^h in

$$(9) \quad \det(xI - BB^T) + \frac{1}{\|b_i\|_2^2} \sum_{j=1}^n \sigma_j^2(B) v_{ij}^2 g_j(x).$$

The cost of this second approach is $O(m^2 n)$.

The problem of computing the Frobenius normal form of a matrix is “numerically not viable” [33]. Also, updating directly the characteristic polynomial as in (9) is prone to numerical cancellation, leading to inaccurate results. For instance, consider the 2×2 matrix

$$A = \begin{bmatrix} 6.583644 \cdot 10^{-7} & 8.113362 \cdot 10^{-3} \\ 8.113362 \cdot 10^{-3} & 100 \end{bmatrix}$$

and the column selection problem for $k = 1$. Algorithm 4 in [14] using (9) selects the first column, giving an error $\|A - A(:, 1)A(:, 1)^\dagger A\|_F \approx 1.2 \cdot 10^{-6} \gg \sqrt{2}\sigma_2(A) = 1.4 \cdot 10^{-10}$.

Therefore, from now on we will avoid updating coefficients of characteristic polynomials and work with singular values instead. More specifically, we will compute the singular values of B_i by updating the SVD of B and then apply the summation algorithm [34, Algorithm 1] to compute the coefficients of the characteristic polynomial of $B_i B_i^T$ from its eigenvalues (that is, the squared singular values of B_i) with $O(m^2)$ operations in a numerically forward stable manner. To describe the updating procedure, consider the (thin) SVD $B = U \Sigma V^T$ as in (8). The (nonzero) singular values of B_i and

$$U^T B_i V = (I - U^T \pi_i U) U^T B V = \left(I - \frac{U^T b_i b_i^T U}{\|b_i\|_2^2} \right) \Sigma = (I - q q^T) \Sigma,$$

with $q = U^T b_i / \|b_i\|_2$, are identical. Using standard bulge chasing algorithms (see, e.g., [41, Algorithm 3.4] and [2]) it is possible to find orthogonal matrices $Q, W \in \mathbb{R}^{m \times m}$ such that $Q^T q = e_1$, where e_1 denotes the first unit vector, and $Q^T \Sigma W$ is upper bidiagonal. In turn, the singular values can be computed from the bidiagonal

matrix

$$Q^T(I - qq^T)\Sigma W = (I - e_1 e_1^T)(Q^T \Sigma W).$$

The matrices Q and W are composed of $O(m^2)$ Givens rotations [21, section 5.1], and the computation of $Q^T \Sigma W$ requires one to apply each of these rotations to at most 3 vectors. In turn, the cost of computing this bidiagonal matrix is $O(m^2)$, which is identical to the cost of computing its singular values [21, section 8.6].

2.3. Overall algorithm. The described variation of the column subset selection algorithm by Deshpande and Rademacher is summarized in Algorithm 1. One execution of line 3 is $O(nm^2)$, lines 6–9 are $O(m^2)$, and lines 14–15 are $O(knm)$. In summary, the overall complexity of Algorithm 1 is $O(knm^2)$. This is identical to the complexity of [14, Algorithm 4] combined with [14, Algorithm 3], and it is better than [14, Algorithm 4] combined with [14, Algorithm 2].

Note that instead of lines 14–15 we could have updated $B \leftarrow B - \pi_{i_t} B$. However, we noticed that recomputing B in lines 14–15 tends to improve accuracy and does not change the overall asymptotic complexity.

Algorithm 1 Column subset selection.

Input: $A \in \mathbb{R}^{m \times n}$, rank $1 \leq k < m$

Output: Column indices $S \in \{1, \dots, n\}^k$

```

1: Initialize  $S = \emptyset$  and  $B = A$ 
2: for  $t = 1, \dots, k$  do
3:   Compute  $U$  and  $\Sigma$  from the thin SVD of  $B = U\Sigma V^T$ 
4:   minRatio =  $+\infty$ 
5:   for  $i = 1, \dots, n$  do
6:      $q = U^T b_i / \|b_i\|_2$ 
7:      $D = Q^T \Sigma W$  bidiagonal matrix obtained by bulge chasing [41, Algorithm 3.4]
8:     Compute singular values  $\sigma_1, \dots, \sigma_m$  of  $(I - e_1 e_1^T)D$ 
9:     Apply the summation algorithm [34, Algorithm 1] to compute  $c_{m-k+t-1}(B_i B_i^T)$ 
       and  $c_{m-k+t}(B_i B_i^T)$  from eigenvalues  $\sigma_1^2, \dots, \sigma_m^2$ 
10:    Set ratio =  $c_{m-k+t-1}(B_i B_i^T) / c_{m-k+t}(B_i B_i^T)$ 
11:    if ratio < minRatio then Set minRatio = ratio and  $i_t = i$  end if
12:  end for
13:  Append index  $S \leftarrow (S, i_t)$ 
14:  Compute orthonormal basis  $Q$  of  $A(:, S)$ 
15:   $B = A - QQ^T A$ 
16: end for
```

2.4. Early stopping of column search. For each column index, Algorithm 1 needs to traverse $O(n)$ columns in order to find the one that minimizes the coefficient ratio or, equivalently, the conditional expectation. This column search can be shortened. To describe the idea, we revisit the argument from section 2.1 that has led to Algorithm 1. Recall that (4) states $\mathbb{E}[\|A - \pi_{X_1, \dots, X_k} A\|_F^2] \leq (k+1)(\sigma_{k+1}^2 + \dots + \sigma_m^2)$. This implies that there exists i_1 such that

$$\mathbb{E}[\|A - \pi_{X_1, \dots, X_k} A\|_F^2 \mid X_1 = i_1] \leq (k+1)(\sigma_{k+1}^2 + \dots + \sigma_m^2).$$

In particular, an index i_1 that minimizes the left-hand side will satisfy the bound, which is the choice made in Algorithm 1. However, there may be other choices of i_1

that satisfy the bound. *Any* such i_1 is a suitable choice. More generally, suppose that i_1, \dots, i_{t-1} have already been selected such that

$$\mathbb{E} [\|A - \pi_{X_1, \dots, X_k} A\|_F^2 \mid X_1 = i_1, \dots, X_{t-1} = i_{t-1}] \leq (k+1)(\sigma_{k+1}^2 + \dots + \sigma_m^2)$$

holds. This implies the existence of i_t such that

$$(10) \quad \mathbb{E} [\|A - \pi_{X_1, \dots, X_k} A\|_F^2 \mid X_1 = i_1, \dots, X_t = i_t] \leq (k+1)(\sigma_{k+1}^2 + \dots + \sigma_m^2).$$

Again, there is no need to choose an index i_t that minimizes the left-hand side; *any* i_t such that (10) holds is a suitable choice. By induction, choosing in every step an index such that (10) is verified implies that the error bound (1) holds.

The discussion above suggests modifying Algorithm 1 such that it computes

$$\text{bound} = (k+1) \cdot (\sigma_{k+1}^2 + \dots + \sigma_m^2)$$

in the beginning and substituting line 11 with

11: **if** $(k-t+1) \cdot \text{ratio} \leq \text{bound}$ **then** Set $i_t = i$ and break **end if**

To be able to stop the search early, it is important to test the columns in a suitable order. We found it beneficial to test the columns of B in descending Euclidean norm.² For each step t , computing the norms of all columns of B and sorting them has complexity $O(mn + n \log n)$.

We remark that selecting one of the columns with largest Euclidean norm may not always be the best choice. Consider, for example, the matrix

$$A = \begin{bmatrix} 1 & 0 & 10^{-b} \\ 0 & 1 & 10^{-b} \\ 0 & 0 & 10^{-2b} \end{bmatrix}$$

for some $b > 1$, say $b = 16$. For $k = 1$, the optimal choice is the third column, which is the one of smallest norm. This matrix also nicely illustrates that the obvious greedy approach (in order to get k columns of A , one first chooses the column i_1 that minimizes $\|A - \pi_{i_1} A\|_F$, then the column i_2 that minimizes $\|A - \pi_{i_1, i_2} A\|_F$, and so on) comes with no guarantees and may, in fact, utterly fail. For $k = 2$ the optimal choice consists of the first two columns. On the other hand, the greedy approach for $k = 2$ first selects the third column and then the first column, resulting in the arbitrarily bad error ratio $\frac{\text{error greedy}}{\text{error best}} \approx 10^b$.

This example also shows that, given a column subset of cardinality $k-1$ selected by Algorithm 1, one cannot obtain a suitable selection of k columns by simply performing another step of Algorithm 1. In order to ensure that (1) holds, Algorithm 1 needs to be re-run from scratch with k instead of $k-1$.

2.5. Numerical experiments. Both variants of Algorithm 1, with and without early stopping, have been implemented in MATLAB version R2019a. As the bulge chasing algorithm in line 7 would perform poorly in MATLAB, this part has been implemented in C++ and is called via a MEX interface.³ All numerical experiments in this work have been run on an eight-core Intel Core i7-8650U 1.90 GHz CPU, with

²For some discussion on the justification of this idea, see Lemma 1 in the arXiv version of this paper at <https://arxiv.org/abs/1908.06059>.

³Our implementation is available at <https://github.com/Alice94/ApproxFrobNorm> together with the codes to reproduce the figures in this paper.

256 KB of level 2 Cache and 16 GB of RAM. Multithreading has been turned off in order to not distort the findings.

We have applied the algorithm to the following three matrices:

1. the Hilbert matrix $A_{\text{hilb}} \in \mathbb{R}^{200 \times 200}$ given by $A_{\text{hilb}}(i, j) = \frac{1}{i+j-1}$;
2. $A_{\text{exp}} \in \mathbb{R}^{100 \times 200}$ given by $A_{\text{exp}}(i, j) = \exp(-0.3 \cdot |i - j|/200)$;
3. $A_{\text{poly}} \in \mathbb{R}^{100 \times 200}$ given by $A_{\text{poly}}(i, j) = \left(\left(\frac{i}{200} \right)^{20} + \left(\frac{j}{200} \right)^{20} \right)^{1/20}$.

The obtained results are shown in Figures 1, 2, and 3 respectively. Each left plot contains, for different values of k , the approximation error $\|A - A(:, S)A(:, S)^\dagger A\|_F$ returned by Algorithm 1, with and without early stopping. We make comparisons with the best rank- k approximation error $\sqrt{\sigma_{k+1}^2 + \dots + \sigma_m^2}$ and the upper bound (1), that is, $\sqrt{(k+1)(\sigma_{k+1}^2 + \dots + \sigma_m^2)}$. It can be seen that both variants of our algorithm stay below the upper bound until it reaches the level of roundoff error. Interestingly, for the matrix A_{exp} , which features the slowest singular value decay, the observed approximation error is much closer to the best approximation error than to the upper bound. The right plots of the figures show, for different values of k , the ratio between the total execution times of Algorithm 1 without early stopping and with early stopping. For example, for the matrix A_{hilb} , using early stopping in Algorithm 1 reduces the time for constructing an approximation of rank $k = 15$ by a factor 22. For the variant with early stopping, we also plot the number of columns that were examined. In the most optimistic scenario, only k columns need to be examined, which means that in every step of the algorithm already the first verifies the desired criterion. The plots reveal that our algorithm actually stays pretty close to this ideal situation, at least for the matrices considered. Note that for values of k larger than the numerical rank of the matrix, Algorithm 1 starts computing ratios (6) from singular values of the order of machine precision. In turn, the computations are severely affected by roundoff error, and it may, in fact, happen that the early stopping criterion is never satisfied. This leads to meaningless results, and we therefore truncate the plots before this happens. A proper implementation of Algorithm 1 needs to detect such a situation and reduce k accordingly.

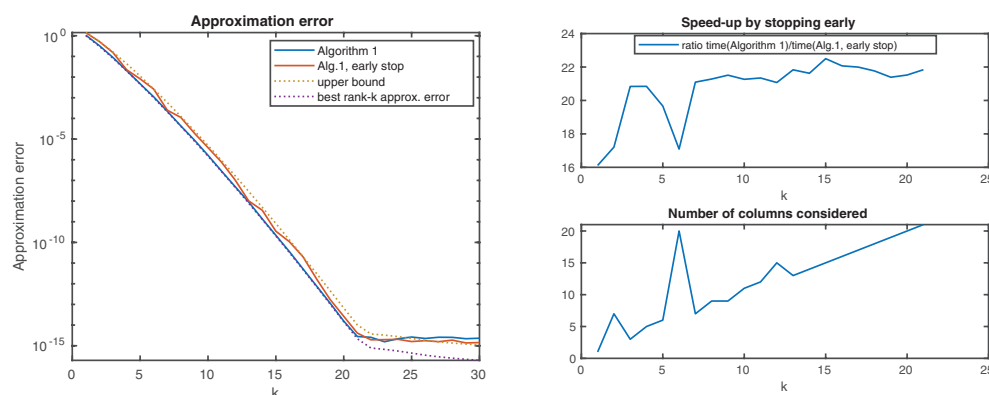
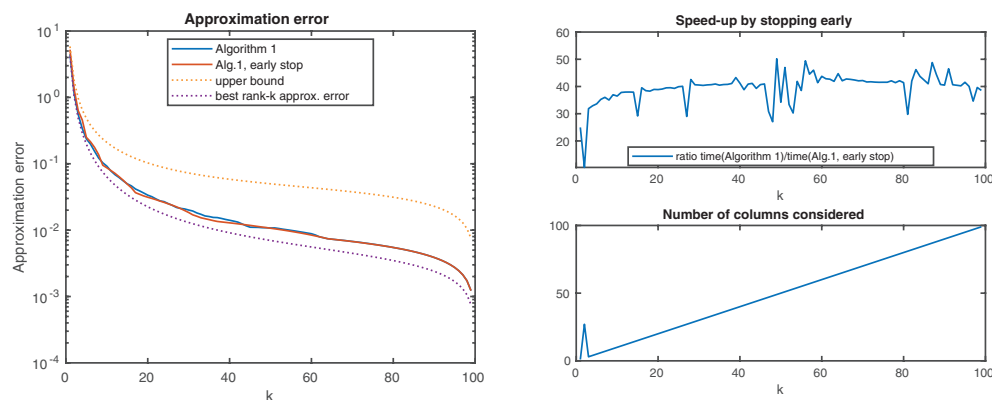
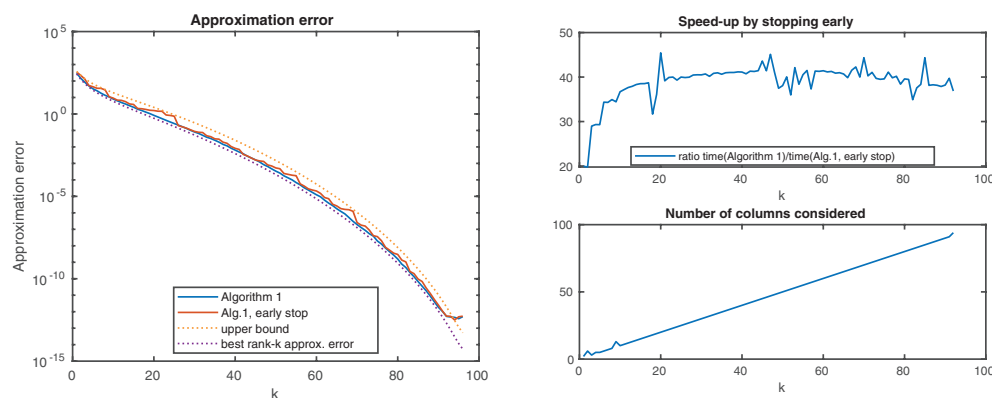


FIG. 1. Results for matrix A_{hilb} .

3. Matrix approximation. In this section, we extend the developments from section 2 on column subset selection to compute certain low-rank matrix approximations of a matrix $A \in \mathbb{R}^{m \times n}$ with $m \leq n$. As already discussed in the introduction, we will pursue two ways. First, in section 3.1, we discuss a general CUR approxima-

FIG. 2. Results for matrix A_{exp} .FIG. 3. Results for matrix A_{poly} .

tion obtained from applying column subset selection to the columns and rows of the matrix. Second, in section 3.2, we present a novel approach to cross approximation, a specific type of CUR approximation, with guaranteed error bounds.

3.1. CUR approximation induced by column subset selection. Suppose that $C \in \mathbb{R}^{m \times k}$ and $R \in \mathbb{R}^{k \times n}$ have been chosen. Then the matrix $U \in \mathbb{R}^{k \times k}$ that minimizes $\|A - CUR\|_F$ is given by the projection $U = C^\dagger A R^\dagger$; see [37, p. 320]. The following corollary provides an error bound for the case when C and R are determined by the techniques from section 2, leading to Algorithm 2. Closely related results can be found in the literature; see, for example, [17, Theorem 4], [35, Corollary 3.5], and [36, Theorem 4.1].

Algorithm 2 Matrix approximation by column subset selection.

Input: $A \in \mathbb{R}^{m \times n}$, rank k

Output: Rank- k CUR approximation, with C, R containing columns and rows of A

- 1: Compute C by applying Algorithm 1 to select k columns of A
 - 2: Compute R by applying Algorithm 1 to select k columns of A^T
 - 3: Compute $U = C^\dagger A R^\dagger$
-

COROLLARY 1. Let $A \in \mathbb{R}^{m \times n}$, with $1 \leq k \leq m \leq n$. Then the CUR approximation returned by Algorithm 2 satisfies

$$\|A - CUR\|_F \leq \sqrt{2k+2} \sqrt{\sigma_{k+1}^2(A) + \cdots + \sigma_m^2(A)}.$$

Proof. Using inequality (1) twice and the fact that CC^\dagger is an orthogonal projection, we obtain

$$\begin{aligned} \|A - CUR\|_F^2 &= \|A - CC^\dagger AR^\dagger R\|_F^2 \\ &= \|A - CC^\dagger A\|_F^2 + \|CC^\dagger(A - AR^\dagger R)\|_F^2 \\ &\leq \|(I - CC^\dagger)A\|_F^2 + \|A(I - R^\dagger R)\|_F^2 \\ &\leq 2(k+1) (\sigma_{k+1}^2(A) + \cdots + \sigma_m^2(A)). \quad \square \end{aligned}$$

3.1.1. Numerical experiments. We have tested a MATLAB implementation of Algorithm 2 in the setting and for the matrices A_{hilb} , A_{exp} , A_{poly} described in section 2.5. Figure 4 displays the obtained approximation errors $\|A - CUR\|_F$ for different values of k . Again, we have tested both variants of Algorithm 1, with and without early stopping, within Algorithm 2. The speedups obtained from early stopping are very similar to the ones reported in section 2.5, and therefore we refrain from providing details.

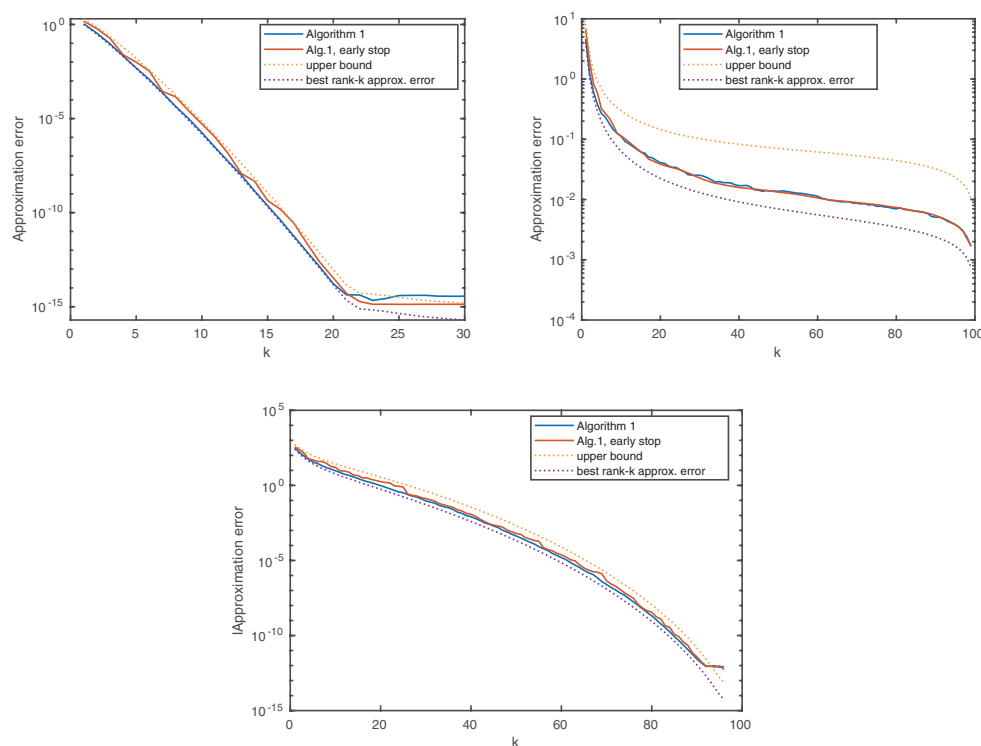


FIG. 4. Approximation errors for matrices A_{hilb} (top left), A_{exp} (top right), and A_{poly} (bottom).

We also consider, for $0 < \alpha < 1$, the $n \times n$ matrix

$$A = Q \cdot \text{diag}(1, \alpha, \alpha^2, \dots, \alpha^{n-1}) \cdot Q^T,$$

where $Q \in \mathbb{R}^{n \times n}$ is determined as the orthogonal factor from the QR decomposition of

$$\begin{bmatrix} 1 & & & & \\ -1 & 1 & & & \\ -1 & -1 & 1 & & \\ \vdots & \vdots & & \ddots & \\ -1 & -1 & -1 & \cdots & 1 \end{bmatrix}.$$

This is known to be a challenging example for the CUR approximation induced by DEIM (discrete empirical interpolation method); see [36, section 4.2], which determines the row and column indices by greedily choosing a maximum volume submatrix of U_k and V_k containing the first k left and right singular vectors of A , respectively. For the example above, the DEIM-induced CUR approximation always chooses $1, \dots, k$ for the column and row indices. For $\alpha = 0.1$, $n = 6$, $k = 5$, the error resulting from this choice is given by

$$\|A - A(:, 1:5)A(:, 1:5)^\dagger AA(1:5, :)^\dagger A(1:5, :)\|_F \approx 2.6 \cdot 10^{-9},$$

which is a magnitude larger than the upper bound $\sqrt{2(k+1)}\sigma_6 \approx 3.5 \cdot 10^{-10}$ guaranteed by Algorithm 2. Note that the latter algorithm selects the last 5 rows and columns for this example, leading to an error of $\approx 1.3 \cdot 10^{-10}$.

3.2. Cross approximation. We now consider *cross approximations*, which take the form

$$(11) \quad A \approx A(:, J)A(I, J)^{-1}A(I, :)$$

for row/column index tuples

$$(I, J) \in \Omega := \{1, \dots, m\}^k \times \{1, \dots, n\}^k$$

such that $A(I, J)$ is invertible. The different choice of the middle matrix makes a fundamental difference. In particular, as the following example shows, choosing the indices I, J as in Algorithm 2 may lead to poor approximation error.

Example 2. Consider $A = \begin{bmatrix} 2\varepsilon & 1 \\ 1 & \varepsilon \end{bmatrix}$ for $\varepsilon > 0$ and $k = 1$. Clearly, the first column and row satisfy the bound (1) for $k = 1$ with respect to A and A^T , respectively. However, the error of the corresponding cross approximation, $\|A - A(:, 1)A(1, 1)^{-1}A(1, :)\|_F = \frac{1}{2\varepsilon} - \varepsilon$, becomes arbitrarily large as $\varepsilon \rightarrow 0$.

Zamarashkin and Osinsky [42] have shown the existence of a cross approximation that satisfies a polynomial error bound in the Frobenius norm. To summarize their result, let

$$(X, Y) = (X_1, \dots, X_k, Y_1, \dots, Y_k)$$

be a $(2k)$ -tuple of random variables with values in Ω such that

$$(12) \quad \mathbb{P}(X = I, Y = J) := \frac{\text{Vol}^2(A(I, J))}{\sum_{(I', J') \in \Omega} \text{Vol}^2(A(I', J'))}.$$

Note that $\text{Vol}(A(I, J)) = 0$ whenever i_1, \dots, i_k or j_1, \dots, j_k contain repeated indices. Then [42, Theorem 1] shows that

$$(13) \quad \mathbb{E}[\|A - A(:, Y)A(X, Y)^{-1}A(X, :)\|_F^2] \leq (k+1)^2 (\sigma_{k+1}^2 + \dots + \sigma_m^2).$$

In particular, this implies that there exists $(I, J) \in \Omega$ such that

$$(14) \quad \|A - A(:, J)A(I, J)^{-1}A(I, :)\|_F^2 \leq (k+1)^2 (\sigma_{k+1}^2 + \cdots + \sigma_m^2).$$

In analogy to section 2.1 and [14], we will now derandomize this result, producing a polynomial-time deterministic algorithm that returns a cross approximation satisfying (14). The key for doing so is to find an expression for the conditional expectations that is easy to work with.

3.2.1. Conditional expectations.

LEMMA 3. Let $1 \leq t \leq k$ and $(i_1, \dots, i_t, j_1, \dots, j_t)$ be such that

$$\mathbb{P} \left(\begin{smallmatrix} X_1=i_1, \dots, X_t=i_t \\ Y_1=j_1, \dots, Y_t=j_t \end{smallmatrix} \right) > 0$$

for a random $(2k)$ -tuple (X, Y) with the probability distribution defined by (12). Consider

$$B = A - A(:, [j_1 \ \cdots \ j_t])A([i_1 \ \cdots \ i_t], [j_1 \ \cdots \ j_t])^{-1}A([i_1 \ \cdots \ i_t], :),$$

the remainder of cross approximation after choosing row indices i_1, \dots, i_t and column indices j_1, \dots, j_t . Then

$$\begin{aligned} & \mathbb{E} [\|A - A(:, Y)A(X, Y)^{-1}A(X, :)\|_F^2 \mid \begin{smallmatrix} X_1=i_1, \dots, X_t=i_t \\ Y_1=j_1, \dots, Y_t=j_t \end{smallmatrix}] \\ &= (k-t+1)^2 \cdot \frac{c_{m-k+t-1}(BB^T)}{c_{m-k+t}(BB^T)}, \end{aligned}$$

where the coefficients $c_{m-k+t}, c_{m-k+t-1}$ are defined as in (6) and the expectation is taken with respect to the distribution (12) defined on the $(2k)$ -tuples in Ω .

Proof. To simplify notation, we let $I_1 = (i_1, \dots, i_t)$, $I_2 = (i_{t+1}, \dots, i_k)$, and $I = (I_1, I_2) = (i_1, \dots, i_k)$ and define J_1, J_2, J analogously. In the following, we always use the convention that row and column summation indices range from 1 to m and from 1 to n , respectively. We have that

$$\begin{aligned} & \mathbb{E} [\|A - A(:, Y)A(X, Y)^{-1}A(X, :)\|_F^2 \mid \begin{smallmatrix} X_1=i_1, \dots, X_t=i_t \\ Y_1=j_1, \dots, Y_t=j_t \end{smallmatrix}] \\ &= \sum_{\substack{i_{t+1}, \dots, i_k \\ j_{t+1}, \dots, j_k}} \|A - A(:, J)A(I, J)^{-1}A(I, :)\|_F^2 \cdot \mathbb{P}(X = I, Y = J \mid \begin{smallmatrix} X_1=i_1, \dots, X_t=i_t \\ Y_1=j_1, \dots, Y_t=j_t \end{smallmatrix}) \\ (15) \quad &= \frac{1}{\gamma} \cdot \sum_{\substack{i_{t+1}, \dots, i_k, i_{k+1} \\ j_{t+1}, \dots, j_k, j_{k+1}}} \text{Vol}^2(A((I, i_{k+1}), (J, j_{k+1}))), \end{aligned}$$

with

$$\gamma = \sum_{\substack{i_{t+1}, \dots, i_k \\ j_{t+1}, \dots, j_k}} \text{Vol}^2(A(I, J)).$$

For establishing the equality in (15) we used from [42, Lemma 1] that

$$\|A - A(:, J)A(I, J)^{-1}A(I, :)\|_F^2 = \frac{\sum_{i_{k+1}, j_{k+1}} \text{Vol}^2(A((I, i_{k+1}), (J, j_{k+1})))}{\text{Vol}^2(A(I, J))}$$

and, from (12), that

$$\mathbb{P}(X = I, Y = J \mid \substack{X_1=i_1, \dots, X_t=i_t \\ Y_1=j_1, \dots, Y_t=j_t}) = \frac{\mathbb{P}(X = I, Y = J)}{\mathbb{P}(\substack{X_1=i_1, \dots, X_t=i_t \\ Y_1=j_1, \dots, Y_t=j_t})} = \frac{1}{\gamma} \cdot \text{Vol}^2(A(I, J)).$$

We now aim at simplifying the expression (15). For this purpose, we assume without loss of generality that $i_1 = 1, \dots, i_t = t$ and $j_1 = 1, \dots, j_t = t$. This allows us to partition

$$A(I, J) = \begin{bmatrix} A(I_1, J_1) & A(I_1, J_2) \\ A(I_2, J_1) & A(I_2, J_2) \end{bmatrix}, \quad B(I, J) = \begin{bmatrix} 0 & 0 \\ 0 & B(I_2, J_2) \end{bmatrix},$$

where $B(I_2, J_2) = A(I_2, J_2) - A(I_2, J_1)A(I_1, J_1)^{-1}A(I_1, J_2)$ by the definition of B . By the relation between determinants and Schur complements [26, equation (0.8.5.1)], $\text{Vol}(A(I, J)) = \text{Vol}(A(I_1, J_1)) \cdot \text{Vol}(B(I_2, J_2))$. Therefore,

$$\gamma = \sum_{\substack{i_{t+1}, \dots, i_k \\ j_{t+1}, \dots, j_k}} \text{Vol}^2(A(I, J)) = \sum_{\substack{i_{t+1}, \dots, i_k \\ j_{t+1}, \dots, j_k}} \text{Vol}^2(B(I_2, J_2)) \cdot \text{Vol}^2(A(I_1, J_1)).$$

Analogously, one shows

$$\begin{aligned} \sum_{\substack{i_{t+1}, \dots, i_{k+1} \\ j_{t+1}, \dots, j_{k+1}}} \text{Vol}^2(A((I, i_{k+1}), (J, j_{k+1}))) \\ = \sum_{\substack{i_{t+1}, \dots, i_{k+1} \\ j_{t+1}, \dots, j_{k+1}}} \text{Vol}^2(B((I_2, i_{k+1}), (J_2, j_{k+1}))) \text{Vol}^2(A(I_1, J_1)). \end{aligned}$$

Inserting these expressions into (15) yields

$$\frac{\sum_{\substack{i_{t+1}, \dots, i_{k+1} \\ j_{t+1}, \dots, j_{k+1}}} \text{Vol}^2(B((I_2, i_{k+1}), (J_2, j_{k+1})))}{\sum_{\substack{i_{t+1}, \dots, i_k \\ j_{t+1}, \dots, j_k}} \text{Vol}^2(B(I_2, J_2))}.$$

By [28, Theorem 7] this ratio is equal to

$$\frac{c_{m-k+t-1}(BB^T) \cdot ((k-t+1)!)^2}{c_{m-k+t}(BB^T) \cdot ((k-t)!)^2} = (k-t+1)^2 \cdot \frac{c_{m-k+t-1}(BB^T)}{c_{m-k+t}(BB^T)}. \quad \square$$

3.2.2. Derandomized cross-approximation algorithm. With Lemma 3 at hand, we can proceed analogously to section 2.1 and sequentially find k pairs of row/column indices such that (14) is satisfied. Suppose that $t-1$ index pairs $(i_1, j_1), \dots, (i_{t-1}, j_{t-1})$ have been determined. Then the t th step of the algorithm proceeds by choosing (i_t, j_t) such that

$$(16) \quad \mathbb{E} [\|A - A(:, Y)A(X, Y)^{-1}A(X, :)\|_F^2 \mid \substack{X_1=i_1, \dots, X_t=i_t \\ Y_1=j_1, \dots, Y_t=j_t}]$$

is minimized. We will show in Theorem 4 below that this choice of index pairs leads to a cross approximation satisfying the desired error bound (14). In view of Lemma 3, the minimization of (16) means that in each step of the algorithm we need to compute the ratios

$$(17) \quad \frac{c_{m-k+t-1}(C_{ij}C_{ij}^T)}{c_{m-k+t}(C_{ij}C_{ij}^T)}, \quad i = 1, \dots, m, \quad j = 1, \dots, n,$$

where

$$C_{ij} = A - A(:, [j_1, \dots, j_{t-1}, j])A([i_1, \dots, i_{t-1}, i], [j_1, \dots, j_{t-1}, j])^{-1}A([i_1, \dots, i_{t-1}, i], :).$$

Parallelizing the developments in section 2.2, we now show how the coefficients in (17) can be computed via updating the singular values of C_{ij} . Let us denote the remainder from the previous step by

$$B = A - A(:, [j_1, \dots, j_{t-1}])A([i_1, \dots, i_{t-1}], [j_1, \dots, j_{t-1}])^{-1}A([i_1, \dots, i_{t-1}], :).$$

Then it follows that

$$(18) \quad C_{ij} = B - \frac{1}{B(i, j)}B(:, j)B(i, :);$$

see, e.g., [4]. We compute a thin SVD $B = U\Sigma V^T$ such that $U \in \mathbb{R}^{m \times m}$, $V \in \mathbb{R}^{m \times n}$ have orthonormal columns and $\Sigma \in \mathbb{R}^{m \times m}$ is diagonal. Note that

$$B(:, j) = U\Sigma V(j, :)^T, \quad B(i, :) = U(i, :)\Sigma V^T.$$

Inserted into (18), this shows that the nonzero singular values of C_{ij} match the singular values of

$$U^T C_{ij} V = \Sigma - \Sigma V(j, :)^T \cdot \frac{U(i, :)\Sigma}{B(i, j)} = \Sigma - xy^T,$$

where $x = \Sigma V(j, :)^T$ and $y = \frac{1}{B(i, j)}\Sigma U(i, :)^T$ are vectors of length m and can be computed with $O(m^2)$ operations.

Similarly as in section 2.2, we transform $\Sigma - xy^T$ into bidiagonal form, after which its singular values can be computed with $O(m^2)$ operations. This transformation proceeds in three steps:

1. We compute orthogonal matrices Q and W such that $Q^T \Sigma W$ is upper bidiagonal and $Q^T x = \pm \|x\|_2 \cdot e_1$ using, for example, [41, Algorithm 3.4]. In turn, the matrix

$$(19) \quad D_1 := Q^T (\Sigma - xy^T) W$$

is bidiagonal with an additional nonzero first row; see the first plot in Figure 5 for an illustration.

2. By a bulge chasing algorithm, we transform D_1 into an upper banded matrix D_2 with two superdiagonals using $O(m^2)$ Givens rotations. We refrain from giving a detailed description of the algorithm and refer the reader to Figure 5 for an illustration.

3. The banded matrix D_2 is reduced to a bidiagonal matrix D_3 using the LAPACK [1] routine `dgbbrd`.

The overall procedure described above can be implemented by means of $O(m^2)$ Givens rotations, each of which is applied to a small matrix of size independent of m, n . Hence, it has complexity $O(m^2)$.

Algorithm 3 summarizes our newly proposed method for cross approximation. The SVD needed at the beginning of each outer loop is of complexity $O(m^2 n)$, and each of the mn inner loops costs $O(m^2)$ operations; the total complexity of Algorithm 3 is therefore $O(km^3 n)$.

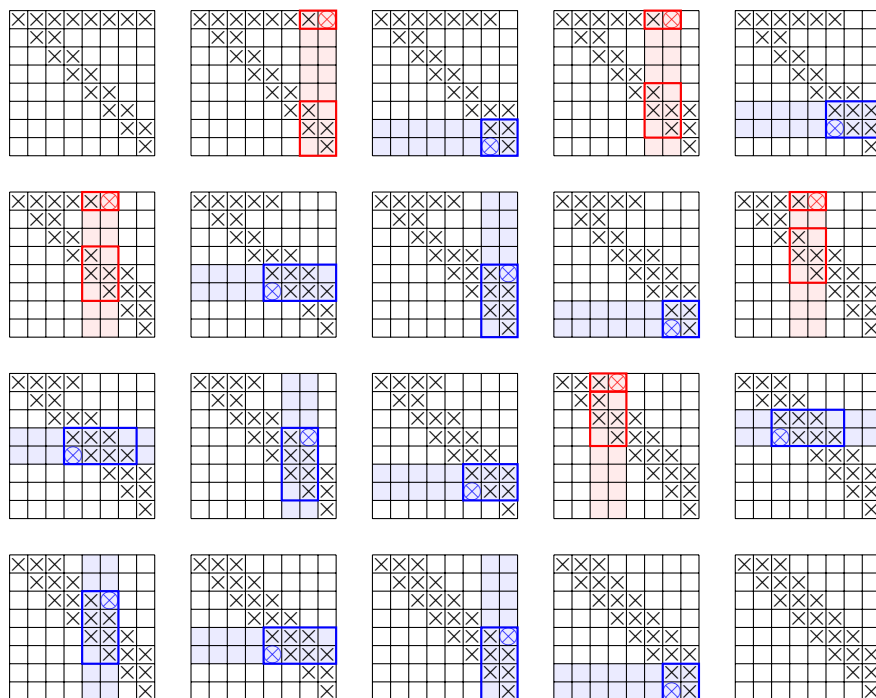


FIG. 5. Illustration of bulge chasing algorithm to transform a bidiagonal matrix with an additional nonzero first row to an upper banded matrix. In each plot, except for the first and last ones, a Givens rotation is applied to a pair of rows or columns to zero out the entry denoted by \otimes .

THEOREM 4. For a matrix A of rank at least k , Algorithm 3 returns index sets I and J such that (14) is satisfied.

Proof. Let $B_{\{X,Y\}} = A - A(:, Y)A(X, Y)^{-1}A(X, :)$. For $t = 1, \dots, k$ we have that

$$\begin{aligned} & \mathbb{E} \left[\|B_{\{X,Y\}}\|_F^2 \mid \begin{smallmatrix} X_1=i_1, \dots, X_{t-1}=i_{t-1} \\ Y_1=j_1, \dots, Y_{t-1}=j_{t-1} \end{smallmatrix} \right] \\ &= \sum_{i,j} \mathbb{E} \left[\|B_{\{X,Y\}}\|_F^2 \mid \begin{smallmatrix} X_1=i_1, \dots, X_{t-1}=i_{t-1}, X_t=i \\ Y_1=j_1, \dots, Y_{t-1}=j_{t-1}, Y_t=j \end{smallmatrix} \right] \mathbb{P}(X_t = i, Y_t = j \mid \begin{smallmatrix} X_1=i_1, \dots, X_{t-1}=i_{t-1} \\ Y_1=j_1, \dots, Y_{t-1}=j_{t-1} \end{smallmatrix}). \end{aligned}$$

Therefore, as (13) holds, the choice (16) inductively ensures that

$$\begin{aligned} \mathbb{E} \left[\|B_{\{X,Y\}}\|_F^2 \mid \begin{smallmatrix} X_1=i_1, \dots, X_t=i_t \\ Y_1=j_1, \dots, Y_t=j_t \end{smallmatrix} \right] &\leq \mathbb{E} \left[\|B_{\{X,Y\}}\|_F^2 \mid \begin{smallmatrix} X_1=i_1, \dots, X_{t-1}=i_{t-1} \\ Y_1=j_1, \dots, Y_{t-1}=j_{t-1} \end{smallmatrix} \right] \\ &\leq (k+1)^2 (\sigma_{k+1}^2 + \dots + \sigma_m^2). \end{aligned}$$

Therefore, the index sets I and J computed by Algorithm 3 satisfy the bound (14). \square

In analogy to the discussion in section 2.4, let us emphasize that it is not necessary to select the pair (i_t, j_t) that minimizes the ratio r . Any pair (i, j) for which the inequality

$$(20) \quad (k-t+1)^2 \frac{c_{m-k+t-1}(C_{ij}C_{ij}^T)}{c_{m-k+t}(C_{ij}C_{ij}^T)} \leq (k+1)^2 (\sigma_{k+1}^2(A) + \dots + \sigma_m^2(A))$$

holds will lead to index sets I and J such that (14) is satisfied. Inspired by adaptive cross approximation with full pivoting [4], we traverse the entries of B from the largest

Algorithm 3 Derandomized cross approximation.**Input:** $A \in \mathbb{R}^{m \times n}$ with $m \leq n$, integer $k \leq m$ **Output:** Index sets I, J of cardinality k defining the cross approximation (11)

```

1: Initialize  $I \leftarrow \emptyset, J \leftarrow \emptyset$ , and  $B \leftarrow A$ 
2: for  $t = 1, \dots, k$  do
3:    $[U, \Sigma, V] = \text{thin SVD of } B$ 
4:    $\text{minRatio} = +\infty$ 
5:   for  $i = 1, \dots, m$  do
6:     for  $j = 1, \dots, n$  do
7:        $x \leftarrow \Sigma V(j, :)^T, y \leftarrow \frac{1}{B(i, j)} \Sigma U(i, :)^T$ 
8:       Compute matrix  $D_1$  defined in (19) using [41, Algorithm 3.4]
9:       Transform  $D_1$  into upper banded form  $D_2$  using the bulge chasing algorithm
10:      Transform  $D_2$  into bidiagonal matrix  $D_3$  using LAPACK's dgbbrd
11:      Compute singular values  $\sigma_1, \dots, \sigma_m$  of  $D_3$ 
12:      Apply the summation algorithm [34, Algorithm 1] to obtain  $c_{m-k+t-1}(C_{ij}C_{ij}^T)$ 
        and  $c_{m-k+t}(C_{ij}C_{ij}^T)$  from eigenvalues  $\sigma_1^2, \dots, \sigma_m^2$ 
13:      Set  $r = \frac{c_{m-k+t-1}(C_{ij}C_{ij}^T)}{c_{m-k+t}(C_{ij}C_{ij}^T)}$ 
14:      if  $r < \text{minRatio}$  then  $i_t \leftarrow i, j_t \leftarrow j, \text{minRatio} = r$  end if
15:    end for
16:  end for
17:   $I \leftarrow I \cup \{i_t\}, J \leftarrow J \cup \{j_t\}$ 
18:   $B \leftarrow B - \frac{B(:, j_t) \cdot B(i_t, :)}{B(i_t, j_t)}$ 
19: end for

```

to the smallest (in magnitude) and stop the search once we have found an index pair (i_t, j_t) satisfying (20).

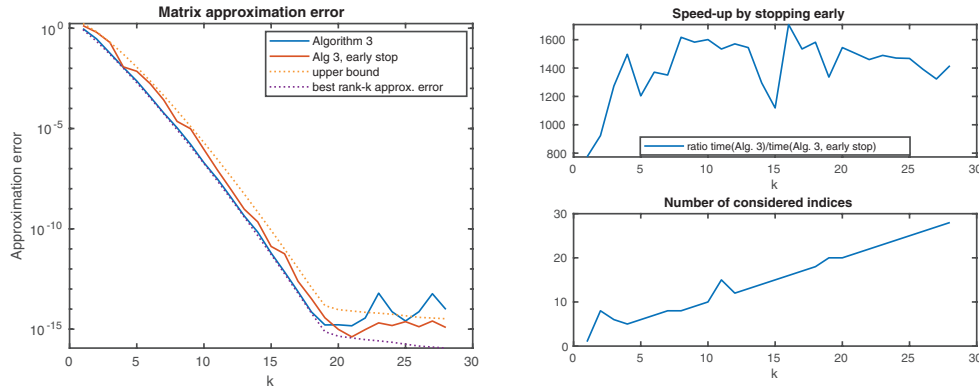
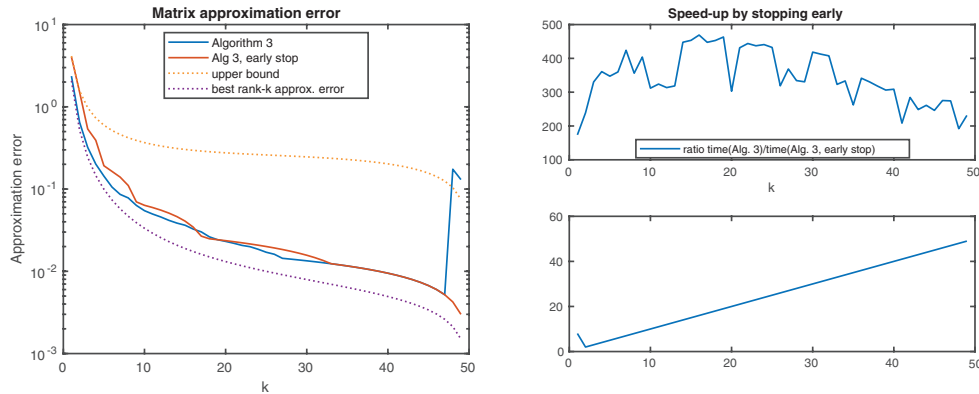
3.2.3. Numerical experiments. We have implemented both variants of Algorithm 3, with and without early stopping, in MATLAB. Again, the two inner loops have been implemented in a C++ function that is called via a MEX interface. The computational environment is the one described in section 2.5, but the test matrices are smaller because Algorithm 3 without early stopping is significantly slower. We choose A_{hilb} to be 100×100 , and choose A_{exp} to be 50×100 , and the matrix $A_{\text{poly}} \in \mathbb{R}^{50 \times 100}$ is given by

$$(21) \quad A_{\text{poly}}(i, j) = \left(\left(\frac{i}{100} \right)^{10} + \left(\frac{j}{100} \right)^{10} \right)^{1/10}.$$

The left plots of Figures 6, 7, and 8 display the approximation error $\|A - A(:, J)A(I, J)^{-1}A(I, :)\|_F$ for the index sets returned by both variants of Algorithm 3. The right plots display the ratios between the execution time of Algorithm 3 with and without early stopping, as well as the total number of index pairs that needed to be tested in Algorithm 3 with early stopping. It can be observed that early stopping dramatically accelerates the computation and is thus the preferred variant.

It can be seen that the approximation errors often stay close to the best rank- k approximation error $\sqrt{\sigma_{k+1}^2 + \dots + \sigma_m^2}$ and do not exceed the upper bound (14), modulo roundoff error. However, for larger values of k , Algorithm 3 without early stopping appears to encounter stability issues; the approximation error is distorted well above the level of roundoff error.

To better understand the numerical instability of our cross-approximation algo-

FIG. 6. Results for matrix A_{hilb} .FIG. 7. Results for matrix A_{exp} .

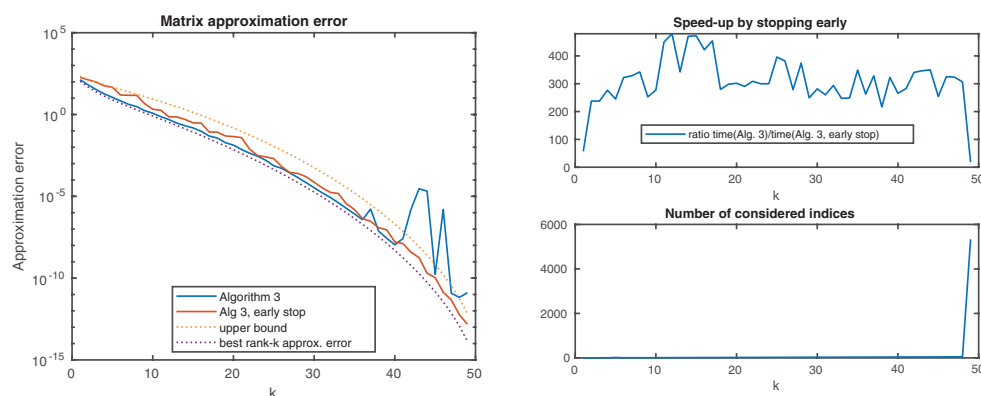
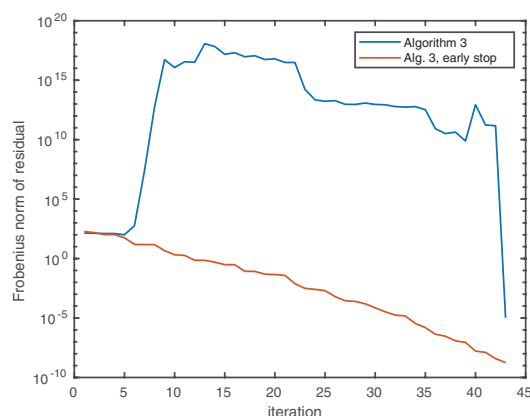
rithm, we analyzed what happens for the matrix A_{poly} for rank $k = 43$, for which Algorithm 3 without early stopping gives a cross-approximation error out of the bounds predicted by (14). We computed the Frobenius norm of the intermediate residuals

$$(22) \quad A - A(:, [j_1 \dots j_t]) \cdot A([i_1 \dots i_t], [j_1 \dots j_t])^{-1} \cdot A([i_1 \dots i_t], :)$$

for $t = 1, \dots, 43$ for the index sets given by Algorithm 3 with and without early stopping. We used MATLAB's `vpa` with 200 digits of accuracy to compute the Frobenius norm of the residual. The results are shown in Figure 9. The intermediate residuals grow significantly in Algorithm 3 without early stopping, which may completely spoil the accuracy of the singular values computed in lines 7–11. As these are needed to determine the indices to select at each step, this selection is not guaranteed to satisfy the bound (20). For A_{poly} and $k = 43$, this happens for the first time at iteration $t = 34$.

Such an intermediate growth of the residual can already happen for small matrices. For example, consider

$$A = \begin{bmatrix} -10^{-4} & 3 & -4 \\ 4 & 1 & 2 \\ 8 & -1 & 1 \end{bmatrix}.$$

FIG. 8. Results for matrix A_{poly} .FIG. 9. Frobenius norm of the residual (22) for $t = 1, \dots, k$ for the matrix A_{poly} with target rank $k = 43$ when using Algorithm 3 with and without early stopping.

When aiming at a rank-2 approximation, the choice that minimizes the expectation at the first step is the pivot in position $(1, 1)$, which results in a residual more than 10^4 larger than the norm of the original matrix.

The intermediate growth of the residuals may explain why Algorithm 3 with early stopping shows more stability in the examples we considered: If possible, it chooses one of the largest entries of the residual, which, in turn, should prevent the residuals from becoming too large. However, there are no results that ensure that we can take a “large entry” at each step of the algorithm; further investigation would be needed to understand the stability of Algorithm 3 with early stopping.

We also consider the $n \times n$ matrix $A = LDL^T$, where

$$L = \begin{bmatrix} 1 & & & & \\ -c & 1 & & & \\ -c & -c & 1 & & \\ \vdots & \vdots & & \ddots & \\ -c & -c & -c & \cdots & 1 \end{bmatrix}, \quad D = \begin{bmatrix} 1 & & & & \\ & s^2 & & & \\ & & s^4 & & \\ & & & \ddots & \\ & & & & s^{2(n-1)} \end{bmatrix},$$

with $s = \sin(\theta)$, $c = \cos(\theta)$ for some $0 < \theta < \pi$. This is known to be a challenging example for greedy cross approximation [25]: When $k = n - 1$ the greedy algorithm selects the leading $k \times k$ submatrix and returns an approximation error that is exponentially larger than the best approximation error. In contrast, Algorithm 3, with and without early stopping, makes the correct choice by selecting the last $n - 1$ rows and columns. For instance, for $n = 6$ and $\theta = 0.1$, we obtain the error

$$\|A - A(:, 2 : 6)A(2 : 6, 2 : 6)^{-1}A(2 : 6, :)\|_F \approx 3.9 \cdot 10^{-13} < 1.8 \cdot 10^{-12} \approx \sqrt{6}\sigma_n.$$

Selecting the first 5 rows and columns results in an error of $9.8 \cdot 10^{-11}$.

Finally, we would like to point out an interesting observation concerning the preservation of structure. In joint work with Massei [12], we have shown that for a symmetric positive definite matrix A there is always a symmetric choice of indices, $J = I$, leading to a symmetric cross approximation such that the favorable error bound of Goreinov and Tyrtyshnikov [22] is attained. For cross approximation in the Frobenius norm, the situation appears to be more complicated; it is generally not true that a symmetric choice of indices achieves the error bound (14) even when A is symmetric positive definite. For instance, for $n = 3$ and $k = 1$ consider

$$A = \begin{bmatrix} 1.87 & -1.82 & -2.11 \\ -1.82 & 1.87 & 2.11 \\ -2.11 & 2.11 & 2.54 \end{bmatrix}.$$

The best symmetric choice is $I = J = \{3\}$, but this leads to an error $\approx 0.1911 > 2\sqrt{\sigma_2^2 + \sigma_3^2} \approx 0.1821$.

4. Tensor approximation. As shown, for example, in [16, 23, 35], column subset selection can be used to approximate tensors as well. In the following, we demonstrate the use of the algorithm from section 2 to obtain approximations of low multilinear rank constructed from the fibers of a third-order tensor $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$.

First, we briefly recall some basic definitions for tensors and refer the reader to [29] for more details. Generalizing the notion of rows and columns of a matrix, the vectors obtained from \mathcal{A} by fixing all indices but the μ th one are called μ -mode fibers. The matrix $A^{(\mu)} \in \mathbb{R}^{n_\mu \times (n_1 n_2 n_3)/n_\mu}$ containing all μ -mode fibers as columns is called the μ -mode matricization of \mathcal{A} . The μ -mode product of a matrix $B \in \mathbb{R}^{m \times n_\mu}$ with \mathcal{A} is denoted by $\mathcal{A} \times_\mu B$, and it is the tensor such that its μ -mode matricization is given by $B \cdot A^{(\mu)}$. We use the Frobenius norm of a tensor defined by

$$\|\mathcal{A}\|_F^2 := \sum_{i_1=1}^{n_1} \sum_{i_2=1}^{n_2} \sum_{i_3=1}^{n_3} \mathcal{A}(i_1, i_2, i_3)^2$$

and recall that $\|\mathcal{A}\|_F = \|A^{(\mu)}\|_F$ for $\mu = 1, 2, 3$. The tuple (k_1, k_2, k_3) defined by $k_\mu = \text{rank}(A^{(\mu)})$ is called the multilinear rank of \mathcal{A} , and we can decompose \mathcal{A} as

$$\mathcal{A} = \mathcal{C} \times_1 B_1 \times_2 B_2 \times_3 B_3$$

for coefficient matrices $B_\mu \in \mathbb{R}^{n_\mu \times k_\mu}$ for $\mu = 1, 2, 3$ and a so-called *core tensor* $\mathcal{C} \in \mathbb{R}^{k_1 \times k_2 \times k_3}$. This so-called *Tucker decomposition* is particularly beneficial when the multilinear rank is much smaller than the size of a tensor.

Algorithm 4 produces an approximate Tucker decomposition for a given tensor such that each coefficient matrix B_μ is composed of μ -mode fibers. The following result

shows that the obtained approximation error remains close to the best approximation error.

Algorithm 4 Approximation of tensors by column selection.

Input: Tensor $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, integers k_1, k_2, k_3

Output: Approximate Tucker decomposition of multilinear rank (k_1, k_2, k_3) in terms of coefficient matrices B_1, B_2, B_3 and core tensor \mathcal{C}

- 1: **for** $\mu = 1, 2, 3$ **do**
 - 2: Compute $B_\mu = A^{(\mu)}(:, S_\mu)$ by applying Algorithm 1 to select k_μ columns from $A^{(\mu)}$
 - 3: **end for**
 - 4: Compute $\mathcal{C} = \mathcal{A} \times_1 B_1^\dagger \times_2 B_2^\dagger \times_3 B_3^\dagger$
-

COROLLARY 5. Consider $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ and integers k_1, k_2, k_3 such that $1 \leq k_\mu \leq n_\mu$ for $\mu = 1, 2, 3$. Then the output of Algorithm 4 satisfies

$$\|\mathcal{A} - \mathcal{C} \times_1 B_1 \times_2 B_2 \times_3 B_3\|_F \leq \sqrt{k_1 + k_2 + k_3 + 3} \cdot \|\mathcal{A} - \mathcal{A}_{\text{best}}\|_F,$$

where $\mathcal{A}_{\text{best}}$ is the best Tucker approximation of \mathcal{A} of multilinear rank at most (k_1, k_2, k_3) .

Proof. The proof is similar to existing proofs on the quasi-optimality of the higher-order SVD [13] and related results in [16, 23, 35].

Using (1) and setting $\pi_\mu = B_\mu B_\mu^+$, the result of Algorithm 1 applied to $A^{(\mu)}$ satisfies

$$\begin{aligned} \|A^{(\mu)} - \pi_\mu(A^{(\mu)})\|_F^2 &\leq (k_\mu + 1)(\sigma_{k_\mu+1}^2(A^{(\mu)}) + \cdots + \sigma_{n_\mu}^2(A^{(\mu)})) \\ &\leq (k_\mu + 1)\|A^{(\mu)} - A_{\text{best}}^{(\mu)}\|_F^2 = (k_\mu + 1)\|\mathcal{A} - \mathcal{A}_{\text{best}}\|_F^2, \end{aligned}$$

where the second inequality follows from the fact that $A_{\text{best}}^{(\mu)}$, the μ -mode matricization of $\mathcal{A}_{\text{best}}$, has rank at most k_μ . Using the orthogonality of the projections π_μ , we obtain

$$\begin{aligned} \|\mathcal{A} - \mathcal{C} \times_1 B_1 \times_2 B_2 \times_3 B_3\|_F^2 &= \|\mathcal{A} - \mathcal{A} \times_1 \pi_1 \times_2 \pi_2 \times_3 \pi_3\|_F^2 \\ &= \|\mathcal{A} - \mathcal{A} \times_1 \pi_1\|_F^2 + \|(\mathcal{A} - \mathcal{A} \times_1 \pi_1) \times_2 \pi_2\|_F^2 + \|(\mathcal{A} - \mathcal{A} \times_1 \pi_1) \times_2 \pi_2 \times_3 \pi_3\|_F^2 \\ &\leq \sum_{\mu=1}^3 \|\mathcal{A} - \mathcal{A} \times_\mu \pi_\mu\|_F^2 = \sum_{\mu=1}^3 \|A^{(\mu)} - \pi_\mu(A^{(\mu)})\|_F^2 \leq \sum_{\mu=1}^3 (k_\mu + 1) \|\mathcal{A} - \mathcal{A}_{\text{best}}\|_F^2 \\ &= (k_1 + k_2 + k_3 + 3) \|\mathcal{A} - \mathcal{A}_{\text{best}}\|_F^2, \end{aligned}$$

where the second equality follows from [39, Theorem 5.1]. \square

Remark 6. Algorithm 4 easily generalizes to tensors of arbitrary order. Given a tensor $\mathcal{A} \in \mathbb{R}^{n_1 \times \cdots \times n_d}$ and integers k_1, \dots, k_d , this generalization constructs subsets of fibers B_1, \dots, B_d and a core tensor \mathcal{C} such that

$$\|\mathcal{A} - \mathcal{C} \times_1 B_1 \times_2 \cdots \times_d B_d\|_F \leq \sqrt{k_1 + \cdots + k_d + d} \cdot \|\mathcal{A} - \mathcal{A}_{\text{best}}\|_F.$$

This compares favorably with other existence results in the literature, which feature much larger constants that grow exponentially with the order; see [23], [32, Theorem 3.1], and [35, Theorem 3.1].

4.1. Numerical experiments. We have implemented Algorithm 4 in MATLAB and tested it on two $50 \times 50 \times 50$ tensors, given by $\mathcal{A}_{\text{hilb}}(i, j, h) = \frac{1}{i+j+h-1}$ and

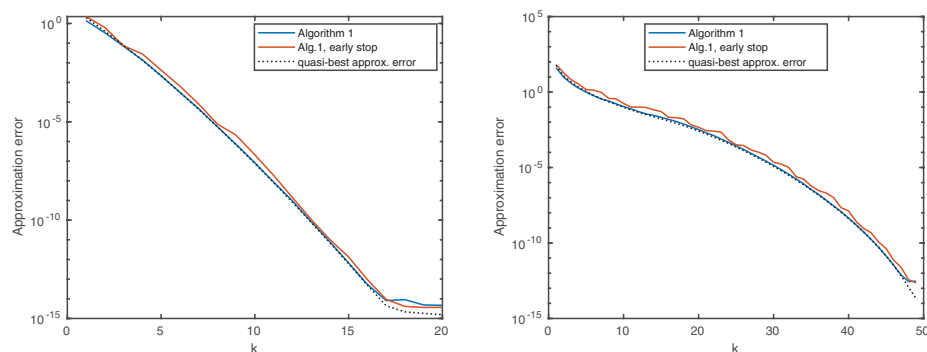


FIG. 10. Results for tensors $\mathcal{A}_{\text{hilb}}$ (left) and $\mathcal{A}_{\text{poly}}$ (right).

$\mathcal{A}_{\text{poly}}(i, j, h) = (i^{10} + j^{10} + h^{10})^{1/10} / 50$. We choose $k_1 = k_2 = k_3 = k$ and report in Figure 10 the obtained approximation errors $\|\mathcal{A} - \mathcal{C} \times_1 B_1 \times_2 B_2 \times_3 B_3\|_F$ for different values of k , where $B_1, B_2, B_3, \mathcal{C}$ are returned by Algorithm 4, with and without early stopping in the column selection part. We make comparisons with the quantity

$$\left(\sum_{\mu=1}^3 \sigma_{k_{\mu}+1}^2(A^{(\mu)}) + \cdots + \sigma_{n_{\mu}}^2(A^{(\mu)}) \right)^{1/2},$$

which provides a (tight) upper bound on the best approximation error. It can be seen that the errors obtained from Algorithm 4 remain close to this quasi-best approximation error.

5. Conclusions. In this work, we have proposed several improvements to the column selection algorithm by Deshpande and Rademacher [14]. The numerical experiments indicate that updating singular values (instead of characteristic polynomials) leads to numerical robustness in the sense that the approximation error obtained in finite precision arithmetic is not affected unduly by roundoff error. We have also developed an extension of [14] to produce cross approximations of matrices, and, to the best of our knowledge, this extension constitutes the first deterministic polynomial-time algorithm that yields a cross approximation with a guaranteed polynomial error bound. We have introduced a mechanism for stopping early the search for indices in column subset selection or cross approximation. Although relatively simple, this mechanism tremendously reduces the execution time for all examples tested.

A number of issues remain for future study, such as the numerical stability analysis of our algorithms. In particular, it would be desirable to study the numerical robustness of the cross approximation returned by Algorithm 3 with early stopping. Also, combining early stopping with a more aggressive reuse of the SVD might lead to further complexity reduction, but a rigorous complexity analysis would require deeper understanding of early stopping. Finally, we would like to stress that the algorithms presented in this work are intended for small- to medium-sized matrices and tensors. For large-scale data, the algorithms presented in this paper need to be combined with other, possibly heuristic dimensionality reduction techniques.

Acknowledgments. The authors thank Sergey Dolgov for helpful discussions on topics related to the work presented in this paper, and the referees for helpful remarks that improved the presentation of the paper.

REFERENCES

- [1] E. ANDERSON, Z. BAI, C. BISCHOF, L. S. BLACKFORD, J. DEMMEL, J. DONGARRA, J. DU CROZ, A. GREENBAUM, S. HAMMARLING, A. MCKENNEY, AND D. SORESENSEN, *LAPACK Users' Guide*, 3rd ed., SIAM, Philadelphia, 1999, <https://doi.org/10.1137/1.9780898719604>.
- [2] J. L. AURENTZ, T. MACH, L. ROBOL, R. VANDEBRIL, AND D. S. WATKINS, *Core-Chasing Algorithms for the Eigenvalue Problem*, Fundam. Algorithms 13, SIAM, Philadelphia, 2018, <https://doi.org/10.1137/1.9781611975345>.
- [3] M. BARRAULT, Y. MADAY, N. C. NGUYEN, AND A. T. PATERA, An “empirical interpolation” method: Application to efficient reduced-basis discretization of partial differential equations, *C. R. Math. Acad. Sci. Paris*, 339 (2004), pp. 667–672, <https://doi.org/10.1016/j.crma.2004.08.006>.
- [4] M. BEBENDORF, *Approximation of boundary element matrices*, *Numer. Math.*, 86 (2000), pp. 565–589, <https://doi.org/10.1007/PL00005410>.
- [5] C. BOUTSIDIS AND M. MAGDON-ISMAIL, *Deterministic feature selection for k-means clustering*, *IEEE Trans. Inform. Theory*, 59 (2013), pp. 6099–6110, <https://doi.org/10.1109/TIT.2013.2255021>.
- [6] C. BOUTSIDIS, M. W. MAHONEY, AND P. DRINEAS, *An improved approximation algorithm for the column subset selection problem*, in *Proceedings of the Twentieth Annual ACM-SIAM Symposium on Discrete Algorithms*, SIAM, Philadelphia, 2009, pp. 968–977, <https://doi.org/10.1137/1.9781611973068.105>.
- [7] P. BÜRGISSE, M. CLAUSEN, AND M. A. SHOKROLLAHI, *Algebraic Complexity Theory*, Grundlehren Math. Wiss. 315, Springer-Verlag, Berlin, 1997, <https://doi.org/10.1007/978-3-662-03338-8>.
- [8] A. ÇIVRIL AND M. MAGDON-ISMAIL, *On selecting a maximum volume sub-matrix of a matrix and related problems*, *Theoret. Comput. Sci.*, 410 (2009), pp. 4801–4811, <https://doi.org/10.1016/j.tcs.2009.06.018>.
- [9] T. F. CHAN, *Rank revealing QR factorizations*, *Linear Algebra Appl.*, 88/89 (1987), pp. 67–82, [https://doi.org/10.1016/0024-3795\(87\)90103-0](https://doi.org/10.1016/0024-3795(87)90103-0).
- [10] S. CHANDRASEKARAN AND I. C. F. IPSEN, *On rank-revealing factorisations*, *SIAM J. Matrix Anal. Appl.*, 15 (1994), pp. 592–622, <https://doi.org/10.1137/S0895479891223781>.
- [11] S. CHATURANTABUT AND D. C. SORESENSEN, *Nonlinear model reduction via discrete empirical interpolation*, *SIAM J. Sci. Comput.*, 32 (2010), pp. 2737–2764, <https://doi.org/10.1137/090766498>.
- [12] A. CORTINOVIS, D. KRESSNER, AND S. MASSEI, *On maximum volume submatrices and cross approximation for symmetric semidefinite and diagonally dominant matrices*, *Linear Algebra Appl.*, 593 (2020), pp. 251–268, <https://doi.org/10.1016/j.laa.2020.02.010>.
- [13] L. DE LATHAUWER, B. DE MOOR, AND J. VANDEWALLE, *A multilinear singular value decomposition*, *SIAM J. Matrix Anal. Appl.*, 21 (2000), pp. 1253–1278, <https://doi.org/10.1137/S0895479896305696>.
- [14] A. DESHPANDE AND L. RADEMACHER, *Efficient volume sampling for row/column subset selection*, in *2010 IEEE 51st Annual Symposium on Foundations of Computer Science—FOCS 2010*, IEEE Computer Soc., Los Alamitos, CA, 2010, pp. 329–338.
- [15] A. DESHPANDE, L. RADEMACHER, S. VEMPALA, AND G. WANG, *Matrix approximation and projective clustering via volume sampling*, *Theory Comput.*, 2 (2006), pp. 225–247, <https://doi.org/10.4086/toc.2006.v002a012>.
- [16] P. DRINEAS AND M. W. MAHONEY, *A randomized algorithm for a tensor-based generalization of the singular value decomposition*, *Linear Algebra Appl.*, 420 (2007), pp. 553–571, <https://doi.org/10.1016/j.laa.2006.08.023>.
- [17] P. DRINEAS, M. W. MAHONEY, AND S. MUTHUKRISHNAN, *Relative-error CUR matrix decompositions*, *SIAM J. Matrix Anal. Appl.*, 30 (2008), pp. 844–881, <https://doi.org/10.1137/07070471X>.
- [18] Z. DRMAČ AND S. GUGERCIN, *A new selection operator for the discrete empirical interpolation method—improved a priori error bound and extensions*, *SIAM J. Sci. Comput.*, 38 (2016), pp. A631–A648, <https://doi.org/10.1137/15M1019271>.
- [19] A. K. FARAHAT, A. ELGOHARY, A. GHODSI, AND M. S. KAMEL, *Distributed column subset selection on MapReduce*, in *Proceedings of the 2013 IEEE International Conference on Data Mining*, IEEE, 2013, pp. 171–180.
- [20] A. FRIEZE, R. KANNAN, AND S. VEMPALA, *Fast Monte-Carlo algorithms for finding low-rank approximations*, *J. ACM*, 51 (2004), pp. 1025–1041, <https://doi.org/10.1145/1039488.1039494>.
- [21] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, 4th ed., Johns Hopkins Stud. Math.

- Sci., Johns Hopkins University Press, Baltimore, MD, 2013.
- [22] S. A. GOREINOV AND E. E. TYRTYSHNIKOV, *The maximal-volume concept in approximation by low-rank matrices*, in Structured Matrices in Mathematics, Computer Science, and Engineering, I (Boulder, CO, 1999), Contemp. Math. 280, AMS, Providence, RI, 2001, pp. 47–51, <https://doi.org/10.1090/conm/280/4620>.
 - [23] S. A. GOREINOV, *Cross approximation of a multi-index array*, Dokl. Akad. Nauk, 420 (2008), pp. 439–441, <https://doi.org/10.1134/S106456240803023X>.
 - [24] M. GU AND S. C. EISENSTAT, *Efficient algorithms for computing a strong rank-revealing QR factorization*, SIAM J. Sci. Comput., 17 (1996), pp. 848–869, <https://doi.org/10.1137/0917055>.
 - [25] H. HARBRECHT, M. PETERS, AND R. SCHNEIDER, *On the low-rank approximation by the pivoted Cholesky decomposition*, Appl. Numer. Math., 62 (2012), pp. 428–440, <https://doi.org/10.1016/j.apnum.2011.10.001>.
 - [26] R. A. HORN AND C. R. JOHNSON, *Matrix Analysis*, 2nd ed., Cambridge University Press, Cambridge, UK, 2013.
 - [27] S. JOSHI AND S. BOYD, *Sensor selection via convex optimization*, IEEE Trans. Signal Process., 57 (2009), pp. 451–462, <https://doi.org/10.1109/TSP.2008.2007095>.
 - [28] O. KNILL, *Cauchy-Binet for pseudo-determinants*, Linear Algebra Appl., 459 (2014), pp. 522–547, <https://doi.org/10.1016/j.laa.2014.07.013>.
 - [29] T. G. KOLDA AND B. W. BADER, *Tensor decompositions and applications*, SIAM Rev., 51 (2009), pp. 455–500, <https://doi.org/10.1137/07070111X>.
 - [30] C. MAUNG AND H. SCHWEITZER, *Pass-efficient unsupervised feature selection*, in Advances in Neural Information Processing Systems, Curran Associates, Red Hook, NY, 2013, pp. 1628–1636.
 - [31] B. ORDOZGOITI, A. MOZO, AND J. G. LÓPEZ DE LACALLE, *Regularized greedy column subset selection*, Inform. Sci., 486 (2019), pp. 393–418, <https://doi.org/10.1016/j.ins.2019.02.039>.
 - [32] I. V. OSELEDETS, D. V. SAVOSTIANOV, AND E. E. TYRTYSHNIKOV, *Tucker dimensionality reduction of three-dimensional arrays in linear time*, SIAM J. Matrix Anal. Appl., 30 (2008), pp. 939–956, <https://doi.org/10.1137/060655894>.
 - [33] R. REHMAN AND I. C. IPSEN, *La Budde's Method for Computing Characteristic Polynomials*, preprint, <https://arxiv.org/abs/1104.3769>, 2011.
 - [34] R. REHMAN AND I. C. F. IPSEN, *Computing characteristic polynomials from eigenvalues*, SIAM J. Matrix Anal. Appl., 32 (2011), pp. 90–114, <https://doi.org/10.1137/100788392>.
 - [35] A. K. SAIBABA, *HOID: Higher order interpolatory decomposition for tensors based on Tucker representation*, SIAM J. Matrix Anal. Appl., 37 (2016), pp. 1223–1249, <https://doi.org/10.1137/15M1048628>.
 - [36] D. C. SORENSEN AND M. EMBREE, *A DEIM induced CUR factorization*, SIAM J. Sci. Comput., 38 (2016), pp. A1454–A1482, <https://doi.org/10.1137/140978430>.
 - [37] G. W. STEWART, *Four algorithms for the efficient computation of truncated pivoted QR approximations to a sparse matrix*, Numer. Math., 83 (1999), pp. 313–323, <https://doi.org/10.1007/s002110050451>.
 - [38] A. TOWNSEND AND L. N. TREFETHEN, *An extension of Chebfun to two dimensions*, SIAM J. Sci. Comput., 35 (2013), pp. C495–C518, <https://doi.org/10.1137/130908002>.
 - [39] N. VANNIEUWENHOVEN, R. VANDEBRIL, AND K. MEERBERGEN, *A new truncation strategy for the higher-order singular value decomposition*, SIAM J. Sci. Comput., 34 (2012), pp. A1027–A1052, <https://doi.org/10.1137/110836067>.
 - [40] D. P. WOODRUFF, *Sketching as a tool for numerical linear algebra*, Found. Trends Theor. Comput. Sci., 10 (2014), pp. iv+157.
 - [41] P. A.-C. YOON, *Modifying Two-Sided Orthogonal Decompositions: Algorithms, Implementation, and Applications*, Ph.D. Thesis, The Pennsylvania State University, ProQuest LLC, Ann Arbor, MI, 1996, <https://www.proquest.com/docview/304312200>.
 - [42] N. L. ZAMARASHKIN AND A. I. OSINSKY, *On the existence of a nearly optimal skeleton approximation of a matrix in the Frobenius norm*, Dokl. Math., 97 (2018), pp. 164–166.