

# FINDING STRUCTURE WITH RANDOMNESS: PROBABILISTIC ALGORITHMS FOR CONSTRUCTING APPROXIMATE MATRIX DECOMPOSITIONS

N. HALKO<sup>†</sup>, P. G. MARTINSSON<sup>†</sup>, AND J. A. TROPP<sup>‡</sup>

**Abstract.** Low-rank matrix approximations, such as the truncated singular value decomposition and the rank-revealing QR decomposition, play a central role in data analysis and scientific computing. This work surveys and extends recent research which demonstrates that *randomization* offers a powerful tool for performing low-rank matrix approximation. These techniques exploit modern computational architectures more fully than classical methods and open the possibility of dealing with truly massive data sets.

This paper presents a modular framework for constructing randomized algorithms that compute partial matrix decompositions. These methods use random sampling to identify a subspace that captures most of the action of a matrix. The input matrix is then compressed—either explicitly or implicitly—to this subspace, and the reduced matrix is manipulated deterministically to obtain the desired low-rank factorization. In many cases, this approach beats its classical competitors in terms of accuracy, speed, and robustness. These claims are supported by extensive numerical experiments and a detailed error analysis.

The specific benefits of randomized techniques depend on the computational environment. Consider the model problem of finding the  $k$  dominant components of the singular value decomposition of an  $m \times n$  matrix. (i) For a dense input matrix, randomized algorithms require  $O(mn \log(k))$  floating-point operations (flops) in contrast with  $O(mnk)$  for classical algorithms. (ii) For a sparse input matrix, the flop count matches classical Krylov subspace methods, but the randomized approach is more robust and can easily be reorganized to exploit multi-processor architectures. (iii) For a matrix that is too large to fit in fast memory, the randomized techniques require only a constant number of passes over the data, as opposed to  $O(k)$  passes for classical algorithms. In fact, it is sometimes possible to perform matrix approximation with a *single pass* over the data.

**Key words.** Dimension reduction, eigenvalue decomposition, interpolative decomposition, Johnson–Lindenstrauss lemma, matrix approximation, parallel algorithm, pass-efficient algorithm, principal component analysis, randomized algorithm, random matrix, rank-revealing QR factorization, singular value decomposition, streaming algorithm.

**AMS subject classifications.** [MSC2010] Primary: 65F30. Secondary: 68W20, 60B20.

## Part I: Introduction

**1. Overview.** On a well-known list of the “Top 10 Algorithms” that have influenced the practice of science and engineering during the 20th century [40], we find an entry that is not really an algorithm: the *idea* of using matrix factorizations to accomplish basic tasks in numerical linear algebra. In the accompanying article [127], Stewart explains that

The underlying principle of the decompositional approach to matrix computation is that it is not the business of the matrix algorithmists to solve particular problems but to construct computational platforms from which a variety of problems can be solved.

Stewart goes on to argue that this point of view has had many fruitful consequences, including the development of robust software for performing these factorizations in a

---

<sup>†</sup>Department of Applied Mathematics, University of Colorado at Boulder, Boulder, CO 80309-0526. Supported by NSF awards #0748488 and #0610097.

<sup>‡</sup>Computing & Mathematical Sciences, California Institute of Technology, MC 305-16, Pasadena, CA 91125-5000. Supported by ONR award #N000140810883.

highly accurate and provably correct manner.

The decompositional approach to matrix computation remains fundamental, but developments in computer hardware and the emergence of new applications in the information sciences have rendered the classical algorithms for this task inadequate in many situations:

- A salient feature of modern applications, especially in data mining, is that the matrices are stupendously big. Classical algorithms are not always well adapted to solving the type of large-scale problems that now arise.
- In the information sciences, it is common that data are missing or inaccurate. Classical algorithms are designed to produce highly accurate matrix decompositions, but it seems profligate to spend extra computational resources when the imprecision of the data inherently limits the resolution of the output.
- Data transfer now plays a major role in the computational cost of numerical algorithms. Techniques that require few passes over the data may be substantially faster in practice, even if they require as many—or more—floating-point operations.
- As the structure of computer processors continues to evolve, it becomes increasingly important for numerical algorithms to adapt to a range of novel architectures, such as graphics processing units.

The purpose of this paper is make the case that *randomized* algorithms provide a powerful tool for constructing approximate matrix factorizations. These techniques are simple and effective, sometimes impressively so. Compared with standard deterministic algorithms, the randomized methods are often faster and—perhaps surprisingly—more robust. Furthermore, they can produce factorizations that are accurate to any specified tolerance above machine precision, which allows the user to trade accuracy for speed if desired. We present numerical evidence that these algorithms succeed for real computational problems.

In short, our goal is to demonstrate how randomized methods interact with classical techniques to yield effective, modern algorithms supported by detailed theoretical guarantees. We have made a special effort to help practitioners identify situations where randomized techniques may outperform established methods.

Throughout this article, we provide detailed citations to previous work on randomized techniques for computing low-rank approximations. The primary sources that inform our presentation include [17, 46, 58, 91, 105, 112, 113, 118, 137].

**REMARK 1.1.** Our experience suggests that many practitioners of scientific computing view randomized algorithms as a desperate and final resort. Let us address this concern immediately. Classical Monte Carlo methods are highly sensitive to the random number generator and typically produce output with low and uncertain accuracy. In contrast, the algorithms discussed herein are relatively insensitive to the quality of randomness and produce highly accurate results. The probability of failure is a user-specified parameter that can be rendered negligible (say, less than  $10^{-15}$ ) with a nominal impact on the computational resources required.

**1.1. Approximation by low-rank matrices.** The roster of standard matrix decompositions includes the pivoted QR factorization, the eigenvalue decomposition, and the singular value decomposition (SVD), all of which expose the (numerical) range of a matrix. Truncated versions of these factorizations are often used to express a

*low-rank approximation* of a given matrix:

$$\begin{array}{ccc} \mathbf{A} & \approx & \mathbf{B} \mathbf{C}, \\ m \times n & & m \times k \quad k \times n. \end{array} \quad (1.1)$$

The inner dimension  $k$  is sometimes called the *numerical rank* of the matrix. When the numerical rank is much smaller than either dimension  $m$  or  $n$ , a factorization such as (1.1) allows the matrix to be stored inexpensively and to be multiplied rapidly with vectors or other matrices. The factorizations can also be used for data interpretation or to solve computational problems, such as least squares.

Matrices with low numerical rank appear in a wide variety of scientific applications. We list only a few:

- A basic method in statistics and data mining is to compute the directions of maximal variance in vector-valued data by performing *principal component analysis* (PCA) on the data matrix. PCA is nothing other than a low-rank matrix approximation [71, §14.5].
- Another standard technique in data analysis is to perform low-dimensional embedding of data under the assumption that there are fewer degrees of freedom than the ambient dimension would suggest. In many cases, the method reduces to computing a partial SVD of a matrix derived from the data. See [71, §§14.8–14.9] or [30].
- The problem of estimating parameters from measured data via least-squares fitting often leads to very large systems of linear equations that are close to linearly dependent. Effective techniques for factoring the coefficient matrix lead to efficient techniques for solving the least-squares problem, [113].
- Many fast numerical algorithms for solving PDEs and for rapidly evaluating potential fields such as the fast multipole method [66] and  $\mathcal{H}$ -matrices [65], rely on low-rank approximations of continuum operators.
- Models of multiscale physical phenomena often involve PDEs with rapidly oscillating coefficients. Techniques for *model reduction* or *coarse graining* in such environments are often based on the observation that the linear transform that maps the input data to the requested output data can be approximated by an operator of low rank [56].

**1.2. Matrix approximation framework.** The task of computing a low-rank approximation to a given matrix can be split naturally into two computational stages. The first is to construct a low-dimensional subspace that captures the action of the matrix. The second is to restrict the matrix to the subspace and then compute a standard factorization (QR, SVD, etc.) of the reduced matrix. To be slightly more formal, we subdivide the computation as follows.

**Stage A:** Compute an approximate basis for the range of the input matrix  $\mathbf{A}$ . In other words, we require a matrix  $\mathbf{Q}$  for which

$$\mathbf{Q} \text{ has orthonormal columns and } \mathbf{A} \approx \mathbf{Q}\mathbf{Q}^* \mathbf{A}. \quad (1.2)$$

We would like the basis matrix  $\mathbf{Q}$  to contain as few columns as possible, but it is even more important to have an accurate approximation of the input matrix.

**Stage B:** Given a matrix  $\mathbf{Q}$  that satisfies (1.2), we use  $\mathbf{Q}$  to help compute a standard factorization (QR, SVD, etc.) of  $\mathbf{A}$ .

The task in Stage A can be executed very efficiently with random sampling methods, and these methods are the primary subject of this work. In the next subsection, we offer an overview of these ideas. The body of the paper provides details of the algorithms (§4) and a theoretical analysis of their performance (§§8–11).

Stage B can be completed with well-established deterministic methods. Section 3.3.3 contains an introduction to these techniques, and §5 shows how we apply them to produce low-rank factorizations.

At this point in the development, it may not be clear why the output from Stage A facilitates our job in Stage B. Let us illustrate by describing how to obtain an approximate SVD of the input matrix  $\mathbf{A}$  given a matrix  $\mathbf{Q}$  that satisfies (1.2). More precisely, we wish to compute matrices  $\mathbf{U}$  and  $\mathbf{V}$  with orthonormal columns and a nonnegative, diagonal matrix  $\mathbf{\Sigma}$  such that  $\mathbf{A} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ . This goal is achieved after three simple steps:

1. Form  $\mathbf{B} = \mathbf{Q}^* \mathbf{A}$ , which yields the low-rank factorization  $\mathbf{A} \approx \mathbf{Q}\mathbf{B}$ .
2. Compute an SVD of the small matrix:  $\mathbf{B} = \tilde{\mathbf{U}}\mathbf{\Sigma}\mathbf{V}^*$ .
3. Set  $\mathbf{U} = \mathbf{Q}\tilde{\mathbf{U}}$ .

When  $\mathbf{Q}$  has few columns, this procedure is efficient because we can easily construct the reduced matrix  $\mathbf{B}$  and rapidly compute its SVD. In practice, we can often avoid forming  $\mathbf{B}$  explicitly by means of subtler techniques. In some cases, it is not even necessary to revisit the input matrix  $\mathbf{A}$  during Stage B. This observation allows us to develop *single-pass algorithms*, which look at each entry of  $\mathbf{A}$  only once.

Similar manipulations readily yield other standard factorizations, such as the pivoted QR factorization, the eigenvalue decomposition, etc.

**1.3. Randomized algorithms.** This paper describes a class of randomized algorithms for completing Stage A of the matrix approximation framework set forth in §1.2. We begin with some details about the approximation problem these algorithms target (§1.3.1). Afterward, we motivate the random sampling technique with a heuristic explanation (§1.3.2) that leads to a prototype algorithm (§1.3.3).

**1.3.1. Problem formulations.** The basic challenge in producing low-rank matrix approximations is a primitive question that we call the *fixed-precision approximation problem*. Suppose we are given a matrix  $\mathbf{A}$  and a positive error tolerance  $\varepsilon$ . We seek a matrix  $\mathbf{Q}$  with  $k = k(\varepsilon)$  orthonormal columns such that

$$\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^* \mathbf{A}\| \leq \varepsilon, \quad (1.3)$$

where  $\|\cdot\|$  denotes the  $\ell_2$  operator norm. The range of  $\mathbf{Q}$  is a  $k$ -dimensional subspace that captures most of the action of  $\mathbf{A}$ , and we would like  $k$  to be as small as possible.

The singular value decomposition furnishes an optimal answer to the fixed-precision problem [97]. Let  $\sigma_j$  denote the  $j$ th largest singular value of  $\mathbf{A}$ . For each  $j \geq 0$ ,

$$\min_{\text{rank}(\mathbf{X}) \leq j} \|\mathbf{A} - \mathbf{X}\| = \sigma_{j+1}. \quad (1.4)$$

One way to construct a minimizer is to choose  $\mathbf{X} = \mathbf{Q}\mathbf{Q}^* \mathbf{A}$ , where the columns of  $\mathbf{Q}$  are  $k$  dominant left singular vectors of  $\mathbf{A}$ . Consequently, the minimal rank  $k$  where (1.3) holds equals the number of singular values of  $\mathbf{A}$  that exceed the tolerance  $\varepsilon$ .

To simplify the development of algorithms, it is convenient to assume that the desired rank  $k$  is specified in advance. We call the resulting problem the *fixed-rank approximation problem*. Given a matrix  $\mathbf{A}$ , a target rank  $k$ , and an oversampling

parameter  $p$ , we seek to construct a matrix  $\mathbf{Q}$  with  $k + p$  orthonormal columns such that

$$\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^*\mathbf{A}\| \approx \min_{\text{rank}(\mathbf{X}) \leq k} \|\mathbf{A} - \mathbf{X}\|. \quad (1.5)$$

Although there exists a minimizer  $\mathbf{Q}$  that solves the fixed rank problem for  $p = 0$ , the opportunity to use a small number of additional columns provides a flexibility that is crucial for the effectiveness of the computational methods we discuss.

We will demonstrate that algorithms for the fixed-rank problem can be adapted to solve the fixed-precision problem. The connection is based on the observation that we can build the basis matrix  $\mathbf{Q}$  incrementally and, at any point in the computation, we can inexpensively estimate the residual error  $\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^*\mathbf{A}\|$ . Refer to §4.4 for the details of this reduction.

**1.3.2. Intuition.** To understand how randomness helps us solve the fixed-rank problem, it is helpful to consider some motivating examples.

First, suppose that we seek a basis for the range of a matrix  $\mathbf{A}$  with *exact* rank  $k$ . Draw a random vector  $\boldsymbol{\omega}$ , and form the product  $\mathbf{y} = \mathbf{A}\boldsymbol{\omega}$ . For now, the precise distribution of the random vector is unimportant; just think of  $\mathbf{y}$  as a random sample from the range of  $\mathbf{A}$ . Let us repeat this sampling process  $k$  times:

$$\mathbf{y}^{(i)} = \mathbf{A}\boldsymbol{\omega}^{(i)}, \quad i = 1, 2, \dots, k. \quad (1.6)$$

Owing to the randomness, the set  $\{\boldsymbol{\omega}^{(i)} : i = 1, 2, \dots, k\}$  of random vectors is likely to be in general linear position. In particular, the random vectors form a linearly independent set and no linear combination falls in the null space of  $\mathbf{A}$ . As a result, the set  $\{\mathbf{y}^{(i)} : i = 1, 2, \dots, k\}$  of sample vectors is also linearly independent, so it spans the range of  $\mathbf{A}$ . Therefore, to produce an orthonormal basis for the range of  $\mathbf{A}$ , we just need to orthonormalize the sample vectors.

Now, imagine that  $\mathbf{A} = \mathbf{B} + \mathbf{E}$  where  $\mathbf{B}$  is a rank- $k$  matrix containing the information we seek and  $\mathbf{E}$  is a small perturbation. Our priority is to obtain a basis that covers as much of the range of  $\mathbf{B}$  as possible, rather than to minimize the number of basis vectors. Therefore, we fix a small number  $p$ , and we generate  $k + p$  samples

$$\mathbf{y}^{(i)} = \mathbf{A}\boldsymbol{\omega}^{(i)} = \mathbf{B}\boldsymbol{\omega}^{(i)} + \mathbf{E}\boldsymbol{\omega}^{(i)}, \quad i = 1, 2, \dots, k + p. \quad (1.7)$$

The perturbation  $\mathbf{E}$  shifts the direction of each sample vector outside the range of  $\mathbf{B}$ , which can prevent the span of  $\{\mathbf{y}^{(i)} : i = 1, 2, \dots, k\}$  from covering the entire range of  $\mathbf{B}$ . In contrast, the enriched set  $\{\mathbf{y}^{(i)} : i = 1, 2, \dots, k + p\}$  of samples has a much better chance of spanning the required subspace.

Just how many extra samples do we need? Remarkably, for certain types of random sampling schemes, the failure probability decreases superexponentially with the oversampling parameter  $p$ ; see (1.9). As a practical matter, setting  $p = 5$  or  $p = 10$  often gives superb results. This observation is one of the principal facts supporting the randomized approach to numerical linear algebra.

**1.3.3. A prototype algorithm.** The intuitive approach of §1.3.2 can be applied to general matrices. Omitting computational details for now, we formalize the procedure in the figure labeled Proto-Algorithm.

This simple algorithm is by no means new. It is essentially the first step of a subspace iteration with a random initial subspace [61, §7.3.2]. The novelty comes from the additional observation that the initial subspace should have a slightly higher

PROTO-ALGORITHM: SOLVING THE FIXED-RANK PROBLEM

*Given an  $m \times n$  matrix  $\mathbf{A}$ , a target rank  $k$ , and an oversampling parameter  $p$ , this procedure computes an  $m \times (k + p)$  matrix  $\mathbf{Q}$  whose columns are orthonormal and whose range approximates the range of  $\mathbf{A}$ .*

- 1 Draw a random  $n \times (k + p)$  test matrix  $\mathbf{\Omega}$ .
- 2 Form the matrix product  $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$ .
- 3 Construct a matrix  $\mathbf{Q}$  whose columns form an orthonormal basis for the range of  $\mathbf{Y}$ .

dimension than the invariant subspace we are trying to approximate. With this revision, it is often the case that *no further iteration is required* to obtain a high-quality solution to (1.5). We believe this idea can be traced to [91, 105, 118].

In order to invoke the proto-algorithm with confidence, we must address several practical and theoretical issues:

- What random matrix  $\mathbf{\Omega}$  should we use? How much oversampling do we need?
- The matrix  $\mathbf{Y}$  is likely to be ill-conditioned. How do we orthonormalize its columns to form the matrix  $\mathbf{Q}$ ?
- What are the computational costs?
- How can we solve the fixed-precision problem (1.3) when the numerical rank of the matrix is not known in advance?
- How can we use the basis  $\mathbf{Q}$  to compute other matrix factorizations?
- Does the randomized method work for problems of practical interest? How does its speed/accuracy/robustness compare with standard techniques?
- What error bounds can we expect? With what probability?

The next few sections provide a summary of the answers to these questions. We describe several problem regimes where the proto-algorithm can be implemented efficiently, and we present a theorem that describes the performance of the most important instantiation. Finally, we elaborate on how these ideas can be applied to approximate the truncated SVD of a large data matrix. The rest of the paper contains a more exhaustive treatment—including pseudocode, numerical experiments, and a detailed theory.

**1.4. A comparison between randomized and traditional techniques.** To select an appropriate computational method for finding a low-rank approximation to a matrix, the practitioner must take into account the properties of the matrix. Is it dense or sparse? Does it fit in fast memory or is it stored out of core? Does the singular spectrum decay quickly or slowly? The behavior of a numerical linear algebra algorithm may depend on all these factors [13, 61, 132]. To facilitate a comparison between classical and randomized techniques, we summarize their relative performance in each of three representative environments. Section 6 contains a more in-depth treatment.

We focus on the task of computing an approximate SVD of an  $m \times n$  matrix  $\mathbf{A}$  with numerical rank  $k$ . For randomized schemes, Stage A generally dominates the cost of Stage B in our matrix approximation framework (§1.2). Within Stage A, the computational bottleneck is usually the matrix–matrix product  $\mathbf{A}\mathbf{\Omega}$  in Step 2 of the proto-algorithm (§1.3.3). The power of randomized algorithms stems from the fact

that we can reorganize this matrix multiplication for maximum efficiency in a variety of computational architectures.

**1.4.1. A general dense matrix that fits in fast memory.** A standard deterministic technique for computing an approximate SVD is to perform a rank-revealing QR factorization of the matrix, and then to manipulate the factors to obtain the final decomposition. The cost of this approach is typically  $O(kmn)$  floating-point operations, or *flops*, although these methods require slightly longer running times in rare cases [68].

In contrast, randomized schemes can produce an approximate SVD using only  $O(mn \log(k) + (m+n)k^2)$  flops. The gain in asymptotic complexity is achieved by using a random matrix  $\Omega$  that has some internal structure, which allows us to evaluate the product  $A\Omega$  rapidly. For example, randomizing and subsampling the discrete Fourier transform works well. Sections 4.6 and 11 contain more information on this approach.

**1.4.2. A matrix for which matrix–vector products can be evaluated rapidly.** When the matrix  $A$  is sparse or structured, we may be able to apply it rapidly to a vector. In this case, the classical prescription for computing a partial SVD is to invoke a Krylov subspace method, such as the Lanczos or Arnoldi algorithm. It is difficult to summarize the computational cost of these methods because their performance depends heavily on properties of the input matrix and on the amount of effort spent to stabilize the algorithm. (Inherently, the Lanczos and Arnoldi methods are numerically unstable.) For the same reasons, the error analysis of such schemes is unsatisfactory in many important environments.

At the risk of being overly simplistic, we claim that the typical cost of a Krylov method for approximating the  $k$  leading singular vectors of the input matrix is proportional to  $kT_{\text{mult}} + (m+n)k^2$ , where  $T_{\text{mult}}$  denotes the cost of a matrix–vector multiplication with the input matrix and the constant of proportionality is small. We can also apply randomized methods using a Gaussian test matrix  $\Omega$  to complete the factorization at the same cost,  $O(kT_{\text{mult}} + (m+n)k^2)$  flops.

With a given budget of floating-point operations, Krylov methods sometimes deliver a more accurate approximation than randomized algorithms. Nevertheless, the methods described in this survey have at least two powerful advantages over Krylov methods. First, the randomized schemes are inherently stable, and they come with very strong performance guarantees that do not depend on subtle spectral properties of the input matrix. Second, the matrix–vector multiplies required to form  $A\Omega$  can be performed *in parallel*. This fact allows us to restructure the calculations to take full advantage of the computational platform, which can lead to dramatic accelerations in practice, especially for parallel and distributed machines.

A more detailed comparison of randomized schemes and Krylov subspace methods is given in §6.2.

**1.4.3. A general dense matrix stored in slow memory or streamed.** When the input matrix is too large to fit in core memory, the cost of transferring the matrix from slow memory typically dominates the cost of performing the arithmetic. The standard techniques for low-rank approximation described in §1.4.1 require  $O(k)$  passes over the matrix, which can be prohibitively expensive.

In contrast, the proto-algorithm of §1.3.3 requires only one pass over the data to produce the approximate basis  $Q$  for Stage A of the approximation framework. This straightforward approach, unfortunately, is not accurate enough for matrices whose



singular spectrum decays slowly, but we can address this problem using very few (say, 2 to 4) additional passes over the data [112]. See §1.6 or §4.5 for more discussion.

Typically, Stage B uses one additional pass over the matrix to construct the approximate SVD. With slight modifications, however, the two-stage randomized scheme can be revised so that it only makes a single pass over the data. Refer to §5.5 for information.

**1.5. Performance analysis.** A principal goal of this paper is to provide a detailed analysis of the performance of the proto-algorithm described in §1.3.3. This investigation produces precise error bounds, expressed in terms of the singular values of the input matrix. Furthermore, we determine how several choices of the random matrix  $\Omega$  impact the behavior of the algorithm.

Let us offer a taste of this theory. The following theorem describes the average-case behavior of the proto-algorithm with a Gaussian test matrix, assuming we perform the computation in exact arithmetic. This result is a simplified version of Theorem 10.6.

**THEOREM 1.1.** *Suppose that  $\mathbf{A}$  is a real  $m \times n$  matrix. Select a target rank  $k \geq 2$  and an oversampling parameter  $p \geq 2$ , where  $k + p \leq \min\{m, n\}$ . Execute the proto-algorithm with a standard Gaussian test matrix to obtain an  $m \times (k + p)$  matrix  $\mathbf{Q}$  with orthonormal columns. Then*

$$\mathbb{E} \|\mathbf{A} - \mathbf{Q}\mathbf{Q}^* \mathbf{A}\| \leq \left[ 1 + \frac{4\sqrt{k+p}}{p-1} \cdot \sqrt{\min\{m, n\}} \right] \sigma_{k+1}, \quad (1.8)$$

where  $\mathbb{E}$  denotes expectation with respect to the random test matrix and  $\sigma_{k+1}$  is the  $(k+1)$ th singular value of  $\mathbf{A}$ .

We recall that the term  $\sigma_{k+1}$  appearing in (1.8) is the smallest possible error (1.4) achievable with any basis matrix  $\mathbf{Q}$ . The theorem asserts that, on average, the algorithm produces a basis whose error lies within a small polynomial factor of the theoretical minimum. Moreover, the error bound (1.8) in the randomized algorithm is slightly sharper than comparable bounds for deterministic techniques based on rank-revealing QR algorithms [68].

The reader might be worried about whether the expectation provides a useful account of the approximation error. Fear not: the actual outcome of the algorithm is *almost always* very close to the typical outcome because of measure concentration effects. As we discuss in §10.3, the probability that the error satisfies

$$\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^* \mathbf{A}\| \leq \left[ 1 + 11\sqrt{k+p} \cdot \sqrt{\min\{m, n\}} \right] \sigma_{k+1} \quad (1.9)$$

is at least  $1 - 6 \cdot p^{-p}$  under very mild assumptions on  $p$ . This fact justifies the use of an oversampling term as small as  $p = 5$ . This simplified estimate is very similar to the major results in [91].

The theory developed in this paper provides much more detailed information about the performance of the proto-algorithm.

- When the singular values of  $\mathbf{A}$  decay slightly, the error  $\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^* \mathbf{A}\|$  does not depend on the dimensions of the matrix (§§10.2–10.3).
- We can reduce the size of the bracket in the error bound (1.8) by combining the proto-algorithm with a power iteration (§10.4). For an example, see §1.6 below.



PROTOTYPE FOR RANDOMIZED SVD

Given an  $m \times n$  matrix  $\mathbf{A}$ , a target number  $k$  of singular vectors, and an exponent  $q$  (say  $q = 1$  or  $q = 2$ ), this procedure computes an approximate rank- $2k$  factorization  $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ , where  $\mathbf{U}$  and  $\mathbf{V}$  are orthonormal, and  $\mathbf{\Sigma}$  is nonnegative and diagonal.

**Stage A:**

- 1 Generate an  $n \times 2k$  Gaussian test matrix  $\mathbf{\Omega}$ .
- 2 Form  $\mathbf{Y} = (\mathbf{A}\mathbf{A}^*)^q \mathbf{A}\mathbf{\Omega}$  by multiplying alternately with  $\mathbf{A}$  and  $\mathbf{A}^*$ .
- 3 Construct a matrix  $\mathbf{Q}$  whose columns form an orthonormal basis for the range of  $\mathbf{Y}$ .

**Stage B:**

- 4 Form  $\mathbf{B} = \mathbf{Q}^* \mathbf{A}$ .
- 5 Compute an SVD of the small matrix:  $\mathbf{B} = \tilde{\mathbf{U}}\mathbf{\Sigma}\mathbf{V}^*$ .
- 6 Set  $\mathbf{U} = \mathbf{Q}\tilde{\mathbf{U}}$ .

**Note:** The computation of  $\mathbf{Y}$  in Step 2 is vulnerable to round-off errors. When high accuracy is required, we must incorporate an orthonormalization step between each application of  $\mathbf{A}$  and  $\mathbf{A}^*$ ; see Algorithm 4.4.

- For the structured random matrices we mentioned in §1.4.1, related error bounds are in force (§11).
- We can obtain inexpensive *a posteriori* error estimates to verify the quality of the approximation (§4.3).

**1.6. Example: Randomized SVD.** We conclude this introduction with a short discussion of how these ideas allow us to perform an approximate SVD of a large data matrix, which is a compelling application of randomized matrix approximation [112].

The two-stage randomized method offers a natural approach to SVD computations. Unfortunately, the simplest version of this scheme is inadequate in many applications because the singular spectrum of the input matrix may decay slowly. To address this difficulty, we incorporate  $q$  steps of a power iteration, where  $q = 1$  or  $q = 2$  usually suffices in practice. The complete scheme appears in the box labeled Prototype for Randomized SVD. For most applications, it is important to incorporate additional refinements, as we discuss in §§4–5.

The Randomized SVD procedure requires only  $2(q+1)$  passes over the matrix, so it is efficient even for matrices stored out-of-core. The flop count satisfies

$$T_{\text{randSVD}} = (2q+2)kT_{\text{mult}} + O(k^2(m+n)),$$

where  $T_{\text{mult}}$  is the flop count of a matrix–vector multiply with  $\mathbf{A}$  or  $\mathbf{A}^*$ . We have the following theorem on the performance of this method in exact arithmetic, which is a consequence of Corollary 10.10.

**THEOREM 1.2.** *Suppose that  $\mathbf{A}$  is a real  $m \times n$  matrix. Select an exponent  $q$  and a target number  $k$  of singular vectors, where  $2 \leq k \leq 0.5 \min\{m, n\}$ . Execute the*

Randomized SVD algorithm to obtain a rank- $2k$  factorization  $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ . Then

$$\mathbb{E} \|\mathbf{A} - \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*\| \leq \left[ 1 + 4\sqrt{\frac{2\min\{m,n\}}{k-1}} \right]^{1/(2q+1)} \sigma_{k+1}, \quad (1.10)$$

where  $\mathbb{E}$  denotes expectation with respect to the random test matrix and  $\sigma_{k+1}$  is the  $(k+1)$ th singular value of  $\mathbf{A}$ .

This result is new. Observe that the bracket in (1.10) is essentially the same as the bracket in the basic error bound (1.8). We find that the power iteration drives the leading constant to one exponentially fast as the power  $q$  increases. The rank- $k$  approximation of  $\mathbf{A}$  can never achieve an error smaller than  $\sigma_{k+1}$ , so the randomized procedure computes  $2k$  approximate singular vectors that capture as much of the matrix as the first  $k$  actual singular vectors.

In practice, we can truncate the approximate SVD, retaining only the first  $k$  singular values and vectors. Equivalently, we replace the diagonal factor  $\mathbf{\Sigma}$  by the matrix  $\mathbf{\Sigma}_{(k)}$  formed by zeroing out all but the largest  $k$  entries of  $\mathbf{\Sigma}$ . For this truncated SVD, we have the error bound

$$\mathbb{E} \|\mathbf{A} - \mathbf{U}\mathbf{\Sigma}_{(k)}\mathbf{V}^*\| \leq \sigma_{k+1} + \left[ 1 + 4\sqrt{\frac{2\min\{m,n\}}{k-1}} \right]^{1/(2q+1)} \sigma_{k+1}. \quad (1.11)$$

In words, we pay no more than an additive term  $\sigma_{k+1}$  when we perform the truncation step. Our numerical experience suggests that the error bound (1.11) is pessimistic. See Remark 5.1 and §9.4 for some discussion of truncation.

**1.7. Outline of paper.** The paper is organized into three parts: an introduction (§§1–3), a description of the algorithms (§§4–7), and a theoretical performance analysis (§§8–11). The two latter parts commence with a short internal outline. Each part is more or less self-contained, and after a brief review of our notation in §§3.1–3.2, the reader can proceed to either the algorithms or the theory part.

**2. Related work and historical context.** Randomness has occasionally surfaced in the numerical linear algebra literature; in particular, it is quite standard to initialize iterative algorithms for constructing invariant subspaces with a randomly chosen point. Nevertheless, we believe that sophisticated ideas from random matrix theory have not been incorporated into classical matrix factorization algorithms until very recently. We can trace this development to earlier work in computer science and—especially—to probabilistic methods in geometric analysis. This section presents an overview of the relevant work. We begin with a survey of randomized methods for matrix approximation; then we attempt to trace some of the ideas backward to their sources.

**2.1. Randomized matrix approximation.** Matrices of low numerical rank contain little information relative to their apparent dimension owing to the linear dependency in their columns (or rows). As a result, it is reasonable to expect that these matrices can be approximated with far fewer degrees of freedom. A less obvious fact is that randomized schemes can be used to produce these approximations efficiently.

Several types of approximation techniques build on this idea. These methods all follow the same basic pattern:

1. Preprocess the matrix, usually to calculate sampling probabilities.
2. Take random samples from the matrix, where the term *sample* refers generically to a linear function of the matrix.
3. Postprocess the samples to compute a final approximation, typically with classical techniques from numerical linear algebra. This step may require another look at the matrix.

We continue with a description of the most common approximation schemes.

**2.1.1. Sparsification.** The simplest approach to matrix approximation is the method of *sparsification* or the related technique of *quantization*. The goal of sparsification is to replace the matrix by a surrogate that contains far fewer nonzero entries. Quantization produces an approximation whose components are drawn from a (small) discrete set of values. These methods can be used to limit storage requirements or to accelerate computations by reducing the cost of matrix–vector and matrix–matrix multiplies [94, Ch. 6]. The manuscript [33] describes applications in optimization.

Sparsification typically involves very simple elementwise calculations. Each entry in the approximation is drawn independently at random from a distribution determined from the corresponding entry of the input matrix. The expected value of the random approximation equals the original matrix, but the distribution is designed so that a typical realization is much sparser.

The first method of this form was devised by Achlioptas and McSherry [2], who built on earlier work on graph sparsification due to Karger [75, 76]. Arora–Hazan–Kale presented a different sampling method in [7]. See [60, 123] for some recent work on sparsification.

**2.1.2. Column selection methods.** A second approach to matrix approximation is based on the idea that a small set of columns describes most of the action of a numerically low-rank matrix. Indeed, classical existential results [117] demonstrate that every  $m \times n$  matrix  $\mathbf{A}$  contains a  $k$ -column submatrix  $\mathbf{C}$  for which

$$\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger\mathbf{A}\| \leq \sqrt{1 + k(n - k)} \cdot \|\mathbf{A} - \mathbf{A}_{(k)}\|, \quad (2.1)$$

where  $k$  is a parameter, the dagger  $\dagger$  denotes the pseudoinverse, and  $\mathbf{A}_{(k)}$  is a best rank- $k$  approximation of  $\mathbf{A}$ . It is NP-hard to perform column selection by optimizing natural objective functions, such as the condition number of the submatrix [27]. Nevertheless, there are efficient deterministic algorithms, such as the rank-revealing QR method of [68], that can nearly achieve the error bound (2.1).

There is a class of randomized algorithms that approach the fixed-rank approximation problem (1.5) using this intuition. These methods first compute a sampling probability for each column, either using the squared Euclidean norms of the columns or their *leverage scores*. (Leverage scores reflect the relative importance of the columns to the action of the matrix; they can be calculated easily from the dominant  $k$  right singular vectors of the matrix.) Columns are then selected randomly according to this distribution. Afterward, a postprocessing step is invoked to produce a more refined approximation of the matrix.

We believe that the earliest method of this form appeared in a 1998 paper of Frieze–Kannan–Vempala [57, 58]. This work was refined substantially in the papers [43, 44, 46]. The basic algorithm samples columns from a distribution related to the squared  $\ell_2$  norms of the columns. This sampling step produces a small column submatrix whose range is aligned with the range of the input matrix. The final approximation is obtained from a truncated SVD of the submatrix. Given a target rank

$k$  and a parameter  $\varepsilon > 0$ , this approach samples  $\ell = \ell(k, \varepsilon)$  columns of the matrix to produce a rank- $k$  approximation  $\mathbf{B}$  that satisfies

$$\|\mathbf{A} - \mathbf{B}\|_F \leq \|\mathbf{A} - \mathbf{A}_{(k)}\|_F + \varepsilon \|\mathbf{A}\|_F, \quad (2.2)$$

where  $\|\cdot\|_F$  denotes the Frobenius norm. We note that the algorithm of [46] requires only a constant number of passes over the data.

Rudelson and Vershynin later showed that the same type of column sampling method also yields spectral-norm error bounds [116]. The techniques in their paper have been very influential; their work has found other applications in randomized regression [52], sparse approximation [133], and compressive sampling [19].

Deshpande et al. [37, 38] demonstrated that the error in the column sampling approach can be improved by iteration and adaptive volume sampling. They showed that it is possible to produce a rank- $k$  matrix  $\mathbf{B}$  that satisfies

$$\|\mathbf{A} - \mathbf{B}\|_F \leq (1 + \varepsilon) \|\mathbf{A} - \mathbf{A}_{(k)}\|_F \quad (2.3)$$

using a  $k$ -pass algorithm. Around the same time, Har-Peled [70] independently developed a recursive algorithm that offers the same approximation guarantees. Very recently, Deshpande and Rademacher have improved the running time of volume-based sampling methods [36].

Drineas et al. and Boutsidis et al. have also developed randomized algorithms for the *column subset selection problem*, which requests a column submatrix  $\mathbf{C}$  that achieves a bound of the form (2.1). Via the methods of Rudelson and Vershynin [116], they showed that sampling columns according to their leverage scores is likely to produce the required submatrix [50, 51]. Subsequent work [17, 18] showed that post-processing the sampled columns with a rank-revealing QR algorithm can reduce the number of output columns required (2.1). The argument in [17] explicitly decouples the linear algebraic part of the analysis from the random matrix theory. The theoretical analysis in the present work involves a very similar technique.

**2.1.3. Approximation by dimension reduction.** A third approach to matrix approximation is based on the concept of *dimension reduction*. Since the rows of a low-rank matrix are linearly dependent, they can be embedded into a low-dimensional space without altering their geometric properties substantially. A random linear map provides an efficient, nonadaptive way to perform this embedding. (Column sampling can also be viewed as an adaptive form of dimension reduction.)

The proto-algorithm we set forth in §1.3.3 is simply a dual description of the dimension reduction approach: collecting random samples from the column space of the matrix is equivalent to reducing the dimension of the rows. No precomputation is required to obtain the sampling distribution, but the sample itself takes some work to collect. Afterward, we orthogonalize the samples as preparation for constructing various matrix approximations.

We believe that the idea of using dimension reduction for algorithmic matrix approximation first appeared in a 1998 paper of Papadimitriou et al. [104, 105], who described an application to latent semantic indexing (LSI). They suggested projecting the input matrix onto a random subspace and compressing the original matrix to (a subspace of) the range of the projected matrix. They established error bounds that echo the result (2.2) of Frieze et al. [58]. Although the Euclidean column selection method is a more computationally efficient way to obtain this type of error bound, dimension reduction has other advantages, e.g., in terms of accuracy.

Sarlós argued in [118] that the computational costs of dimension reduction can be reduced substantially by means of the structured random maps proposed by Ailon–Chazelle [3]. Sarlós used these ideas to develop efficient randomized algorithms for least-squares problems; he also studied approximate matrix multiplication and low-rank matrix approximation. The recent paper [102] analyzes a very similar matrix approximation algorithm using Rudelson and Vershynin’s methods [116].

The initial work of Sarlós on structured dimension reduction did not immediately yield algorithms for low-rank matrix approximation that were superior to classical techniques. Woolfe et al. showed how to obtain an improvement in asymptotic computational cost, and they applied these techniques to problems in scientific computing [137]. Related work includes [86, 88].

Martinsson–Rokhlin–Tygert have studied dimension reduction using a Gaussian transform matrix, and they demonstrated that this approach performs much better than earlier analyses had suggested [91]. Their work highlights the importance of over-sampling, and their error bounds are very similar to the estimate (1.9) we presented in the introduction. They also demonstrated that dimension reduction can be used to compute an interpolative decomposition of the input matrix, which is essentially equivalent to performing column subset selection.

Rokhlin–Szlam–Tygert have shown that combining dimension reduction with a power iteration is an effective way to improve its performance [112]. These ideas lead to very efficient randomized methods for large-scale PCA [69]. An efficient, numerically stable version of the power iteration is discussed in §4.5, as well as [92]. Related ideas appear in a paper of Roweis [114].

Very recently, Clarkson and Woodruff [29] have developed one-pass algorithms for performing low-rank matrix approximation, and they have established lower bounds which prove that many of their algorithms have optimal or near-optimal resource guarantees, modulo constants.

**2.1.4. Approximation by submatrices.** The matrix approximation literature contains a subgenre that discusses methods for building an approximation from a submatrix and computed coefficient matrices. For example, we can construct an approximation using a subcollection of columns (the interpolative decomposition), a subcollection of rows and a subcollection of columns (the CUR decomposition), or a square submatrix (the matrix skeleton). This type of decomposition was developed and studied in several papers, including [26, 64, 126]. For data analysis applications, see the recent paper [89].

A number of works develop randomized algorithms for this class of matrix approximations. Drineas et al. have developed techniques for computing CUR decompositions, which express  $\mathbf{A} \approx \mathbf{C}\mathbf{U}\mathbf{R}$ , where  $\mathbf{C}$  and  $\mathbf{R}$  denote small column and row submatrices of  $\mathbf{A}$  and where  $\mathbf{U}$  is a small linkage matrix. These methods identify columns (rows) that approximate the range (corange) of the matrix; the linkage matrix is then computed by solving a small least-squares problem. A randomized algorithm for CUR approximation with controlled absolute error appears in [47]; a relative error algorithm appears in [51]. We also mention a paper on computing a closely related factorization called the *compact matrix decomposition* [129].

It is also possible to produce interpolative decompositions and matrix skeletons using randomized methods, as discussed in [91, 112] and §5.2 of the present work.

**2.1.5. Other numerical problems.** The literature contains a variety of other randomized algorithms for solving standard problems in and around numerical linear algebra. We list some of the basic references.

**Tensor skeletons.** Randomized column selection methods can be used to produce CUR-type decompositions of higher-order tensors [49].

**Matrix multiplication.** Column selection and dimension reduction techniques can be used to accelerate the multiplication of rank-deficient matrices [45, 118]. See also [10].

**Overdetermined linear systems.** The randomized Kaczmarz algorithm is a linearly convergent iterative method that can be used to solve overdetermined linear systems [101, 128].

**Overdetermined least squares.** Fast dimension-reduction maps can sometimes accelerate the solution of overdetermined least-squares problems [52, 118].

**Nonnegative least squares.** Fast dimension reduction maps can be used to reduce the size of nonnegative least-squares problems [16].

**Preconditioned least squares.** Randomized matrix approximations can be used to precondition conjugate gradient to solve least-squares problems [113].

**Other regression problems.** Randomized algorithms for  $\ell_1$  regression are described in [28]. Regression in  $\ell_p$  for  $p \in [1, \infty)$  has also been considered [31].

**Facility location.** The Fermat–Weber facility location problem can be viewed as matrix approximation with respect to a different discrepancy measure. Randomized algorithms for this type of problem appear in [121].

**2.1.6. Compressive sampling.** Although randomized matrix approximation and compressive sampling are based on some common intuitions, it is facile to consider either one as a subspecies of the other. We offer a short overview of the field of compressive sampling—especially the part connected with matrices—so we can highlight some of the differences.

The theory of compressive sampling starts with the observation that many types of vector-space data are *compressible*. That is, the data are approximated well using a short linear combination of basis functions drawn from a fixed collection [42]. For example, natural images are well approximated in a wavelet basis; numerically low-rank matrices are well approximated as a sum of rank-one matrices. The idea behind compressive sampling is that suitably chosen random samples from this type of compressible object carry a large amount of information. Furthermore, it is possible to reconstruct the compressible object from a small set of these random samples, often by solving a convex optimization problem. The initial discovery works of Candès–Romberg–Tao [20] and Donoho [41] were written in 2004.

The earliest work in compressive sampling focused on vector-valued data; soon after, researchers began to study compressive sampling for matrices. In 2007, Recht–Fazel–Parillo demonstrated that it is possible to reconstruct a rank-deficient matrix from Gaussian measurements [111]. More recently, Candès–Recht [22] and Candès–Tao [23] considered the problem of completing a low-rank matrix from a random sample of its entries.

The usual goals of compressive sampling are (i) to design a method for collecting informative, nonadaptive data about a compressible object and (ii) to reconstruct a compressible object given some measured data. In both cases, there is an implicit assumption that we have limited—if any—access to the underlying data.

In the problem of matrix approximation, we typically have a complete representation of the matrix at our disposal. The point is to compute a simpler representation as efficiently as possible under some operational constraints. In particular, we would like to perform as little computation as we can, but we are usually allowed to revisit the input matrix. Because of the different focus, randomized matrix approximation algo-



gorithms require fewer random samples from the matrix and use fewer computational resources than compressive sampling reconstruction algorithms.

**2.2. Origins.** This section attempts to identify some of the major threads of research that ultimately led to the development of the randomized techniques we discuss in this paper.

**2.2.1. Random embeddings.** The field of random embeddings is a major precursor to randomized matrix approximation. In a celebrated 1984 paper [74], Johnson and Lindenstrauss showed that the pairwise distances among a collection of  $N$  points in a Euclidean space are approximately maintained when the points are mapped randomly to a Euclidean space of dimension  $O(\log N)$ . In other words, random embeddings preserve Euclidean geometry. Shortly afterward, Bourgain showed that appropriate random low-dimensional embeddings preserve the geometry of point sets in finite-dimensional  $\ell_1$  spaces [15].

These observations suggest that we might be able to solve some computational problems of a geometric nature more efficiently by translating them into a lower-dimensional space and solving them there. This idea was cultivated by the theoretical computer science community beginning in the late 1980s, with research flowering in the late 1990s. In particular, nearest-neighbor search can benefit from dimension-reduction techniques [73, 78, 80]. The papers [57, 104] were apparently the first to apply this approach to linear algebra.

Around the same time, researchers became interested in simplifying the form of dimension reduction maps and improving the computational cost of applying the map. Several researchers developed refined results on the performance of a Gaussian matrix as a linear dimension reduction map [32, 73, 93]. Achlioptas demonstrated that discrete random matrices would serve nearly as well [1]. In 2006, Ailon and Chazelle proposed the *fast Johnson–Lindenstrauss transform* [3], which combines the speed of the FFT with the favorable embedding properties of a Gaussian matrix. Subsequent refinements appear in [4, 87]. Sarlós then imported these techniques to study several problems in numerical linear algebra, which has led to some of the fastest algorithms currently available [88, 137].

**2.2.2. Data streams.** Muthukrishnan argues that a distinguishing feature of modern data is the manner in which it is *presented* to us. The sheer volume of information and the speed at which it must be processed tax our ability to *transmit* the data elsewhere, to *compute* complicated functions on the data, or to *store* a substantial part of the data [100, §3]. As a result, computer scientists have started to develop algorithms that can address familiar computational problems under these novel constraints. The data stream phenomenon is one of the primary justifications cited by [45] for developing pass-efficient methods for numerical linear algebra problems, and it is also the focus of the recent treatment [29].

One of the methods for dealing with massive data sets is to maintain *sketches*, which are small summaries that allow functions of interest to be calculated. In the simplest case, a sketch is simply a random projection of the data, but it might be a more sophisticated object [100, §5.1]. The idea of sketching can be traced to the work of Alon et al. [5, 6].

**2.2.3. Numerical linear algebra.** Classically, the field of numerical linear algebra has focused on developing deterministic algorithms that produce highly accurate matrix approximations with provable guarantees. Nevertheless, randomized techniques have appeared in several environments.



One of the original examples is the use of random models for arithmetical errors, which was pioneered by von Neumann and Goldstine. Their papers [135, 136] stand among the first works to study the properties of random matrices. The earliest numerical linear algebra algorithm that depends essentially on randomized techniques is probably Dixon’s method for estimating norms and condition numbers [39].

Another situation where randomness commonly arises is the initialization of iterative methods for computing invariant subspaces. For example, most numerical linear algebra texts advocate random selection of the starting vector for the power method because it ensures that the vector has a nonzero component in the direction of a dominant eigenvector. Woźniakowski and coauthors have analyzed the performance of the power method and the Lanczos iteration given a random starting vector [79, 85].

Among other interesting applications of randomness, we mention the work by Parker and Pierce, which applies a randomized FFT to eliminate pivoting in Gaussian elimination [106], work by Demmel et al. who have studied randomization in connection with the stability of fast methods for linear algebra [35], and work by Le and Parker utilizing randomized methods for stabilizing fast linear algebraic computations based on recursive algorithms, such as Strassen’s matrix multiplication [81].

**2.2.4. Scientific computing.** One of the first algorithmic applications of randomness is the method of Monte Carlo integration introduced by Von Neumann and Ulam [95], and its extensions, such as the Metropolis algorithm for simulations in statistical physics. (See [9] for an introduction.) The most basic technique is to estimate an integral by sampling  $m$  points from the measure and computing an empirical mean of the integrand evaluated at the sample locations:

$$\int f(x) d\mu(x) \approx \frac{1}{m} \sum_{i=1}^m f(X_i),$$

where  $X_i$  are independent and identically distributed according to the probability measure  $\mu$ . The law of large numbers (usually) ensures that this approach produces the correct result in the limit as  $m \rightarrow \infty$ . Unfortunately, the approximation error typically has a standard deviation of  $m^{-1/2}$ , and the method provides no certificate of success.

The disappointing computational profile of Monte Carlo integration seems to have inspired a distaste for randomized approaches within the scientific computing community. Fortunately, there are many other types of randomized algorithms—such as the ones in this paper—that do not suffer from the same shortcomings.

**2.2.5. Geometric functional analysis.** There is one more character that plays a central role in our story: the probabilistic method in geometric analysis. Many of the algorithms and proof techniques ultimately come from work in this beautiful but recondite corner of mathematics.

Dvoretzky’s theorem [53] states (roughly) that every infinite-dimensional Banach space contains an  $n$ -dimensional subspace whose geometry is essentially the same as an  $n$ -dimensional Hilbert space, where  $n$  is an arbitrary natural number. In 1971, V. D. Milman developed a striking proof of this result by showing that a *random*  $n$ -dimensional subspace of an  $N$ -dimensional Banach space has this property with exceedingly high probability, provided that  $N$  is large enough [96]. Milman’s article debuted the *concentration of measure phenomenon*, which is a geometric interpretation of the classical idea that regular functions of independent random variables rarely

deviate far from their mean. This work opened a new era in geometric analysis where the probabilistic method became a basic instrument.

Another prominent example of measure concentration is Kashin's computation of the Gel'fand widths of the  $\ell_1$  ball [77], subsequently refined in [59]. This work showed that a *random*  $(N-n)$ -dimensional projection of the  $N$ -dimensional  $\ell_1$  ball has an astonishingly small Euclidean diameter: approximately  $\sqrt{(1 + \log(N/n))/n}$ . In contrast, a nonzero projection of the  $\ell_2$  ball always has Euclidean diameter one. This basic geometric fact undergirds recent developments in compressive sampling [21].

We have already described a third class of examples: the randomized embeddings of Johnson–Lindenstrauss [74] and of Bourgain [15].

Finally, we mention Maurey's technique of empirical approximation. The original work was unpublished; one of the earliest applications appears in [24, §1]. Although Maurey's idea has not received as much press as the examples above, it can lead to simple and efficient algorithms for sparse approximation. For some examples in machine learning, consider [8, 84, 110, 119].

The importance of random constructions in the geometric analysis community has led to the development of powerful techniques for studying random matrices. Classical random matrix theory focuses on a detailed asymptotic analysis of the spectral properties of special classes of random matrices. In contrast, geometric analysts know methods for determining the approximate behavior of rather complicated finite-dimensional random matrices. See [34] for a fairly current survey article. We also mention the works of Rudelson [115] and Rudelson–Vershynin [116], which describe powerful tools for studying random matrices drawn from certain discrete distributions. Their papers are rooted deeply in the field of geometric functional analysis, but they reach out toward computational applications.

**3. Linear algebraic preliminaries.** This section summarizes the background we need for the detailed description of randomized algorithms in §§4–6 and the analysis in §§8–11. We introduce notation in §3.1, describe some standard matrix decompositions in §3.2, and briefly review standard techniques for computing matrix factorizations in §3.3.

**3.1. Basic definitions.** The standard Hermitian geometry for  $\mathbb{C}^n$  is induced by the inner product

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x} \cdot \mathbf{y} = \sum_j x_j \overline{y_j}.$$

The associated norm is

$$\|\mathbf{x}\|^2 = \langle \mathbf{x}, \mathbf{x} \rangle = \sum_j |x_j|^2.$$

We usually measure the magnitude of a matrix  $\mathbf{A}$  with the operator norm

$$\|\mathbf{A}\| = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|},$$

which is often referred to as the *spectral norm*. The Frobenius norm is given by

$$\|\mathbf{A}\|_F = \left[ \sum_{jk} |a_{jk}|^2 \right]^{1/2}.$$

The conjugate transpose, or *adjoint*, of a matrix  $\mathbf{A}$  is denoted  $\mathbf{A}^*$ . The important identities

$$\|\mathbf{A}\|^2 = \|\mathbf{A}^* \mathbf{A}\| = \|\mathbf{A} \mathbf{A}^*\|$$

hold for each matrix  $\mathbf{A}$ .

We say that a matrix  $\mathbf{U}$  is *orthonormal* if its columns form an orthonormal set with respect to the Hermitian inner product. An orthonormal matrix  $\mathbf{U}$  preserves geometry in the sense that  $\|\mathbf{U}\mathbf{x}\| = \|\mathbf{x}\|$  for every vector  $\mathbf{x}$ . A *unitary* matrix is a square orthonormal matrix, and an *orthogonal* matrix is a real unitary matrix. Unitary matrices satisfy the relations  $\mathbf{U}\mathbf{U}^* = \mathbf{U}^*\mathbf{U} = \mathbf{I}$ . Both the operator norm and the Frobenius norm are *unitarily invariant*, which means that

$$\|\mathbf{U}\mathbf{A}\mathbf{V}^*\| = \|\mathbf{A}\| \quad \text{and} \quad \|\mathbf{U}\mathbf{A}\mathbf{V}^*\|_{\text{F}} = \|\mathbf{A}\|_{\text{F}}$$

for every matrix  $\mathbf{A}$  and all orthonormal matrices  $\mathbf{U}$  and  $\mathbf{V}$ .

We use the notation of [61] to denote submatrices. If  $\mathbf{A}$  is a matrix with entries  $a_{ij}$ , and if  $I = [i_1, i_2, \dots, i_p]$  and  $J = [j_1, j_2, \dots, j_q]$  are two index vectors, then the associated  $p \times q$  submatrix is expressed as

$$\mathbf{A}_{(I,J)} = \begin{bmatrix} a_{i_1,j_1} & \cdots & a_{i_1,j_q} \\ \vdots & & \vdots \\ a_{i_p,j_1} & \cdots & a_{i_p,j_q} \end{bmatrix}.$$

For column- and row-submatrices, we use the standard abbreviations

$$\mathbf{A}_{(:,J)} = \mathbf{A}_{([1, 2, \dots, m], J)}, \quad \text{and} \quad \mathbf{A}_{(I,:)} = \mathbf{A}_{(I, [1, 2, \dots, n])}.$$

**3.2. Standard matrix factorizations.** This section defines three basic matrix decompositions. Methods for computing them are described in §3.3.

**3.2.1. The pivoted QR factorization.** Each  $m \times n$  matrix  $\mathbf{A}$  of rank  $k$  admits a decomposition

$$\mathbf{A} = \mathbf{Q}\mathbf{R},$$

where  $\mathbf{Q}$  is an  $m \times k$  orthonormal matrix, and  $\mathbf{R}$  is a  $k \times n$  *weakly upper-triangular* matrix. That is, there exists a permutation  $J$  of the numbers  $\{1, 2, \dots, n\}$  such that  $\mathbf{R}_{(:,J)}$  is upper triangular. Moreover, the diagonal entries of  $\mathbf{R}_{(:,J)}$  are weakly decreasing. See [61, §5.4.1] for details.

**3.2.2. The singular value decomposition (SVD).** Each  $m \times n$  matrix  $\mathbf{A}$  of rank  $k$  admits a factorization

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*,$$

where  $\mathbf{U}$  is an  $m \times k$  orthonormal matrix,  $\mathbf{V}$  is an  $n \times k$  orthonormal matrix, and  $\mathbf{\Sigma}$  is a  $k \times k$  nonnegative, diagonal matrix

$$\mathbf{\Sigma} = \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_k \end{bmatrix}.$$

The numbers  $\sigma_j$  are called the *singular values* of  $\mathbf{A}$ . They are arranged in weakly decreasing order:

$$\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_k \geq 0.$$

The columns of  $\mathbf{U}$  and  $\mathbf{V}$  are called *left singular vectors* and *right singular vectors*, respectively.

Singular values are connected with the approximability of matrices. For each  $j$ , the number  $\sigma_{j+1}$  equals the spectral-norm discrepancy between  $\mathbf{A}$  and an optimal rank- $j$  approximation [97]. That is,

$$\sigma_{j+1} = \min\{\|\mathbf{A} - \mathbf{B}\| : \mathbf{B} \text{ has rank } j\}. \quad (3.1)$$

In particular,  $\sigma_1 = \|\mathbf{A}\|$ . See [61, §2.5.3 and §5.4.5] for additional details.

**3.2.3. The interpolative decomposition (ID).** Our final factorization identifies a collection of  $k$  columns from a rank- $k$  matrix  $\mathbf{A}$  that span the range of  $\mathbf{A}$ . To be precise, we can compute an index set  $J = [j_1, \dots, j_k]$  such that

$$\mathbf{A} = \mathbf{A}_{(:, J)} \mathbf{X},$$

where  $\mathbf{X}$  is a  $k \times n$  matrix that satisfies  $\mathbf{X}_{(:, J)} = \mathbf{I}_k$ . Furthermore, no entry of  $\mathbf{X}$  has magnitude larger than two. In other words, this decomposition expresses each column of  $\mathbf{A}$  using a linear combination of  $k$  fixed columns with *bounded* coefficients. Stable and efficient algorithms for computing the ID appear in the papers [26, 68].

It is also possible to compute a two-sided ID

$$\mathbf{A} = \mathbf{W} \mathbf{A}_{(J', J)} \mathbf{X},$$

where  $J'$  is an index set identifying  $k$  of the rows of  $\mathbf{A}$ , and  $\mathbf{W}$  is an  $m \times k$  matrix that satisfies  $\mathbf{W}_{(J', :)} = \mathbf{I}_k$  and whose entries are all bounded by two.

REMARK 3.1. There always exists an ID where the entries in the factor  $\mathbf{X}$  have magnitude bounded by one. Known proofs of this fact are constructive, e.g., [103, Lem. 3.3], but they require us to find a collection of  $k$  columns that has “maximum volume.” It is NP-hard to identify a subset of columns with this type of extremal property [27]. We find it remarkable that ID computations are possible as soon as the bound on  $\mathbf{X}$  is relaxed.

**3.3. Techniques for computing standard factorizations.** This section discusses some established deterministic techniques for computing the factorizations presented in §3.2. The material on pivoted QR and SVD can be located in any major text on numerical linear algebra, such as [61, 132]. References for the ID include [26, 68].

**3.3.1. Computing the full decomposition.** It is possible to compute the full QR factorization or the full SVD of an  $m \times n$  matrix to double-precision accuracy with  $O(mn \min\{m, n\})$  flops. Techniques for computing the SVD are iterative by necessity, but they converge so fast that we can treat them as finite for practical purposes.

**3.3.2. Computing partial decompositions.** Suppose that an  $m \times n$  matrix has numerical rank  $k$ , where  $k$  is substantially smaller than  $m$  and  $n$ . In this case, it is possible to produce a structured low-rank decomposition that approximates the matrix well. Sections 4 and 5 describe a set of randomized techniques for obtaining these partial decompositions. This section briefly reviews the classical techniques, which also play a role in developing randomized methods.

To compute a partial QR decomposition, the classical device is the Businger–Golub algorithm, which performs successive orthogonalization with pivoting on the columns of the matrix. The procedure halts when the Frobenius norm of the remaining

columns is less than a computational tolerance  $\varepsilon$ . Letting  $\ell$  denote the number of steps required, the process results in a partial factorization

$$\mathbf{A} = \mathbf{Q}\mathbf{R} + \mathbf{E}, \quad (3.2)$$

where  $\mathbf{Q}$  is an  $m \times \ell$  orthonormal matrix,  $\mathbf{R}$  is a  $\ell \times n$  weakly upper-triangular matrix, and  $\mathbf{E}$  is a residual that satisfies  $\|\mathbf{E}\|_F \leq \varepsilon$ . The computational cost is  $O(\ell mn)$ , and the number  $\ell$  of steps taken is typically close to the minimal rank  $k$  for which precision  $\varepsilon$  (in the Frobenius norm) is achievable. The Businger–Golub algorithm can in principle significantly overpredict the rank, but in practice this problem is very rare provided that orthonormality is maintained scrupulously.

Subsequent research has led to strong rank-revealing QR algorithms that succeed for all matrices. For example, the Gu–Eisenstat algorithm [68] (setting their parameter  $f = 2$ ) produces an QR decomposition of the form (3.2), where

$$\|\mathbf{E}\| \leq \sqrt{1 + 4k(n - k)} \cdot \sigma_{k+1}.$$

Recall that  $\sigma_{k+1}$  is the minimal error possible in a rank- $k$  approximation [97]. The cost of the Gu–Eisenstat algorithm is typically  $O(kmn)$ , but it can be slightly higher in rare cases. The algorithm can also be used to obtain an approximate ID [26].

To compute an approximate SVD of a general  $m \times n$  matrix, the most straightforward technique is to compute the full SVD and truncate it. This procedure is stable and accurate, but it requires  $O(mn \min\{m, n\})$  flops. A more efficient approach is to compute a partial QR factorization and postprocess the factors to obtain a partial SVD using the methods described below in §3.3.3. This scheme takes only  $O(kmn)$  flops. Krylov subspace methods can also compute partial SVDs at a comparable cost of  $O(kmn)$ , but they are less robust.

Note that all the techniques described in this section require extensive random access to the matrix, and they can be very slow when the matrix is stored out-of-core.

**3.3.3. Converting from one partial factorization to another.** Suppose that we have obtained a partial decomposition of a matrix  $\mathbf{A}$  by some means:

$$\|\mathbf{A} - \mathbf{C}\mathbf{B}\| \leq \varepsilon,$$

where  $\mathbf{B}$  and  $\mathbf{C}$  have rank  $k$ . Given this information, we can efficiently compute any of the basic factorizations.

We construct a partial QR factorization using the following three steps:

1. Compute a QR factorization of  $\mathbf{C}$  so that  $\mathbf{C} = \mathbf{Q}_1\mathbf{R}_1$ .
2. Form the product  $\mathbf{D} = \mathbf{R}_1\mathbf{B}$ , and compute a QR factorization:  $\mathbf{D} = \mathbf{Q}_2\mathbf{R}$ .
3. Form the product  $\mathbf{Q} = \mathbf{Q}_1\mathbf{Q}_2$ .

The result is an orthonormal matrix  $\mathbf{Q}$  and a weakly upper-triangular matrix  $\mathbf{R}$  such that  $\|\mathbf{A} - \mathbf{Q}\mathbf{R}\| \leq \varepsilon$ .

An analogous technique yields a partial SVD:

1. Compute a QR factorization of  $\mathbf{C}$  so that  $\mathbf{C} = \mathbf{Q}_1\mathbf{R}_1$ .
2. Form the product  $\mathbf{D} = \mathbf{R}_1\mathbf{B}$ , and compute an SVD:  $\mathbf{D} = \mathbf{U}_2\mathbf{\Sigma}\mathbf{V}^*$ .
3. Form the product  $\mathbf{U} = \mathbf{Q}_1\mathbf{U}_2$ .

The result is a diagonal matrix  $\mathbf{\Sigma}$  and orthonormal matrices  $\mathbf{U}$  and  $\mathbf{V}$  such that  $\|\mathbf{A} - \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*\| \leq \varepsilon$ .

Converting  $\mathbf{B}$  and  $\mathbf{C}$  into a partial ID is a one-step process:

1. Compute  $J$  and  $\mathbf{X}$  such that  $\mathbf{B} = \mathbf{B}_{(:,J)}\mathbf{X}$ .

Then  $\mathbf{A} \approx \mathbf{A}_{(:,J)}\mathbf{X}$ , but the approximation error may deteriorate from the initial estimate. For example, if we compute the ID using the Gu–Eisenstat algorithm [68] with the parameter  $f = 2$ , then the error  $\|\mathbf{A} - \mathbf{A}_{(:,J)}\mathbf{X}\| \leq (1 + \sqrt{1 + 4k(n-k)}) \cdot \varepsilon$ . Compare this bound with Lemma 5.1 below.

**3.3.4. Krylov-subspace methods.** Suppose that the matrix  $\mathbf{A}$  can be applied rapidly to vectors, as happens when  $\mathbf{A}$  is sparse or structured. Then Krylov subspace techniques can very effectively and accurately compute partial spectral decompositions. For concreteness, assume that  $\mathbf{A}$  is Hermitian. The idea of these techniques is to fix a starting vector  $\boldsymbol{\omega}$  and to seek approximations to the eigenvectors within the corresponding *Krylov subspace*

$$\mathcal{V}_q(\boldsymbol{\omega}) = \text{span}\{\boldsymbol{\omega}, \mathbf{A}\boldsymbol{\omega}, \mathbf{A}^2\boldsymbol{\omega}, \dots, \mathbf{A}^{q-1}\boldsymbol{\omega}\}.$$

Krylov methods also come in blocked versions, in which the starting *vector*  $\boldsymbol{\omega}$  is replaced by a starting *matrix*  $\boldsymbol{\Omega}$ . A common recommendation is to draw a starting vector  $\boldsymbol{\omega}$  (or starting matrix  $\boldsymbol{\Omega}$ ) from a standardized Gaussian distribution, which indicates a significant overlap between Krylov methods and the methods in this paper.

The most basic versions of Krylov methods for computing spectral decompositions are numerically unstable. High-quality implementations require that we incorporate restarting strategies, techniques for maintaining high-quality bases for the Krylov subspaces, etc. The diversity and complexity of such methods make it hard to state a precise computational cost, but in the environment we consider in this paper, a typical cost for a fully stable implementation would be

$$T_{\text{Krylov}} \sim k T_{\text{mult}} + k^2(m+n), \quad (3.3)$$

where  $T_{\text{mult}}$  is the cost of a matrix–vector multiplication.

## Part II: Algorithms

This part of the paper, §§4–7, provides detailed descriptions of randomized algorithms for constructing low-rank approximations to matrices. As discussed in §1.2, we split the problem into two stages. In Stage A, we construct a subspace that captures the action of the input matrix. In Stage B, we use this subspace to obtain an approximate factorization of the matrix.

Section 4 develops randomized methods for completing Stage A, and §5 describes deterministic methods for Stage B. Section 6 compares the computational costs of the resulting two-stage algorithm with the classical approaches outlined in §3. Finally, §7 illustrates the performance of the randomized schemes via numerical examples.

**4. Stage A: Randomized schemes for approximating the range.** This section outlines techniques for constructing a subspace that captures most of the action of a matrix. We begin with a recapitulation of the proto-algorithm that we introduced in §1.3. We discuss how it can be implemented in practice (§4.1) and then consider the question of how many random samples to acquire (§4.2). Afterward, we

## ALGORITHM 4.1: RANDOMIZED RANGE FINDER

*Given an  $m \times n$  matrix  $\mathbf{A}$ , and an integer  $\ell$ , this scheme computes an  $m \times \ell$  orthonormal matrix  $\mathbf{Q}$  whose range approximates the range of  $\mathbf{A}$ .*

- 1 Draw an  $n \times \ell$  Gaussian random matrix  $\mathbf{\Omega}$ .
- 2 Form the  $m \times \ell$  matrix  $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$ .
- 3 Construct an  $m \times \ell$  matrix  $\mathbf{Q}$  whose columns form an orthonormal basis for the range of  $\mathbf{Y}$ , e.g., using the QR factorization  $\mathbf{Y} = \mathbf{Q}\mathbf{R}$ .

present several ways in which the basic scheme can be improved. Sections 4.3 and 4.4 explain how to address the situation where the numerical rank of the input matrix is not known in advance. Section 4.5 shows how to modify the scheme to improve its accuracy when the singular spectrum of the input matrix decays slowly. Finally, §4.6 describes how the scheme can be accelerated by using a structured random matrix.

**4.1. The proto-algorithm revisited.** The most natural way to implement the proto-algorithm from §1.3 is to draw a random test matrix  $\mathbf{\Omega}$  from the standard Gaussian distribution. That is, each entry of  $\mathbf{\Omega}$  is an independent Gaussian random variable with mean zero and variance one. For reference, we formulate the resulting scheme as Algorithm 4.1.

The number  $T_{\text{basic}}$  of flops required by Algorithm 4.1 satisfies

$$T_{\text{basic}} \sim \ell n T_{\text{rand}} + \ell T_{\text{mult}} + \ell^2 m \quad (4.1)$$

where  $T_{\text{rand}}$  is the cost of generating a Gaussian random number and  $T_{\text{mult}}$  is the cost of multiplying  $\mathbf{A}$  by a vector. The three terms in (4.1) correspond directly with the three steps of Algorithm 4.1.

Empirically, we have found that the performance of Algorithm 4.1 depends very little on the quality of the random number generator used in Step 1.

The actual cost of Step 2 depends substantially on the matrix  $\mathbf{A}$  and the computational environment that we are working in. The estimate (4.1) suggests that Algorithm 4.1 is especially efficient when the matrix–vector product  $\mathbf{x} \mapsto \mathbf{A}\mathbf{x}$  can be evaluated rapidly. In particular, the scheme is appropriate for approximating sparse or structured matrices. Turn to §6 for more details.

The most important implementation issue arises when performing the basis calculation in Step 3. Typically, the columns of the sample matrix  $\mathbf{Y}$  are almost linearly dependent, so it is imperative to use stable methods for performing the orthonormalization. We have found that the Gram–Schmidt procedure, augmented with the *double orthogonalization* described in [12], is both convenient and reliable. Methods based on Householder reflectors or Givens rotations also work very well. Note that very little is gained by pivoting because the columns of the random matrix  $\mathbf{Y}$  are independent samples drawn from the same distribution.

**4.2. The number of samples required.** The goal of Algorithm 4.1 is to produce an orthonormal matrix  $\mathbf{Q}$  with few columns that achieves

$$\|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)\mathbf{A}\| \leq \varepsilon, \quad (4.2)$$

where  $\varepsilon$  is a specified tolerance. The number of columns  $\ell$  that the algorithm needs to reach this threshold is usually slightly larger than the minimal rank  $k$  of the smallest



basis that verifies (4.2). We refer to this discrepancy  $p = \ell - k$  as the *oversampling parameter*. The size of the oversampling parameter depends on several factors:

**The matrix dimensions.** Very large matrices may require more oversampling.

**The singular spectrum.** The more rapid the decay of the singular values, the less oversampling is needed. In the extreme case that the matrix has exact rank  $k$ , it is not necessary to oversample.

**The random test matrix.** Gaussian matrices succeed with very little oversampling, but are not always the most cost-effective option. The structured random matrices discussed in §4.6 may require substantial oversampling, but they still yield computational gains in certain settings.

The theoretical results in Part III provide detailed information about how the behavior of randomized schemes depends on these factors. For the moment, we limit ourselves to some general remarks on implementation issues.

For Gaussian test matrices, it is adequate to choose the oversampling parameter to be a small constant, such as  $p = 5$  or  $p = 10$ . There is rarely any advantage to select  $p > k$ . This observation, first presented in [91], demonstrates that a Gaussian test matrix results in a negligible amount of extra computation.

In practice, the target rank  $k$  is rarely known in advance. Randomized algorithms are usually implemented in an adaptive fashion where the number of samples is increased until the error satisfies the desired tolerance. In other words, the user never *chooses* the oversampling parameter. Theoretical results that bound the amount of oversampling are valuable primarily as aids for designing algorithms. We develop an adaptive approach in §§4.3–4.4.

The computational bottleneck in Algorithm 4.1 is usually the formation of the product  $\mathbf{A}\mathbf{\Omega}$ . As a result, it often pays to draw a larger number  $\ell$  of samples than necessary because the user can minimize the cost of the matrix multiplication with tools such as blocking of operations, high-level linear algebra subroutines, parallel processors, etc. This approach may lead to an ill-conditioned sample matrix  $\mathbf{Y}$ , but the orthogonalization in Step 3 of Algorithm 4.1 can easily identify the numerical rank of the sample matrix and ignore the excess samples. Furthermore, Stage B of the matrix approximation process succeeds even when the basis matrix  $\mathbf{Q}$  has a larger dimension than necessary.

**4.3. A posteriori error estimation.** Algorithm 4.1 is designed for solving the fixed-rank problem, where the target rank of the input matrix is specified in advance. To handle the fixed-precision problem, where the parameter is the computational tolerance, we need a scheme for estimating how well a putative basis matrix  $\mathbf{Q}$  captures the action of the matrix  $\mathbf{A}$ . To do so, we develop a probabilistic error estimator. These methods are inspired by work of Dixon [39]; our treatment follows [88, 137].

The exact approximation error is  $\|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)\mathbf{A}\|$ . It is intuitively plausible that we can obtain some information about this quantity by computing  $\|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)\mathbf{A}\boldsymbol{\omega}\|$ , where  $\boldsymbol{\omega}$  is a standard Gaussian vector. This notion leads to the following method. Draw a sequence  $\{\boldsymbol{\omega}^{(i)} : i = 1, 2, \dots, r\}$  of standard Gaussian vectors, where  $r$  is a small integer that balances computational cost and reliability. Then

$$\|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)\mathbf{A}\| \leq 10\sqrt{\frac{2}{\pi}} \max_{i=1, \dots, r} \|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)\mathbf{A}\boldsymbol{\omega}^{(i)}\| \quad (4.3)$$

with probability at least  $1 - 10^{-r}$ . This statement follows by setting  $\mathbf{B} = (\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)\mathbf{A}$  and  $\alpha = 10$  in the following lemma, whose proof appears in [137, §3.4].

LEMMA 4.1. *Let  $\mathbf{B}$  be a real  $m \times n$  matrix. Fix a positive integer  $r$  and a real number  $\alpha > 1$ . Draw an independent family  $\{\boldsymbol{\omega}^{(i)} : i = 1, 2, \dots, r\}$  of standard Gaussian vectors. Then*

$$\|\mathbf{B}\| \leq \alpha \sqrt{\frac{2}{\pi}} \max_{i=1, \dots, r} \|\mathbf{B}\boldsymbol{\omega}^{(i)}\|$$

except with probability  $\alpha^{-r}$ .

The critical point is that the error estimate (4.3) is computationally inexpensive because it requires only a small number of matrix–vector products. Therefore, we can make a lowball guess for the numerical rank of  $\mathbf{A}$  and add more samples if the error estimate is too large. The asymptotic cost of Algorithm 4.1 is preserved if we double our guess for the rank at each step. For example, we can start with 32 samples, compute another 32, then another 64, etc.

REMARK 4.1. The estimate (4.3) is actually somewhat crude. We can obtain a better estimate at a similar computational cost by initializing a power iteration with a random vector and repeating the process several times [88].

**4.4. Error estimation (almost) for free.** The error estimate described in §4.3 can be combined with any method for constructing an approximate basis for the range of a matrix. In this section, we explain how the error estimator can be incorporated into Algorithm 4.1 at almost no additional cost.

To be precise, let us suppose that  $\mathbf{A}$  is an  $m \times n$  matrix and  $\varepsilon$  is a computational tolerance. We seek an integer  $\ell$  and an  $m \times \ell$  orthonormal matrix  $\mathbf{Q}^{(\ell)}$  such that

$$\|(\mathbf{I} - \mathbf{Q}^{(\ell)}(\mathbf{Q}^{(\ell)})^*)\mathbf{A}\| \leq \varepsilon. \quad (4.4)$$

The size  $\ell$  of the basis will typically be slightly larger than the size  $k$  of the smallest basis that achieves this error.

The basic observation behind the adaptive scheme is that we can generate the basis in Step 3 of Algorithm 4.1 incrementally. Starting with an empty basis matrix  $\mathbf{Q}^{(0)}$ , the following scheme generates an orthonormal matrix whose range captures the action of  $\mathbf{A}$ :

```

for  $i = 1, 2, 3, \dots$ 
    Draw an  $n \times 1$  Gaussian random vector  $\boldsymbol{\omega}^{(i)}$ , and set  $\mathbf{y}^{(i)} = \mathbf{A}\boldsymbol{\omega}^{(i)}$ .
    Compute  $\tilde{\mathbf{q}}^{(i)} = (\mathbf{I} - \mathbf{Q}^{(i-1)}(\mathbf{Q}^{(i-1)})^*)\mathbf{y}^{(i)}$ .
    Normalize  $\mathbf{q}^{(i)} = \tilde{\mathbf{q}}^{(i)} / \|\tilde{\mathbf{q}}^{(i)}\|$ , and form  $\mathbf{Q}^{(i)} = [\mathbf{Q}^{(i-1)} \ \mathbf{q}^{(i)}]$ .
end for
```

How do we know when we have reached a basis  $\mathbf{Q}^{(\ell)}$  that verifies (4.4)? The answer becomes apparent once we observe that the vectors  $\tilde{\mathbf{q}}^{(i)}$  are precisely the vectors that appear in the error bound (4.3). The resulting rule is that we break the loop once we observe  $r$  consecutive vectors  $\tilde{\mathbf{q}}^{(i)}$  whose norms are smaller than  $\varepsilon / (10\sqrt{2/\pi})$ .

A formal description of the resulting algorithm appears as Algorithm 4.2. A potential complication of the method is that the vectors  $\tilde{\mathbf{q}}^{(i)}$  become small as the basis starts to capture most of the action of  $\mathbf{A}$ . In finite-precision arithmetic, their direction is extremely unreliable. To address this problem, we simply reproject the normalized vector  $\mathbf{q}^{(i)}$  onto  $\text{range}(\mathbf{Q}^{(i-1)})^\perp$  in steps 7 and 8 of Algorithm 4.2.

The CPU time requirements of Algorithms 4.2 and 4.1 are essentially identical. Although Algorithm 4.2 computes the last few samples purely to obtain the error

estimate, this apparent extra cost is offset by the fact that Algorithm 4.1 always includes an oversampling factor. The failure probability stated for Algorithm 4.2 is pessimistic because it is derived from a simple union bound argument. In practice, the error estimator is reliable in a range of circumstances when we take  $r = 10$ .

ALGORITHM 4.2: ADAPTIVE RANDOMIZED RANGE FINDER

*Given an  $m \times n$  matrix  $\mathbf{A}$ , a tolerance  $\varepsilon$ , and an integer  $r$  (e.g.  $r = 10$ ), the following scheme computes an orthonormal matrix  $\mathbf{Q}$  such that (4.2) holds with probability at least  $1 - \min\{m, n\}10^{-r}$ .*

```

1  Draw standard Gaussian vectors  $\omega^{(1)}, \dots, \omega^{(r)}$  of length  $n$ .
2  For  $i = 1, 2, \dots, r$ , compute  $\mathbf{y}^{(i)} = \mathbf{A}\omega^{(i)}$ .
3   $j = 0$ .
4   $\mathbf{Q}^{(0)} = [\ ]$ , the  $m \times 0$  empty matrix.
5  while  $\max\{\|\mathbf{y}^{(j+1)}\|, \|\mathbf{y}^{(j+2)}\|, \dots, \|\mathbf{y}^{(j+r)}\|\} > \varepsilon/(10\sqrt{2/\pi})$ ,
6       $j = j + 1$ .
7      Overwrite  $\mathbf{y}^{(j)}$  by  $(\mathbf{I} - \mathbf{Q}^{(j-1)}(\mathbf{Q}^{(j-1)})^*)\mathbf{y}^{(j)}$ .
8       $\mathbf{q}^{(j)} = \mathbf{y}^{(j)} / \|\mathbf{y}^{(j)}\|$ .
9       $\mathbf{Q}^{(j)} = [\mathbf{Q}^{(j-1)} \ \mathbf{q}^{(j)}]$ .
10     Draw a standard Gaussian vector  $\omega^{(j+r)}$  of length  $n$ .
11      $\mathbf{y}^{(j+r)} = (\mathbf{I} - \mathbf{Q}^{(j)}(\mathbf{Q}^{(j)})^*)\mathbf{A}\omega^{(j+r)}$ .
12     for  $i = (j+1), (j+2), \dots, (j+r-1)$ ,
13         Overwrite  $\mathbf{y}^{(i)}$  by  $\mathbf{y}^{(i)} - \mathbf{q}^{(j)}\langle \mathbf{q}^{(j)}, \mathbf{y}^{(i)} \rangle$ .
14     end for
15 end while
16  $\mathbf{Q} = \mathbf{Q}^{(j)}$ .
```

REMARK 4.2. The calculations in Algorithm 4.2 can be organized so that each iteration processes a block of samples simultaneously. This revision can lead to dramatic improvements in speed because it allows us to exploit higher-level linear algebra subroutines (e.g., BLAS3) or parallel processors. Although blocking can lead to the generation of unnecessary samples, this outcome is generally harmless, as noted in §4.2.

#### 4.5. A modified scheme for matrices whose singular values decay slowly.

The techniques described in §4.1 and §4.4 work well for matrices whose singular values exhibit some decay, but they may produce a poor basis when the input matrix has a flat singular spectrum or when the input matrix is very large. In this section, we describe techniques, originally proposed in [67, 112], for improving the accuracy of randomized algorithms in these situations. Related earlier work includes [114] and the literature on classical orthogonal iteration methods [61, p. 332].

The intuition behind these techniques is that the singular vectors associated with small singular values interfere with the calculation, so we reduce their weight relative to the dominant singular vectors by taking powers of the matrix to be analyzed. More precisely, we wish to apply the randomized sampling scheme to the matrix  $\mathbf{B} = (\mathbf{A}\mathbf{A}^*)^q\mathbf{A}$ , where  $q$  is a small integer. The matrix  $\mathbf{B}$  has the same singular

## ALGORITHM 4.3: RANDOMIZED POWER ITERATION

*Given an  $m \times n$  matrix  $\mathbf{A}$  and integers  $\ell$  and  $q$ , this algorithm computes an  $m \times \ell$  orthonormal matrix  $\mathbf{Q}$  whose range approximates the range of  $\mathbf{A}$ .*

- 1 Draw an  $n \times \ell$  Gaussian random matrix  $\mathbf{\Omega}$ .
- 2 Form the  $m \times \ell$  matrix  $\mathbf{Y} = (\mathbf{A}\mathbf{A}^*)^q \mathbf{A}\mathbf{\Omega}$  via alternating application of  $\mathbf{A}$  and  $\mathbf{A}^*$ .
- 3 Construct an  $m \times \ell$  matrix  $\mathbf{Q}$  whose columns form an orthonormal basis for the range of  $\mathbf{Y}$ , e.g., via the QR factorization  $\mathbf{Y} = \mathbf{Q}\mathbf{R}$ .

**Note:** This procedure is vulnerable to round-off errors; see Remark 4.3. The recommended implementation appears as Algorithm 4.4.

vectors as the input matrix  $\mathbf{A}$ , but its singular values decay much more quickly:

$$\sigma_j(\mathbf{B}) = \sigma_j(\mathbf{A})^{2q+1}, \quad j = 1, 2, 3, \dots \quad (4.5)$$

We modify Algorithm 4.1 by replacing the formula  $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$  in Step 2 by the formula  $\mathbf{Y} = \mathbf{B}\mathbf{\Omega} = (\mathbf{A}\mathbf{A}^*)^q \mathbf{A}\mathbf{\Omega}$ , and we obtain Algorithm 4.3.

Algorithm 4.3 requires  $2q + 1$  times as many matrix–vector multiplies as Algorithm 4.1, but is far more accurate in situations where the singular values of  $\mathbf{A}$  decay slowly. A good heuristic is that when the original scheme produces a basis whose approximation error is within a factor  $C$  of the optimum, the power scheme produces an approximation error within  $C^{1/(2q+1)}$  of the optimum. In other words, the power iteration drives the approximation gap to one exponentially fast. See Theorem 9.2 and §10.4 for the details.

Algorithm 4.3 targets the fixed-rank problem. To address the fixed-precision problem, we can incorporate the error estimators described in §4.3 to obtain an adaptive scheme analogous with Algorithm 4.2. In situations where it is critical to achieve near-optimal approximation errors, one can increase the oversampling beyond our standard recommendation  $\ell = k + 5$  all the way to  $\ell = 2k$  without changing the scaling of the asymptotic computational cost. A supporting analysis appears in Corollary 10.10.

**REMARK 4.3.** Unfortunately, when Algorithm 4.3 is executed in floating point arithmetic, rounding errors will extinguish all information pertaining to singular modes associated with singular values that are small compared with  $\|\mathbf{A}\|$ . (Roughly, if machine precision is  $\mu$ , then all information associated with singular values smaller than  $\mu^{1/(2q+1)} \|\mathbf{A}\|$  is lost.) This problem can easily be remedied by orthonormalizing the columns of the sample matrix between each application of  $\mathbf{A}$  and  $\mathbf{A}^*$ . The resulting scheme, summarized as Algorithm 4.4, is algebraically equivalent to Algorithm 4.3 when executed in exact arithmetic [92, 124]. We recommend Algorithm 4.4 because its computational costs are similar to Algorithm 4.3, even though the former is substantially more accurate in floating-point arithmetic.

**4.6. An accelerated technique for general dense matrices.** This section describes a set of techniques that allow us to compute an approximate rank- $\ell$  factorization of a general dense  $m \times n$  matrix in roughly  $O(mn \log(\ell))$  flops, in contrast to the asymptotic cost  $O(mn\ell)$  required by earlier methods. We can tailor this scheme for the real or complex case, but we focus on the conceptually simpler complex case. These algorithms were introduced in [137]; similar techniques were proposed in [118].

## ALGORITHM 4.4: RANDOMIZED SUBSPACE ITERATION

Given an  $m \times n$  matrix  $\mathbf{A}$  and integers  $\ell$  and  $q$ , this algorithm computes an  $m \times \ell$  orthonormal matrix  $\mathbf{Q}$  whose range approximates the range of  $\mathbf{A}$ .

```

1  Draw an  $n \times \ell$  standard Gaussian matrix  $\mathbf{\Omega}$ .
2  Form  $\mathbf{Y}_0 = \mathbf{A}\mathbf{\Omega}$  and compute its QR factorization  $\mathbf{Y}_0 = \mathbf{Q}_0\mathbf{R}_0$ .
3  for  $j = 1, 2, \dots, q$ 
4      Form  $\tilde{\mathbf{Y}}_j = \mathbf{A}^*\mathbf{Q}_{j-1}$  and compute its QR factorization  $\tilde{\mathbf{Y}}_j = \tilde{\mathbf{Q}}_j\tilde{\mathbf{R}}_j$ .
5      Form  $\mathbf{Y}_j = \mathbf{A}\tilde{\mathbf{Q}}_j$  and compute its QR factorization  $\mathbf{Y}_j = \mathbf{Q}_j\mathbf{R}_j$ .
6  end
7   $\mathbf{Q} = \mathbf{Q}_q$ .
```

The first step toward this accelerated technique is to observe that the bottleneck in Algorithm 4.1 is the computation of the matrix product  $\mathbf{A}\mathbf{\Omega}$ . When the test matrix  $\mathbf{\Omega}$  is standard Gaussian, the cost of this multiplication is  $O(mn\ell)$ , the same as a rank-revealing QR algorithm [68]. The key idea is to use a *structured* random matrix that allows us to compute the product in  $O(mn \log(\ell))$  flops.

The *subsampled random Fourier transform*, or SRFT, is perhaps the simplest example of a structured random matrix that meets our goals. An SRFT is an  $n \times \ell$  matrix of the form

$$\mathbf{\Omega} = \sqrt{\frac{n}{\ell}} \mathbf{D}\mathbf{F}\mathbf{R}, \quad (4.6)$$

where

- $\mathbf{D}$  is an  $n \times n$  diagonal matrix whose entries are independent random variables uniformly distributed on the complex unit circle,
- $\mathbf{F}$  is the  $n \times n$  unitary discrete Fourier transform (DFT), whose entries take the values  $f_{pq} = n^{-1/2} e^{-2\pi i(p-1)(q-1)/n}$  for  $p, q = 1, 2, \dots, n$ , and
- $\mathbf{R}$  is an  $n \times \ell$  matrix that samples  $\ell$  coordinates from  $n$  uniformly at random, i.e., its  $\ell$  columns are drawn randomly without replacement from the columns of the  $n \times n$  identity matrix.

When  $\mathbf{\Omega}$  is defined by (4.6), we can compute the sample matrix  $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$  using  $O(mn \log(\ell))$  flops via a subsampled FFT [137]. Then we form the basis  $\mathbf{Q}$  by orthonormalizing the columns of  $\mathbf{Y}$ , as described in §4.1. This scheme appears as Algorithm 4.5. The total number  $T_{\text{struct}}$  of flops required by this procedure is

$$T_{\text{struct}} \sim mn \log(\ell) + \ell^2 n \quad (4.7)$$

Note that if  $\ell$  is substantially larger than the numerical rank  $k$  of the input matrix, we can perform the orthogonalization with  $O(k\ell n)$  flops because the columns of the sample matrix are almost linearly dependent.

The test matrix (4.6) is just one choice among many possibilities. Other suggestions that appear in the literature include subsampled Hadamard transforms, chains of Givens rotations acting on randomly chosen coordinates, and many more. See [86] and its bibliography. Empirically, we have found that the transform summarized in Remark 4.6 below performs very well in a variety of environments [113].

## ALGORITHM 4.5: FAST RANDOMIZED RANGE FINDER

*Given an  $m \times n$  matrix  $\mathbf{A}$ , and an integer  $\ell$ , this scheme computes an  $m \times \ell$  orthonormal matrix  $\mathbf{Q}$  whose range approximates the range of  $\mathbf{A}$ .*

- 1 Draw an  $n \times \ell$  SRFT test matrix  $\mathbf{\Omega}$ , as defined by (4.6).
- 2 Form the  $m \times \ell$  matrix  $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$  using a (subsampled) FFT.
- 3 Construct an  $m \times \ell$  matrix  $\mathbf{Q}$  whose columns form an orthonormal basis for the range of  $\mathbf{Y}$ , e.g., using the QR factorization  $\mathbf{Y} = \mathbf{Q}\mathbf{R}$ .

At this point, it is not well understood how to quantify and compare the behavior of structured random transforms. One reason for this uncertainty is that it has been difficult to analyze the amount of oversampling that various transforms require. Section 11 establishes that the random matrix (4.6) can be used to identify a near-optimal basis for a rank- $k$  matrix using  $\ell \sim (k + \log(n)) \log(k)$  samples. In practice, the transforms (4.6) and (4.8) typically require no more oversampling than a Gaussian test matrix requires. (For a numerical example, see §7.4.) As a consequence, setting  $\ell = k + 10$  or  $\ell = k + 20$  is typically more than adequate. Further research on these questions would be valuable.

REMARK 4.4. The structured random matrices discussed in this section do not adapt readily to the fixed-precision problem, where the computational tolerance is specified, because the samples from the range are usually computed in bulk. Fortunately, these schemes are sufficiently inexpensive that we can progressively increase the number of samples computed starting with  $\ell = 32$ , say, and then proceeding to  $\ell = 64, 128, 256, \dots$  until we achieve the desired tolerance.

REMARK 4.5. When using the SRFT (4.6) for matrix approximation, we have a choice whether to use a subsampled FFT or a full FFT. The complete FFT is so inexpensive that it often pays to construct an extended sample matrix  $\mathbf{Y}_{\text{large}} = \mathbf{A}\mathbf{D}\mathbf{F}$  and then generate the actual samples by drawing columns at random from  $\mathbf{Y}_{\text{large}}$  and rescaling as needed. The asymptotic cost increases to  $O(mn \log(n))$  flops, but the full FFT is actually faster for moderate problem sizes because the constant suppressed by the big-O notation is so small. Adaptive rank determination is easy because we just examine extra samples as needed.

REMARK 4.6. Among the structured random matrices that we have tried, one of the strongest candidates involves sequences of random Givens rotations [113]. This matrix takes the form

$$\mathbf{\Omega} = \mathbf{D}'' \mathbf{\Theta}' \mathbf{D}' \mathbf{\Theta} \mathbf{D} \mathbf{F} \mathbf{R}, \quad (4.8)$$

where the prime symbol  $'$  indicates an independent realization of a random matrix. The matrices  $\mathbf{R}$ ,  $\mathbf{F}$ , and  $\mathbf{D}$  are defined after (4.6). The matrix  $\mathbf{\Theta}$  is a chain of random Givens rotations:

$$\mathbf{\Theta} = \mathbf{\Pi} \mathbf{G}(1, 2; \theta_1) \mathbf{G}(2, 3; \theta_2) \cdots \mathbf{G}(n-1, n; \theta_{n-1})$$

where  $\mathbf{\Pi}$  is a random  $n \times n$  permutation matrix; where  $\theta_1, \dots, \theta_{n-1}$  are independent random variables uniformly distributed on the interval  $[0, 2\pi]$ ; and where  $\mathbf{G}(i, j; \theta)$  denotes a rotation on  $\mathbb{C}^n$  by the angle  $\theta$  in the  $(i, j)$  coordinate plane [61, §5.1.8].

## ALGORITHM 5.1: DIRECT SVD

*Given matrices  $\mathbf{A}$  and  $\mathbf{Q}$  such that (5.1) holds, this procedure computes an approximate factorization  $\mathbf{A} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ , where  $\mathbf{U}$  and  $\mathbf{V}$  are orthonormal, and  $\mathbf{\Sigma}$  is a nonnegative diagonal matrix.*

- 1 Form the matrix  $\mathbf{B} = \mathbf{Q}^* \mathbf{A}$ .
- 2 Compute an SVD of the small matrix:  $\mathbf{B} = \tilde{\mathbf{U}}\mathbf{\Sigma}\mathbf{V}^*$ .
- 3 Form the orthonormal matrix  $\mathbf{U} = \mathbf{Q}\tilde{\mathbf{U}}$ .

REMARK 4.7. When the singular values of the input matrix  $\mathbf{A}$  decay slowly, Algorithm 4.5 may perform poorly in terms of accuracy. When randomized sampling is used with a Gaussian random matrix, the recourse is to take a couple of steps of a power iteration; see Algorithm 4.4. However, it is not currently known whether such an iterative scheme can be accelerated to  $O(mn \log(k))$  complexity using “fast” random transforms such as the SRFT.

**5. Stage B: Construction of standard factorizations.** The algorithms for Stage A described in §4 produce an orthonormal matrix  $\mathbf{Q}$  whose range captures the action of an input matrix  $\mathbf{A}$ :

$$\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^* \mathbf{A}\| \leq \varepsilon, \quad (5.1)$$

where  $\varepsilon$  is a computational tolerance. This section describes methods for approximating standard factorizations of  $\mathbf{A}$  using the information in the basis  $\mathbf{Q}$ .

To accomplish this task, we pursue the idea from §3.3.3 that any low-rank factorization  $\mathbf{A} \approx \mathbf{C}\mathbf{B}$  can be manipulated to produce a standard decomposition. When the bound (5.1) holds, the low-rank factors are simply  $\mathbf{C} = \mathbf{Q}$  and  $\mathbf{B} = \mathbf{Q}^* \mathbf{A}$ . The simplest scheme (§5.1) computes the factor  $\mathbf{B}$  directly with a matrix–matrix product to ensure a minimal error in the final approximation. An alternative approach (§5.2) constructs factors  $\mathbf{B}$  and  $\mathbf{C}$  without forming any matrix–matrix product. The approach of §5.2 is often faster than the approach of §5.1 but typically results in larger errors. Both schemes can be streamlined for an Hermitian input matrix (§5.3) and a positive semidefinite input matrix (§5.4). Finally, we develop single-pass algorithms that exploit other information generated in Stage A to avoid revisiting the input matrix (§5.5).

Throughout this section,  $\mathbf{A}$  denotes an  $m \times n$  matrix, and  $\mathbf{Q}$  is an  $m \times k$  orthonormal matrix that verifies (5.1). For purposes of exposition, we concentrate on methods for constructing the partial SVD.

**5.1. Factorizations based on forming  $\mathbf{Q}^* \mathbf{A}$  directly.** The relation (5.1) implies that  $\|\mathbf{A} - \mathbf{Q}\mathbf{B}\| \leq \varepsilon$ , where  $\mathbf{B} = \mathbf{Q}^* \mathbf{A}$ . Once we have computed  $\mathbf{B}$ , we can produce any standard factorization using the methods of §3.3.3. Algorithm 5.1 illustrates how to build an approximate SVD.

The factors produced by Algorithm 5.1 satisfy

$$\|\mathbf{A} - \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*\| \leq \varepsilon. \quad (5.2)$$

In other words, the approximation error does not degrade.

The cost of Algorithm 5.1 is generally dominated by the cost of the product  $\mathbf{Q}^* \mathbf{A}$  in Step 1, which takes  $O(kmn)$  flops for a general dense matrix. Note that this scheme



is particularly well suited to environments where we have a fast method for computing the matrix–vector product  $\mathbf{x} \mapsto \mathbf{A}^* \mathbf{x}$ , for example when  $\mathbf{A}$  is sparse or structured. This approach retains a strong advantage over Krylov-subspace methods and rank-revealing QR because Step 1 can be accelerated using BLAS3, parallel processors, and so forth. Steps 2 and 3 require  $O(k^2 n)$  and  $O(k^2 m)$  flops respectively.

**REMARK 5.1.** Algorithm 5.1 produces an approximate SVD with the same rank as the basis matrix  $\mathbf{Q}$ . When the size of the basis exceeds the desired rank  $k$  of the SVD, it may be preferable to retain only the dominant  $k$  singular values and singular vectors. Equivalently, we replace the diagonal matrix  $\mathbf{\Sigma}$  of computed singular values with the matrix  $\mathbf{\Sigma}_{(k)}$  formed by zeroing out all but the largest  $k$  entries of  $\mathbf{\Sigma}$ . In the worst case, this truncation step can increase the approximation error by  $\sigma_{k+1}$ ; see §9.4 for an analysis. Our numerical experience suggests that this error analysis is pessimistic, and the term  $\sigma_{k+1}$  often does not appear in practice.

**5.2. Postprocessing via row extraction.** Given a matrix  $\mathbf{Q}$  such that (5.1) holds, we can obtain a rank- $k$  factorization

$$\mathbf{A} \approx \mathbf{X} \mathbf{B}, \quad (5.3)$$

where  $\mathbf{B}$  is a  $k \times n$  matrix consisting of  $k$  rows extracted from  $\mathbf{A}$ . The approximation (5.3) can be produced without computing any matrix–matrix products, which makes this approach to postprocessing very fast. The drawback comes because the error  $\|\mathbf{A} - \mathbf{X} \mathbf{B}\|$  is usually larger than the initial error  $\|\mathbf{A} - \mathbf{Q} \mathbf{Q}^* \mathbf{A}\|$ , especially when the dimensions of  $\mathbf{A}$  are large. See Remark 5.3 for more discussion.

To obtain the factorization (5.3), we simply construct the interpolative decomposition (§3.2.3) of the matrix  $\mathbf{Q}$ :

$$\mathbf{Q} = \mathbf{X} \mathbf{Q}_{(J, :)}. \quad (5.4)$$

The index set  $J$  marks  $k$  rows of  $\mathbf{Q}$  that span the row space of  $\mathbf{Q}$ , and  $\mathbf{X}$  is an  $m \times k$  matrix whose entries are bounded in magnitude by two and contains the  $k \times k$  identity as a submatrix:  $\mathbf{X}_{(J, :)} = \mathbf{I}_k$ . Combining (5.4) and (5.1), we reach

$$\mathbf{A} \approx \mathbf{Q} \mathbf{Q}^* \mathbf{A} = \mathbf{X} \mathbf{Q}_{(J, :)} \mathbf{Q}^* \mathbf{A}. \quad (5.5)$$

Since  $\mathbf{X}_{(J, :)} = \mathbf{I}_k$ , equation (5.5) implies that  $\mathbf{A}_{(J, :)} \approx \mathbf{Q}_{(J, :)} \mathbf{Q}^* \mathbf{A}$ . Therefore, (5.3) follows when we put  $\mathbf{B} = \mathbf{A}_{(J, :)}$ .

Provided with the factorization (5.3), we can obtain any standard factorization using the techniques of §3.3.3. Algorithm 5.2 illustrates an SVD calculation. This procedure requires  $O(k^2(m+n))$  flops. The following lemma guarantees the accuracy of the computed factors.

**LEMMA 5.1.** *Let  $\mathbf{A}$  be an  $m \times n$  matrix and let  $\mathbf{Q}$  be an  $m \times k$  matrix that satisfy (5.1). Suppose that  $\mathbf{U}$ ,  $\mathbf{\Sigma}$ , and  $\mathbf{V}$  are the matrices constructed by Algorithm 5.2. Then*

$$\|\mathbf{A} - \mathbf{U} \mathbf{\Sigma} \mathbf{V}^*\| \leq \left[1 + \sqrt{1 + 4k(n-k)}\right] \varepsilon. \quad (5.6)$$

*Proof.* The factors  $\mathbf{U}$ ,  $\mathbf{\Sigma}$ ,  $\mathbf{V}$  constructed by the algorithm satisfy

$$\mathbf{U} \mathbf{\Sigma} \mathbf{V}^* = \mathbf{U} \mathbf{\Sigma} \tilde{\mathbf{V}}^* \mathbf{W}^* = \mathbf{Z} \mathbf{W}^* = \mathbf{X} \mathbf{R}^* \mathbf{W}^* = \mathbf{X} \mathbf{A}_{(J, :)}.$$

## ALGORITHM 5.2: SVD VIA ROW EXTRACTION

Given matrices  $\mathbf{A}$  and  $\mathbf{Q}$  such that (5.1) holds, this procedure computes an approximate factorization  $\mathbf{A} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ , where  $\mathbf{U}$  and  $\mathbf{V}$  are orthonormal, and  $\mathbf{\Sigma}$  is a nonnegative diagonal matrix.

- 1 Compute an ID  $\mathbf{Q} = \mathbf{X}\mathbf{Q}_{(J,:)}.$  (The ID is defined in §3.2.3.)
- 2 Extract  $\mathbf{A}_{(J,:)}$ , and compute a QR factorization  $\mathbf{A}_{(J,:)} = \mathbf{R}^*\mathbf{W}^*.$
- 3 Form the product  $\mathbf{Z} = \mathbf{X}\mathbf{R}^*.$
- 4 Compute an SVD  $\mathbf{Z} = \mathbf{U}\mathbf{\Sigma}\tilde{\mathbf{V}}^*.$
- 5 Form the orthonormal matrix  $\mathbf{V} = \mathbf{W}\tilde{\mathbf{V}}.$

**Note:** Algorithm 5.2 is faster than Algorithm 5.1 but less accurate.

**Note:** It is advantageous to replace the basis  $\mathbf{Q}$  by the sample matrix  $\mathbf{Y}$  produced in Stage A, cf. Remark 5.2.

Define the approximation

$$\hat{\mathbf{A}} = \mathbf{Q}\mathbf{Q}^*\mathbf{A}. \quad (5.7)$$

Since  $\hat{\mathbf{A}} = \mathbf{X}\mathbf{Q}_{(J,:)}\mathbf{Q}^*\mathbf{A}$  and since  $\mathbf{X}_{(J,:)} = \mathbf{I}_k$ , it must be that  $\hat{\mathbf{A}}_{(J,:)} = \mathbf{Q}_{(J,:)}\mathbf{Q}^*\mathbf{A}.$  Consequently,

$$\hat{\mathbf{A}} = \mathbf{X}\hat{\mathbf{A}}_{(J,:)}.$$

We have the chain of relations

$$\begin{aligned} \|\mathbf{A} - \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*\| &= \|\mathbf{A} - \mathbf{X}\mathbf{A}_{(J,:)}\| \\ &= \|(\mathbf{A} - \mathbf{X}\hat{\mathbf{A}}_{(J,:)}) + (\mathbf{X}\hat{\mathbf{A}}_{(J,:)} - \mathbf{X}\mathbf{A}_{(J,:)})\| \\ &\leq \|\mathbf{A} - \hat{\mathbf{A}}\| + \|\mathbf{X}\hat{\mathbf{A}}_{(J,:)} - \mathbf{X}\mathbf{A}_{(J,:)}\| \\ &\leq \|\mathbf{A} - \hat{\mathbf{A}}\| + \|\mathbf{X}\| \|\mathbf{A}_{(J,:)} - \hat{\mathbf{A}}_{(J,:)}\|. \end{aligned} \quad (5.8)$$

Inequality (5.1) ensures that  $\|\mathbf{A} - \hat{\mathbf{A}}\| \leq \varepsilon.$  Since  $\mathbf{A}_{(J,:)} - \hat{\mathbf{A}}_{(J,:)}$  is a submatrix of  $\mathbf{A} - \hat{\mathbf{A}},$  we must also have  $\|\mathbf{A}_{(J,:)} - \hat{\mathbf{A}}_{(J,:)}\| \leq \varepsilon.$  Thus, (5.8) reduces to

$$\|\mathbf{A} - \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*\| \leq (1 + \|\mathbf{X}\|)\varepsilon. \quad (5.9)$$

The bound (5.6) follows from (5.9) after we observe that  $\mathbf{X}$  contains a  $k \times k$  identity matrix and that the entries of the remaining  $(n - k) \times k$  submatrix are bounded in magnitude by two.  $\square$

**REMARK 5.2.** To maintain a unified presentation, we have formulated all the postprocessing techniques so they take an orthonormal matrix  $\mathbf{Q}$  as input. Recall that, in Stage A of our framework, we construct the matrix  $\mathbf{Q}$  by orthonormalizing the columns of the sample matrix  $\mathbf{Y}.$  With finite-precision arithmetic, it is preferable to adapt Algorithm 5.2 to start directly from the sample matrix  $\mathbf{Y}.$  To be precise, we modify Step 1 to compute  $\mathbf{X}$  and  $J$  so that  $\mathbf{Y} = \mathbf{X}\mathbf{Y}_{(J,:)}.$  This revision is recommended even when  $\mathbf{Q}$  is available from the adaptive rank determination of Algorithm 4.2.

**ALGORITHM 5.3: DIRECT EIGENVALUE DECOMPOSITION**

*Given an Hermitian matrix  $\mathbf{A}$  and a basis  $\mathbf{Q}$  such that (5.1) holds, this procedure computes an approximate eigenvalue decomposition  $\mathbf{A} \approx \mathbf{U}\mathbf{\Lambda}\mathbf{U}^*$ , where  $\mathbf{U}$  is orthonormal, and  $\mathbf{\Lambda}$  is a real diagonal matrix.*

- 1 Form the small matrix  $\mathbf{B} = \mathbf{Q}^* \mathbf{A} \mathbf{Q}$ .
- 2 Compute an eigenvalue decomposition  $\mathbf{B} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^*$ .
- 3 Form the orthonormal matrix  $\mathbf{U} = \mathbf{Q} \mathbf{V}$ .

**REMARK 5.3.** As the inequality (5.6) suggests, the factorization produced by Algorithm 5.2 is potentially less accurate than the basis that it uses as input. This loss of accuracy is problematic when  $\varepsilon$  is not so small or when  $kn$  is large. In such cases, we recommend Algorithm 5.1 over Algorithm 5.2; the former is more costly, but it does not amplify the error, as shown in (5.2).

**5.3. Postprocessing an Hermitian matrix.** When  $\mathbf{A}$  is Hermitian, the postprocessing becomes particularly elegant. In this case, the columns of  $\mathbf{Q}$  form a good basis for both the column space *and* the row space of  $\mathbf{A}$  so that we have  $\mathbf{A} \approx \mathbf{Q} \mathbf{Q}^* \mathbf{A} \mathbf{Q} \mathbf{Q}^*$ . More precisely, when (5.1) is in force, we have

$$\begin{aligned} \|\mathbf{A} - \mathbf{Q} \mathbf{Q}^* \mathbf{A} \mathbf{Q} \mathbf{Q}^*\| &= \|\mathbf{A} - \mathbf{Q} \mathbf{Q}^* \mathbf{A} + \mathbf{Q} \mathbf{Q}^* \mathbf{A} - \mathbf{Q} \mathbf{Q}^* \mathbf{A} \mathbf{Q} \mathbf{Q}^*\| \\ &\leq \|\mathbf{A} - \mathbf{Q} \mathbf{Q}^* \mathbf{A}\| + \|\mathbf{Q} \mathbf{Q}^* (\mathbf{A} - \mathbf{A} \mathbf{Q} \mathbf{Q}^*)\| \leq 2\varepsilon. \end{aligned} \quad (5.10)$$

The last inequality relies on the facts that  $\|\mathbf{Q} \mathbf{Q}^*\| = 1$  and that

$$\|\mathbf{A} - \mathbf{A} \mathbf{Q} \mathbf{Q}^*\| = \|(\mathbf{A} - \mathbf{A} \mathbf{Q} \mathbf{Q}^*)^*\| = \|\mathbf{A} - \mathbf{Q} \mathbf{Q}^* \mathbf{A}\|.$$

Since  $\mathbf{A} \approx \mathbf{Q}(\mathbf{Q}^* \mathbf{A} \mathbf{Q}) \mathbf{Q}^*$  is a low-rank approximation of  $\mathbf{A}$ , we can form any standard factorization using the techniques from §3.3.3.

For Hermitian  $\mathbf{A}$ , it is more common to compute an eigenvalue decomposition than an SVD. We can accomplish this goal using Algorithm 5.3, which adapts the scheme from §5.1. This procedure delivers a factorization that satisfies the error bound  $\|\mathbf{A} - \mathbf{U} \mathbf{\Lambda} \mathbf{U}^*\| \leq 2\varepsilon$ . The calculation requires  $O(kn^2)$  flops.

We can also pursue the row extraction approach from §5.2, which is faster but less accurate. See Algorithm 5.4 for the details. The total cost is  $O(k^2n)$  flops.

**5.4. Postprocessing a positive semidefinite matrix.** When the input matrix  $\mathbf{A}$  is positive semidefinite, the *Nyström method* can be used to improve the quality of standard factorizations at almost no additional cost; see [48] and its bibliography. To describe the main idea, we first recall that the direct method presented in §5.3 manipulates the approximate rank- $k$  factorization

$$\mathbf{A} \approx \mathbf{Q} (\mathbf{Q}^* \mathbf{A} \mathbf{Q}) \mathbf{Q}^*. \quad (5.11)$$

In contrast, the Nyström scheme builds a more sophisticated rank- $k$  approximation, namely

$$\begin{aligned} \mathbf{A} &\approx (\mathbf{A} \mathbf{Q}) (\mathbf{Q}^* \mathbf{A} \mathbf{Q})^{-1} (\mathbf{A} \mathbf{Q})^* \\ &= \left[ (\mathbf{A} \mathbf{Q}) (\mathbf{Q}^* \mathbf{A} \mathbf{Q})^{-1/2} \right] \left[ (\mathbf{A} \mathbf{Q}) (\mathbf{Q}^* \mathbf{A} \mathbf{Q})^{-1/2} \right]^* = \mathbf{F} \mathbf{F}^*, \end{aligned} \quad (5.12)$$

**ALGORITHM 5.4: EIGENVALUE DECOMPOSITION VIA ROW EXTRACTION**

*Given an Hermitian matrix  $\mathbf{A}$  and a basis  $\mathbf{Q}$  such that (5.1) holds, this procedure computes an approximate eigenvalue decomposition  $\mathbf{A} \approx \mathbf{U}\mathbf{\Lambda}\mathbf{U}^*$ , where  $\mathbf{U}$  is orthonormal, and  $\mathbf{\Lambda}$  is a real diagonal matrix.*

- 1 Compute an ID  $\mathbf{Q} = \mathbf{X}\mathbf{Q}_{(J,:)}.$
- 2 Perform a QR factorization  $\mathbf{X} = \mathbf{V}\mathbf{R}.$
- 3 Form the product  $\mathbf{Z} = \mathbf{R}\mathbf{A}_{(J,J)}\mathbf{R}^*.$
- 4 Compute an eigenvalue decomposition  $\mathbf{Z} = \mathbf{W}\mathbf{\Lambda}\mathbf{W}^*.$
- 5 Form the orthonormal matrix  $\mathbf{U} = \mathbf{V}\mathbf{W}.$

**Note:** Algorithm 5.4 is faster than Algorithm 5.3 but less accurate.

**Note:** It is advantageous to replace the basis  $\mathbf{Q}$  by the sample matrix  $\mathbf{Y}$  produced in Stage A, cf. Remark 5.2.

**ALGORITHM 5.5: EIGENVALUE DECOMPOSITION VIA NYSTRÖM METHOD**

*Given a positive semidefinite matrix  $\mathbf{A}$  and a basis  $\mathbf{Q}$  such that (5.1) holds, this procedure computes an approximate eigenvalue decomposition  $\mathbf{A} \approx \mathbf{U}\mathbf{\Lambda}\mathbf{U}^*$ , where  $\mathbf{U}$  is orthonormal, and  $\mathbf{\Lambda}$  is nonnegative and diagonal.*

- 1 Form the matrices  $\mathbf{B}_1 = \mathbf{A}\mathbf{Q}$  and  $\mathbf{B}_2 = \mathbf{Q}^*\mathbf{B}_1.$
- 2 Perform a Cholesky factorization  $\mathbf{B}_2 = \mathbf{C}^*\mathbf{C}.$
- 3 Form  $\mathbf{F} = \mathbf{B}_1\mathbf{C}^{-1}$  using a triangular solve.
- 4 Compute an SVD  $\mathbf{F} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$  and set  $\mathbf{\Lambda} = \mathbf{\Sigma}^2.$

where  $\mathbf{F}$  is an approximate Cholesky factor of  $\mathbf{A}$  with dimension  $n \times k$ . To compute the factor  $\mathbf{F}$  numerically, first form the matrices  $\mathbf{B}_1 = \mathbf{A}\mathbf{Q}$  and  $\mathbf{B}_2 = \mathbf{Q}^*\mathbf{B}_1$ . Then decompose the psd matrix  $\mathbf{B}_2 = \mathbf{C}^*\mathbf{C}$  into its Cholesky factors. Finally compute the factor  $\mathbf{F} = \mathbf{B}_1\mathbf{C}^{-1}$  by performing a triangular solve. The low-rank factorization (5.12) can be converted to a standard decomposition using the techniques from §3.3.3.

The literature contains an explicit expression [48, Lem. 4] for the approximation error in (5.12). This result implies that, in the spectral norm, the Nyström approximation error never exceeds  $\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^*\mathbf{A}\|$ , and it is often substantially smaller. We omit a detailed discussion.

For an example of the Nyström technique, consider Algorithm 5.5, which computes an approximate eigenvalue decomposition of a positive semidefinite matrix. This method should be compared with the scheme for Hermitian matrices, Algorithm 5.3. In both cases, the dominant cost occurs when we form  $\mathbf{A}\mathbf{Q}$ , so the two procedures have roughly the same running time. On the other hand, Algorithm 5.5 is typically much more accurate than Algorithm 5.3. In a sense, we are exploiting the fact that  $\mathbf{A}$  is positive semidefinite to take one step of subspace iteration (Algorithm 4.4) for free.

**5.5. Single-pass algorithms.** The techniques described in §§5.1–5.4 all require us to revisit the input matrix. This may not be feasible in environments where the matrix is too large to be stored. In this section, we develop a method that requires just one pass over the matrix to construct not only an approximate basis but also a complete factorization. Similar techniques appear in [137] and [29].

For motivation, we begin with the case where  $\mathbf{A}$  is Hermitian. Let us recall the proto-algorithm from §1.3.3: Draw a random test matrix  $\mathbf{\Omega}$ ; form the sample matrix  $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$ ; then construct a basis  $\mathbf{Q}$  for the range of  $\mathbf{Y}$ . It turns out that the matrices  $\mathbf{\Omega}$ ,  $\mathbf{Y}$ , and  $\mathbf{Q}$  contain all the information we need to approximate  $\mathbf{A}$ .

To see why, define the (currently unknown) matrix  $\mathbf{B}$  via  $\mathbf{B} = \mathbf{Q}^* \mathbf{A} \mathbf{Q}$ . Postmultiplying the definition by  $\mathbf{Q}^* \mathbf{\Omega}$ , we obtain the identity  $\mathbf{B} \mathbf{Q}^* \mathbf{\Omega} = \mathbf{Q}^* \mathbf{A} \mathbf{Q} \mathbf{Q}^* \mathbf{\Omega}$ . The relationships  $\mathbf{A} \mathbf{Q} \mathbf{Q}^* \approx \mathbf{A}$  and  $\mathbf{A} \mathbf{\Omega} = \mathbf{Y}$  show that  $\mathbf{B}$  must satisfy

$$\mathbf{B} \mathbf{Q}^* \mathbf{\Omega} \approx \mathbf{Q}^* \mathbf{Y}. \quad (5.13)$$

All three matrices  $\mathbf{\Omega}$ ,  $\mathbf{Y}$ , and  $\mathbf{Q}$  are available, so we can solve (5.13) to obtain the matrix  $\mathbf{B}$ . Then the low-rank factorization  $\mathbf{A} \approx \mathbf{Q} \mathbf{B} \mathbf{Q}^*$  can be converted to an eigenvalue decomposition via familiar techniques. The entire procedure requires  $O(k^2 n)$  flops, and it is summarized as Algorithm 5.6.

ALGORITHM 5.6: EIGENVALUE DECOMPOSITION IN ONE PASS

*Given an Hermitian matrix  $\mathbf{A}$ , a random test matrix  $\mathbf{\Omega}$ , a sample matrix  $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$ , and an orthonormal matrix  $\mathbf{Q}$  that verifies (5.1) and  $\mathbf{Y} = \mathbf{Q} \mathbf{Q}^* \mathbf{Y}$ , this algorithm computes an approximate eigenvalue decomposition  $\mathbf{A} \approx \mathbf{U} \mathbf{\Lambda} \mathbf{U}^*$ .*

- 1 Use a standard least-squares solver to find an Hermitian matrix  $\mathbf{B}_{\text{approx}}$  that approximately satisfies the equation  $\mathbf{B}_{\text{approx}}(\mathbf{Q}^* \mathbf{\Omega}) \approx \mathbf{Q}^* \mathbf{Y}$ .
- 2 Compute the eigenvalue decomposition  $\mathbf{B}_{\text{approx}} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^*$ .
- 3 Form the product  $\mathbf{U} = \mathbf{Q} \mathbf{V}$ .

When  $\mathbf{A}$  is not Hermitian, it is still possible to devise single-pass algorithms, but we must modify the initial Stage A of the approximation framework to simultaneously construct bases for the ranges of  $\mathbf{A}$  and  $\mathbf{A}^*$ :

1. Generate random matrices  $\mathbf{\Omega}$  and  $\tilde{\mathbf{\Omega}}$ .
2. Compute  $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$  and  $\tilde{\mathbf{Y}} = \mathbf{A}^* \tilde{\mathbf{\Omega}}$  in a single pass over  $\mathbf{A}$ .
3. Compute QR factorizations  $\mathbf{Y} = \mathbf{Q} \mathbf{R}$  and  $\tilde{\mathbf{Y}} = \tilde{\mathbf{Q}} \tilde{\mathbf{R}}$ .

This procedure results in matrices  $\mathbf{Q}$  and  $\tilde{\mathbf{Q}}$  such that  $\mathbf{A} \approx \mathbf{Q} \mathbf{Q}^* \mathbf{A} \tilde{\mathbf{Q}} \tilde{\mathbf{Q}}^*$ . The reduced matrix we must approximate is  $\mathbf{B} = \mathbf{Q}^* \mathbf{A} \tilde{\mathbf{Q}}$ . In analogy with (5.13), we find that

$$\mathbf{Q}^* \mathbf{Y} = \mathbf{Q}^* \mathbf{A} \mathbf{\Omega} \approx \mathbf{Q}^* \mathbf{A} \tilde{\mathbf{Q}} \tilde{\mathbf{Q}}^* \mathbf{\Omega} = \mathbf{B} \tilde{\mathbf{Q}}^* \mathbf{\Omega}. \quad (5.14)$$

An analogous calculation shows that  $\mathbf{B}$  should also satisfy

$$\tilde{\mathbf{Q}}^* \tilde{\mathbf{Y}} \approx \mathbf{B}^* \mathbf{Q}^* \tilde{\mathbf{\Omega}}. \quad (5.15)$$

Now, the reduced matrix  $\mathbf{B}_{\text{approx}}$  can be determined by finding a minimum-residual solution to the system of relations (5.14) and (5.15).

REMARK 5.4. The single-pass approaches described in this section can degrade the approximation error in the final decomposition significantly. To explain the issue, we focus on the Hermitian case. It turns out that the coefficient matrix  $\mathbf{Q}^* \mathbf{\Omega}$  in the linear system (5.13) is usually ill-conditioned. In a worst-case scenario, the error  $\|\mathbf{A} - \mathbf{U} \mathbf{\Lambda} \mathbf{U}^*\|$  in the factorization produced by Algorithm 5.6 could be larger than

the error resulting from the two-pass method of Section 5.3 by a factor of  $1/\tau_{\min}$ , where  $\tau_{\min}$  is the minimal singular value of the matrix  $\mathbf{Q}^*\mathbf{\Omega}$ .

The situation can be improved by oversampling. Suppose that we seek a rank- $k$  approximate eigenvalue decomposition. Pick a small oversampling parameter  $p$ . Draw an  $n \times (k+p)$  random matrix  $\mathbf{\Omega}$ , and form the sample matrix  $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$ . Let  $\mathbf{Q}$  denote the  $n \times k$  matrix formed by the  $k$  leading left singular vectors of  $\mathbf{Y}$ . Now, the linear system (5.13) has a coefficient matrix  $\mathbf{Q}^*\mathbf{\Omega}$  of size  $k \times (k+p)$ , so it is overdetermined. An approximate solution of this system yields a  $k \times k$  matrix  $\mathbf{B}$ .

**6. Computational costs.** So far, we have postponed a detailed discussion of the computational cost of randomized matrix approximation algorithms because it is necessary to account for both the first stage, where we compute an approximate basis for the range (§4), and the second stage, where we postprocess the basis to complete the factorization (§5). We are now prepared to compare the cost of the two-stage scheme with the cost of traditional techniques.

Choosing an appropriate algorithm, whether classical or randomized, requires us to consider the properties of the input matrix. To draw a nuanced picture, we discuss three representative computational environments in §6.1–6.3. We close with some comments on parallel implementations in §6.4.

For concreteness, we focus on the problem of computing an approximate SVD of an  $m \times n$  matrix  $\mathbf{A}$  with numerical rank  $k$ . The costs for other factorizations are similar.

**6.1. General matrices that fit in core memory.** Suppose that  $\mathbf{A}$  is a general matrix presented as an array of numbers that fits in core memory. In this case, the appropriate method for Stage A is to use a structured random matrix (§4.6), which allows us to find a basis that captures the action of the matrix using  $O(mn \log(k) + k^2m)$  flops. For Stage B, we apply the row-extraction technique (§5.2), which costs an additional  $O(k^2(m+n))$  flops. The total number of operations  $T_{\text{random}}$  for this approach satisfies

$$T_{\text{random}} \sim mn \log(k) + k^2(m+n).$$

As a rule of thumb, the approximation error of this procedure satisfies

$$\|\mathbf{A} - \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*\| \lesssim n \cdot \sigma_{k+1}, \quad (6.1)$$

where  $\sigma_{k+1}$  is the  $(k+1)$ th singular value of  $\mathbf{A}$ . The estimate (6.1), which follows from Theorem 11.2 and Lemma 5.1, reflects the worst-case scenario; actual errors are usually smaller.

This algorithm should be compared with modern deterministic techniques, such as rank-revealing QR followed by postprocessing (§3.3.2) which typically require

$$T_{\text{RRQR}} \sim kmn$$

operations to achieve a comparable error.

In this setting, the randomized algorithm can be several times faster than classical techniques even for problems of moderate size, say  $m, n \sim 10^3$  and  $k \sim 10^2$ . See §7.4 for numerical evidence.

**REMARK 6.1.** In case row extraction is impractical, there is an alternative  $O(mn \log(k))$  technique described in [137, §5.2]. When the error (6.1) is unacceptably large, we can use the direct method (§5.1) for Stage B, which brings the total cost to  $O(kmn)$  flops.

**6.2. Matrices for which matrix–vector products can be rapidly evaluated.** In many problems in data mining and scientific computing, the cost  $T_{\text{mult}}$  of performing the matrix–vector multiplication  $\mathbf{x} \mapsto \mathbf{A}\mathbf{x}$  is substantially smaller than the nominal cost  $O(mn)$  for the dense case. It is not uncommon that  $O(m+n)$  flops suffice. Standard examples include (i) very sparse matrices; (ii) structured matrices, such as Töplitz operators, that can be applied using the FFT or other means; and (iii) matrices that arise from physical problems, such as discretized integral operators, that can be applied via, e.g., the fast multipole method [66].

Suppose that both  $\mathbf{A}$  and  $\mathbf{A}^*$  admit fast multiplies. The appropriate randomized approach for this scenario completes Stage A using Algorithm 4.1 with  $p$  constant (for the fixed-rank problem) or Algorithm 4.2 (for the fixed-precision problem) at a cost of  $(k+p)T_{\text{mult}} + O(k^2m)$  flops. For Stage B, we invoke Algorithm 5.1, which requires  $(k+p)T_{\text{mult}} + O(k^2(m+n))$  flops. The total cost  $T_{\text{sparse}}$  satisfies

$$T_{\text{sparse}} = 2(k+p)T_{\text{mult}} + O(k^2(m+n)). \quad (6.2)$$

As a rule of thumb, the approximation error of this procedure satisfies

$$\|\mathbf{A} - \mathbf{U}\Sigma\mathbf{V}^*\| \lesssim \sqrt{kn} \cdot \sigma_{k+1}. \quad (6.3)$$

The estimate (6.3) follows from Corollary 10.9 and the discussion in §5.1. Actual errors are usually smaller.

When the singular spectrum of  $\mathbf{A}$  decays slowly, we can incorporate  $q$  iterations of the power method (Algorithm 4.3) to obtain superior solutions to the fixed-rank problem. The computational cost increases to, cf. (6.2),

$$T_{\text{sparse}} = (2q+2)(k+p)T_{\text{mult}} + O(k^2(m+n)), \quad (6.4)$$

while the error (6.3) improves to

$$\|\mathbf{A} - \mathbf{U}\Sigma\mathbf{V}^*\| \lesssim (kn)^{1/2(2q+1)} \cdot \sigma_{k+1}. \quad (6.5)$$

The estimate (6.5) takes into account the discussion in §10.4. The power scheme can also be adapted for the fixed-precision problem (§4.5).

In this setting, the classical prescription for obtaining a partial SVD is some variation of a Krylov-subspace method; see §3.3.4. These methods exhibit great diversity, so it is hard to specify a “typical” computational cost. To a first approximation, it is fair to say that in order to obtain an approximate SVD of rank  $k$ , the cost of a numerically stable implementation of a Krylov method is no less than the cost (6.2) with  $p$  set to zero. At this price, the Krylov method often obtains better accuracy than the basic randomized method obtained by combining Algorithms 4.1 and 5.1, especially for matrices whose singular values decay slowly. On the other hand, the randomized schemes are inherently more robust and allow much more freedom in organizing the computation to suit a particular application or a particular hardware architecture. The latter point is in practice of crucial importance because it is usually much faster to apply a matrix to  $k$  vectors simultaneously than it is to execute  $k$  matrix–vector multiplications consecutively. In practice, blocking and parallelism can lead to enough gain that a few steps of the power method (Algorithm 4.3) can be performed more quickly than  $k$  steps of a Krylov method.

**REMARK 6.2.** Any comparison between randomized sampling schemes and Krylov variants becomes complicated because of the fact that “basic” Krylov schemes such



as Lanczos [61, p. 473] or Arnoldi [61, p. 499] are inherently unstable. To obtain numerical robustness, we must incorporate sophisticated modifications such as restarts, reorthogonalization procedures, etc. Constructing a high-quality implementation is sufficiently hard that the authors of a popular book on “numerical recipes” qualify their treatment of spectral computations as follows [109, p. 567]:

You have probably gathered by now that the solution of eigensystems is a fairly complicated business. It is. It is one of the few subjects covered in this book for which we do *not* recommend that you avoid canned routines. On the contrary, the purpose of this chapter is precisely to give you some appreciation of what is going on inside such canned routines, so that you can make intelligent choices about using them, and intelligent diagnoses when something goes wrong.

Randomized sampling does not eliminate the difficulties referred to in this quotation; however it reduces the task of computing a *partial* spectral decomposition of a very large matrix to the task of computing a *full* decomposition of a small dense matrix. (For example, in Algorithm 5.1, the input matrix  $\mathbf{A}$  is large and  $\mathbf{B}$  is small.) The latter task is much better understood and is eminently suitable for using canned routines. Random sampling schemes interact with the large matrix only through matrix–matrix products, which can easily be implemented by a user in a manner appropriate to the application and to the available hardware.

The comparison is further complicated by the fact that there is significant overlap between the two sets of ideas. Algorithm 4.3 is conceptually similar to a “block Lanczos method” [61, p. 485] with a random starting matrix. Indeed, we believe that there are significant opportunities for cross-fertilization in this area. Hybrid schemes that combine the best ideas from both fields may perform very well.

**6.3. General matrices stored in slow memory or streamed.** The traditional metric for numerical algorithms is the number of floating-point operations they require. When the data does not fit in fast memory, however, the computational time is often dominated by the cost of memory access. In this setting, a more appropriate measure of algorithmic performance is *pass-efficiency*, which counts how many times the data needs to be cycled through fast memory. Flop counts become largely irrelevant.

All the classical matrix factorization techniques that we discuss in §3.2—including dense SVD, rank-revealing QR, Krylov methods, and so forth—require at least  $k$  passes over the the matrix, which is prohibitively expensive for huge data matrices. A desire to reduce the pass count of matrix approximation algorithms served as one of the early motivations for developing randomized schemes [46, 58, 105]. Detailed recent work appears in [29].

For many matrices, randomized techniques can produce an accurate approximation using just one pass over the data. For Hermitian matrices, we obtain a single-pass algorithm by combining Algorithm 4.1, which constructs an approximate basis, with Algorithm 5.6, which produces an eigenvalue decomposition without any additional access to the matrix. Section 5.5 describes the analogous technique for general matrices.

For the huge matrices that arise in applications such as data mining, it is common that the singular spectrum decays slowly. Relevant applications include image processing (see §§7.2–7.3 for numerical examples), statistical data analysis, and network monitoring. To compute approximate factorizations in these environments, it is

crucial to enhance the accuracy of the randomized approach using the power scheme, Algorithm 4.3, or some other device. This approach increases the pass count somewhat, but in our experience it is very rare that more than five passes are required.

**6.4. Gains from parallelization.** As mentioned in §§6.2–6.3, randomized methods often outperform classical techniques not because they involve fewer floating-point operations but rather because they allow us to reorganize the calculations to exploit the matrix properties and the computer architecture more fully. In addition, these methods are well suited for parallel implementation. For example, in Algorithm 4.1, the computational bottleneck is the evaluation of the matrix product  $\mathbf{A}\mathbf{\Omega}$ , which is embarrassingly parallelizable.

**7. Numerical examples.** By this time, the reader has surely formulated a pointed question: Do these randomized matrix approximation algorithms actually work in practice? In this section, we attempt to address this concern by illustrating how the algorithms perform on a diverse collection of test cases.

Section 7.1 starts with two examples from the physical sciences involving discrete approximations to operators with exponentially decaying spectra. Sections 7.2 and 7.3 continue with two examples of matrices arising in “data mining.” These are large matrices whose singular spectra decay slowly; one is sparse and fits in RAM, one is dense and is stored out-of-core. Finally, §7.4 investigates the performance of randomized methods based on structured random matrices.

Sections 7.1–7.3 focus on the algorithms for Stage A that we presented in §4 because we wish to isolate the performance of the randomized step.

Computational examples illustrating truly large data matrices have been reported elsewhere, for instance in [69].

**7.1. Two matrices with rapidly decaying singular values.** We first illustrate the behavior of the adaptive range approximation method, Algorithm 4.2. We apply it to two matrices associated with the numerical analysis of differential and integral operators. The matrices in question have rapidly decaying singular values and our intent is to demonstrate that in this environment, the approximation error of a bare-bones randomized method such as Algorithm 4.2 is *very* close to the minimal error achievable by any method. We observe that the approximation error of a randomized method is itself a random variable (it is a function of the random matrix  $\mathbf{\Omega}$ ) so what we need to demonstrate is not only that the error is small in a typical realization, but also that it clusters tightly around the mean value.

We first consider a  $200 \times 200$  matrix  $\mathbf{A}$  that results from discretizing the following single-layer operator associated with the Laplace equation:

$$[S\sigma](x) = \text{const} \cdot \int_{\Gamma_1} \log|x-y| \sigma(y) dA(y), \quad x \in \Gamma_2, \quad (7.1)$$

where  $\Gamma_1$  and  $\Gamma_2$  are the two contours in  $\mathbb{R}^2$  illustrated in Figure 7.1(a). We approximate the integral with the trapezoidal rule, which converges superalgebraically because the kernel is smooth. In the absence of floating-point errors, we estimate that the discretization error would be less than  $10^{-20}$  for a smooth source  $\sigma$ . The leading constant is selected so the matrix  $\mathbf{A}$  has unit operator norm.

We implement Algorithm 4.2 in Matlab v6.5. Gaussian test matrices are generated using the `randn` command. For each number  $\ell$  of samples, we compare the following three quantities:

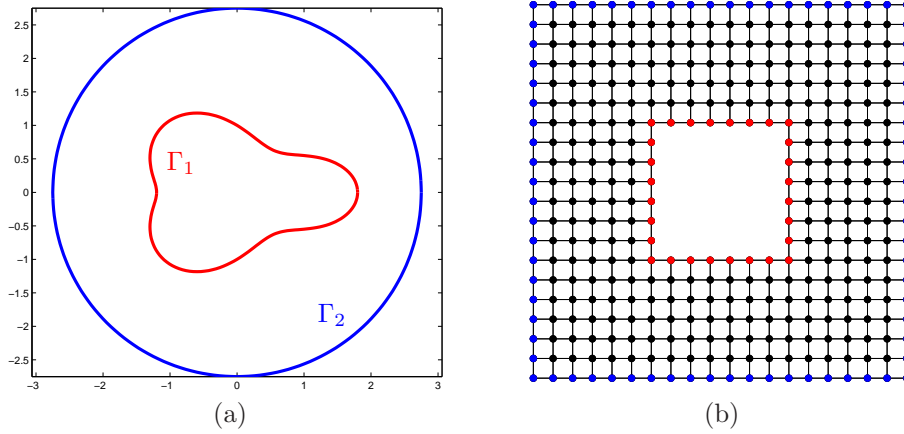


FIG. 7.1. Configurations for physical problems. (a) The contours  $\Gamma_1$  (red) and  $\Gamma_2$  (blue) for the integral operator (7.1). (b) Geometry of the lattice problem associated with matrix  $\mathbf{B}$  in §7.1.

1. The minimum rank- $\ell$  approximation error  $\sigma_{\ell+1}$  is determined using `svd`.
2. The actual error  $e_\ell = \|(\mathbf{I} - \mathbf{Q}^{(\ell)}(\mathbf{Q}^{(\ell)})^*)\mathbf{A}\|$  is computed with `norm`.
3. A random estimator  $f_\ell$  for the actual error  $e_\ell$  is obtained from (4.3), with the parameter  $r$  set to 5.

Note that any values less than  $10^{-15}$  should be considered numerical artifacts.

Figure 7.2 tracks a characteristic execution of Algorithm 4.2. We make three observations: (i) The error  $e_\ell$  incurred by the algorithm is remarkably close to the theoretical minimum  $\sigma_{\ell+1}$ . (ii) The error estimate always produces an upper bound for the actual error. Without the built-in  $10\times$  safety margin, the estimate would track the actual error almost exactly. (iii) The basis constructed by the algorithm essentially reaches full double-precision accuracy.

How typical is the trial documented in Figure 7.2? To answer this question, we examine the empirical performance of the algorithm over 2000 independent trials. Figure 7.3 charts the error estimate versus the actual error at four points during the course of execution:  $\ell = 25, 50, 75, 100$ . We offer four observations: (i) The initial run detailed in Figure 7.2 is entirely typical. (ii) Both the actual and estimated error concentrate about their mean value. (iii) The actual error drifts slowly away from the optimal error as the number  $\ell$  of samples increases. (iv) The error estimator is *always* pessimistic by a factor of about ten, which means that the algorithm *never* produces a basis with lower accuracy than requested. The only effect of selecting an unlucky sample matrix  $\mathbf{\Omega}$  is that the algorithm proceeds for a few additional steps.

We next consider a matrix  $\mathbf{B}$  which is defined implicitly in the sense that we cannot access its elements directly; we can only evaluate the map  $\mathbf{x} \mapsto \mathbf{B}\mathbf{x}$  for a given vector  $\mathbf{x}$ . To be precise,  $\mathbf{B}$  represents a transfer matrix for a network of resistors like the one shown in Figure 7.1(b). The vector  $\mathbf{x}$  represents a set of electric potentials specified on the red nodes in the figure. These potentials induce a unique equilibrium field on the network in which the potential of each black and blue node is the average of the potentials of its three or four neighbors. The vector  $\mathbf{B}\mathbf{x}$  is then the restriction of the potential to the blue exterior nodes. Given a vector  $\mathbf{x}$ , the vector  $\mathbf{B}\mathbf{x}$  can be obtained by solving a large sparse linear system whose coefficient matrix is the classical five-point stencil approximating the 2D Laplace operator.

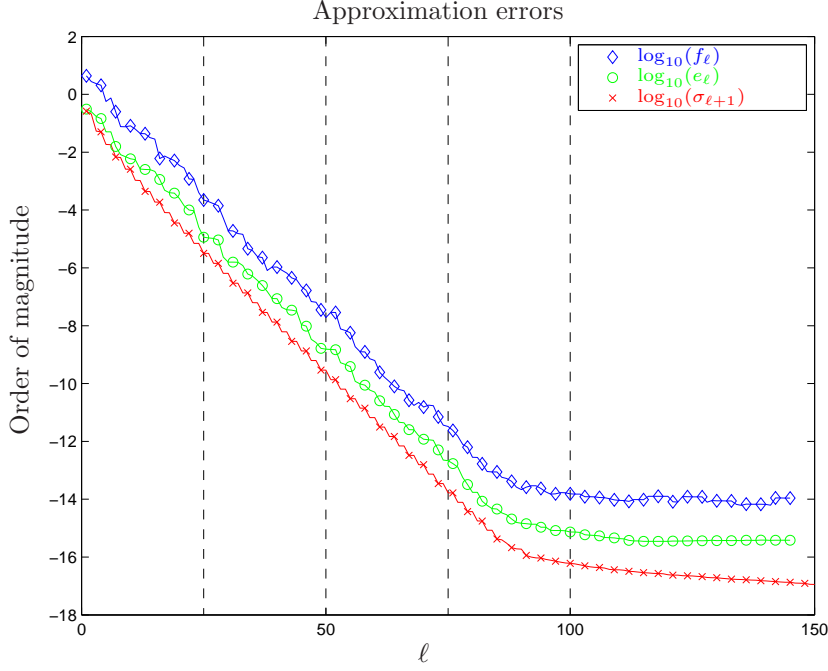


FIG. 7.2. Approximating a Laplace integral operator. One execution of Algorithm 4.2 for the  $200 \times 200$  input matrix  $\mathbf{A}$  described in §7.1. The number  $\ell$  of random samples varies along the horizontal axis; the vertical axis measures the base-10 logarithm of error magnitudes. The dashed vertical lines mark the points during execution at which Figure 7.3 provides additional statistics.

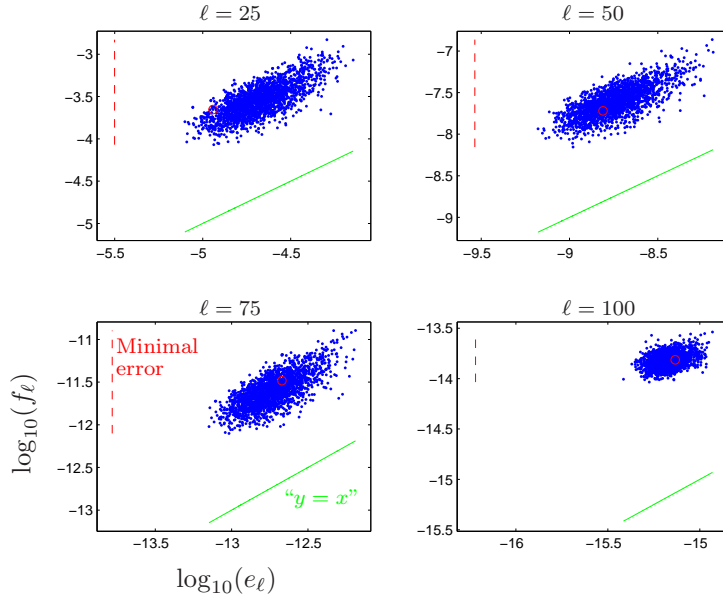


FIG. 7.3. Error statistics for approximating a Laplace integral operator. 2,000 trials of Algorithm 4.2 applied to a  $200 \times 200$  matrix approximating the integral operator (7.1). The panels isolate the moments at which  $\ell = 25, 50, 75, 100$  random samples have been drawn. Each solid point compares the estimated error  $f_\ell$  versus the actual error  $e_\ell$  in one trial; the open circle indicates the trial detailed in Figure 7.2. The dashed line identifies the minimal error  $\sigma_{\ell+1}$ , and the solid line marks the contour where the error estimator would equal the actual error.

We applied Algorithm 4.2 to the  $1596 \times 532$  matrix  $\mathbf{B}$  associated with a lattice in which there were 532 nodes (red) on the “inner ring” and 1596 nodes on the (blue) “outer ring.” Each application of  $\mathbf{B}$  to a vector requires the solution of a sparse linear system of size roughly  $140\,000 \times 140\,000$ . We implemented the scheme in Matlab using the “backslash” operator for the linear solve. The results of a typical trial appear in Figure 7.4. Qualitatively, the performance matches the results in Figure 7.3.

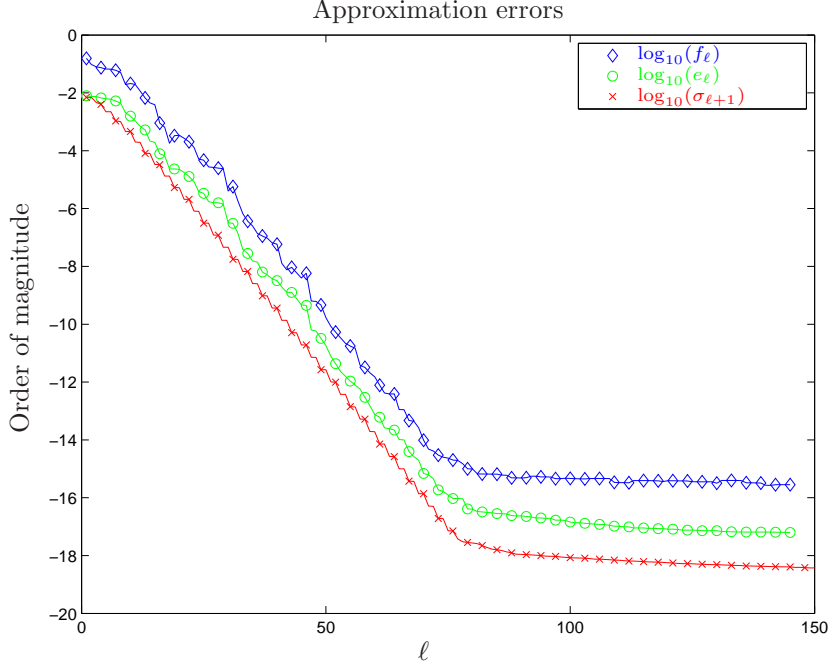


FIG. 7.4. Approximating the inverse of a discrete Laplacian. One execution of Algorithm 4.2 for the  $1596 \times 532$  input matrix  $\mathbf{B}$  described in §7.1. See Figure 7.2 for notations.

**7.2. A large, sparse, noisy matrix arising in image processing.** Our next example involves a matrix that arises in image processing. A recent line of work uses information about the local geometry of an image to develop promising new algorithms for standard tasks, such as denoising, inpainting, and so forth. These methods are based on approximating a *graph Laplacian* associated with the image. The dominant eigenvectors of this matrix provide “coordinates” that help us smooth out noisy image patches [120, 131].

We begin with a  $95 \times 95$  pixel grayscale image. The intensity of each pixel is represented as an integer in the range 0 to 4095. We form for each pixel  $i$  a vector  $\mathbf{x}^{(i)} \in \mathbb{R}^{25}$  by gathering the 25 intensities of the pixels in a  $5 \times 5$  neighborhood centered at pixel  $i$  (with appropriate modifications near the edges). Next, we form the  $9025 \times 9025$  *weight matrix*  $\widetilde{\mathbf{W}}$  that reflects the similarities between patches:

$$\widetilde{w}_{ij} = \exp \left\{ - \left\| \mathbf{x}^{(i)} - \mathbf{x}^{(j)} \right\|^2 / \sigma^2 \right\},$$

where the parameter  $\sigma = 50$  controls the level of sensitivity. We obtain a sparse weight matrix  $\mathbf{W}$  by zeroing out all entries in  $\widetilde{\mathbf{W}}$  except the seven largest ones in each row. The object is then to construct the low frequency eigenvectors of the graph

Laplacian matrix

$$\mathbf{L} = \mathbf{I} - \mathbf{D}^{-1/2} \mathbf{W} \mathbf{D}^{-1/2},$$

where  $\mathbf{D}$  is the diagonal matrix with entries  $d_{ii} = \sum_j w_{ij}$ . These are the eigenvectors associated with the dominant eigenvalues of the auxiliary matrix  $\mathbf{A} = \mathbf{D}^{-1/2} \mathbf{W} \mathbf{D}^{-1/2}$ .

The matrix  $\mathbf{A}$  is large, and its eigenvalues decay slowly, so we use the power scheme summarized in Algorithm 4.3 to approximate it. Figure 7.5[left] illustrates how the approximation error  $e_\ell$  declines as the number  $\ell$  of samples increases. When we set the exponent  $q = 0$ , which corresponds with the basic Algorithm 4.1, the approximation is rather poor. The graph illustrates that increasing the exponent  $q$  slightly results in a tremendous improvement in the accuracy of the power scheme.

Next, we illustrate the results of using the two-stage approach to approximate the eigenvalues of  $\mathbf{A}$ . In Stage A, we construct a basis for  $\mathbf{A}$  using Algorithm 4.3 with  $\ell = 100$  samples for different values of  $q$ . In Stage B, we apply the Hermitian variant of Algorithm 5.1 described in §5.3 to compute an approximate eigenvalue decomposition. Figure 7.5[right] shows the approximate eigenvalues and the actual eigenvalues of  $\mathbf{A}$ . Once again, we see that the minimal exponent  $q = 0$  produces miserable results, but the largest eigenvalues are quite accurate even for  $q = 1$ .

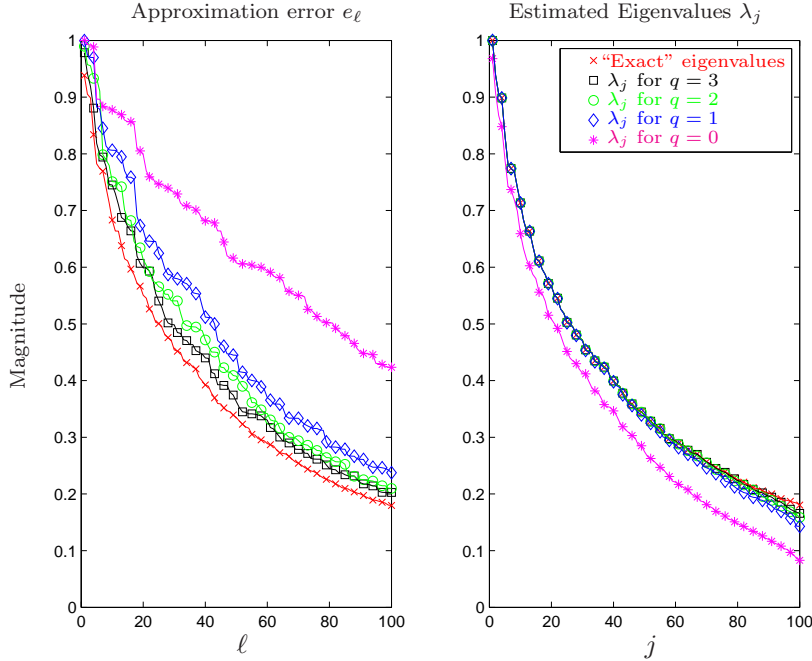


FIG. 7.5. Approximating a graph Laplacian. For varying exponent  $q$ , one trial of the power scheme, Algorithm 4.3, applied to the  $9025 \times 9025$  matrix  $\mathbf{A}$  described in §7.2. [Left] Approximation errors as a function of the number  $\ell$  of random samples. [Right] Estimates for the 100 largest eigenvalues given  $\ell = 100$  random samples compared with the 100 largest eigenvalues of  $\mathbf{A}$ .

**7.3. Eigenfaces.** Our next example involves a large, dense matrix derived from the FERET databank of face images [107, 108]. A simple method for performing face recognition is to identify the principal directions of the image data, which are called *eigenfaces*. Each of the original photographs can be summarized by its components

along these principal directions. To identify the subject in a new picture, we compute its decomposition in this basis and use a classification technique, such as nearest neighbors, to select the closest image in the database [122].

We construct a data matrix  $\mathbf{A}$  as follows: The FERET database contains 7254 images, and each  $384 \times 256$  image contains 98 304 pixels. First, we build a  $98\,304 \times 7254$  matrix  $\tilde{\mathbf{A}}$  whose columns are the images. We form  $\mathbf{A}$  by centering each column of  $\tilde{\mathbf{A}}$  and scaling it to unit norm, so that the images are roughly comparable. The eigenfaces are the dominant left singular vectors of this matrix.

Our goal then is to compute an approximate SVD of the matrix  $\mathbf{A}$ . Represented as an array of double-precision real numbers,  $\mathbf{A}$  would require 5.4 GB of storage, which does not fit within the fast memory of many machines. It is possible to compress the database down to at 57 MB or less (in JPEG format), but then the data would have to be uncompressed with each sweep over the matrix. Furthermore, the matrix  $\mathbf{A}$  has slowly decaying singular values, so we need to use the power scheme, Algorithm 4.3, to capture the range of the matrix accurately.

To address these concerns, we implemented the power scheme to run in a pass-efficient manner. An additional difficulty arises because the size of the data makes it expensive to calculate the actual error  $e_\ell$  incurred by the approximation or to determine the minimal error  $\sigma_{\ell+1}$ . To estimate the errors, we use the technique described in Remark 4.1.

Figure 7.6 describes the behavior of the power scheme, which is similar to its performance for the graph Laplacian in §7.2. When the exponent  $q = 0$ , the approximation of the data matrix is very poor, but it improves quickly as  $q$  increases. Likewise, the estimate for the spectrum of  $\mathbf{A}$  appears to converge rapidly; the largest singular values are already quite accurate when  $q = 1$ . We see essentially no improvement in the estimates after the first 3–5 passes over the matrix.

**7.4. Performance of structured random matrices.** Our final set of experiments illustrates that the structured random matrices described in §4.6 lead to matrix approximation algorithms that are both fast and accurate.

First, we compare the computational speeds of four methods for computing an approximation to the  $\ell$  dominant terms in the SVD of an  $n \times n$  matrix  $\mathbf{A}$ . For now, we are interested in execution time only (not accuracy), so the choice of matrix is irrelevant and we have selected  $\mathbf{A}$  to be a Gaussian matrix. The four methods are summarized in the following table; Remark 7.1 provides more details on the implementation.

Method	Stage A	Stage B
direct	Rank-revealing QR executed using column pivoting and Householder reflectors	Algorithm 5.1
gauss	Algorithm 4.1 with a Gaussian random matrix	Algorithm 5.1
srft	Algorithm 4.1 with the modified SRFT (4.8)	Algorithm 5.2
svd	Full SVD with LAPACK routine <code>dgesdd</code>	Truncate to $\ell$ terms

Table 7.1 lists the measured runtime of a single execution of each algorithm for various choices of the dimension  $n$  of the input matrix and the rank  $\ell$  of the approximation. Of course, the cost of the full SVD does not depend on the number  $\ell$  of components required. A more informative way to look at the runtime data is to compare the *relative* cost of the algorithms. The `direct` method is the best determin-



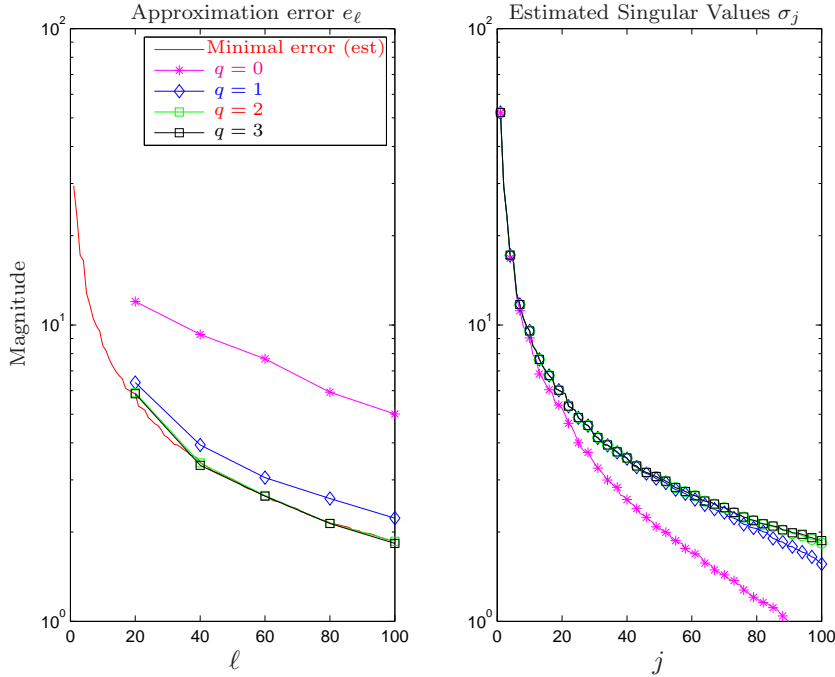


FIG. 7.6. Computing eigenfaces. For varying exponent  $q$ , one trial of the power scheme, Algorithm 4.3, applied to the  $98\,304 \times 7254$  matrix  $\mathbf{A}$  described in §7.3. (Left) Approximation errors as a function of the number  $\ell$  of random samples. The red line indicates the minimal errors as estimated by the singular values computed using  $\ell = 100$  and  $q = 3$ . (Right) Estimates for the 100 largest eigenvalues given  $\ell = 100$  random samples.

istic approach for dense matrices, so we calculate the factor by which the randomized methods improve on this benchmark. Figure 7.7 displays the results. We make two observations: (i) Using an SRFT often leads to a dramatic speed-up over classical techniques, even for moderate problem sizes. (ii) Using a standard Gaussian test matrix typically leads to a moderate speed-up over classical methods, primarily because performing a matrix–matrix multiplication is faster than a QR factorization.

Second, we investigate how the choice of random test matrix influences the error in approximating an input matrix. For these experiments, we return to the  $200 \times 200$  matrix  $\mathbf{A}$  defined in Section 7.1. Consider variations of Algorithm 4.1 obtained when the random test matrix  $\mathbf{\Omega}$  is drawn from the following four distributions:

- Gauss:** The standard Gaussian distribution.
- Ortho:** The uniform distribution on  $n \times \ell$  orthonormal matrices.
- SRFT:** The SRFT distribution defined in (4.6).
- GSFT:** The modified SRFT distribution defined in (4.8).

Intuitively, we expect that **Ortho** should provide the best performance.

For each distribution, we perform 100 000 trials of the following experiment. Apply the corresponding version of Algorithm 4.1 to the matrix  $\mathbf{A}$ , and calculate the approximation error  $e_\ell = \|\mathbf{A} - \mathbf{Q}_\ell \mathbf{Q}_\ell^* \mathbf{A}\|$ . Figure 7.8 displays the empirical probability density function for the error  $e_\ell$  obtained with each algorithm. We offer three observations: (i) The SRFT actually performs slightly *better* than a Gaussian ran-

TABLE 7.1

Computational times for a partial SVD. The time, in seconds, required to compute the  $\ell$  leading components in the SVD of an  $n \times n$  matrix using each of the methods from §7.4. The last row indicates the time needed to obtain a full SVD.

$\ell$	$n = 1024$			$n = 2048$			$n = 4096$		
	direct	gauss	srft	direct	gauss	srft	direct	gauss	srft
10	1.08e-1	5.63e-2	9.06e-2	4.22e-1	2.16e-1	3.56e-1	1.70e 0	8.94e-1	1.45e 0
20	1.97e-1	9.69e-2	1.03e-1	7.67e-1	3.69e-1	3.89e-1	3.07e 0	1.44e 0	1.53e 0
40	3.91e-1	1.84e-1	1.27e-1	1.50e 0	6.69e-1	4.33e-1	6.03e 0	2.64e 0	1.63e 0
80	7.84e-1	4.00e-1	2.19e-1	3.04e 0	1.43e 0	6.64e-1	1.20e 1	5.43e 0	2.08e 0
160	1.70e 0	9.92e-1	6.92e-1	6.36e 0	3.36e 0	1.61e 0	2.46e 1	1.16e 1	3.94e 0
320	3.89e 0	2.65e 0	2.98e 0	1.34e 1	7.45e 0	5.87e 0	5.00e 1	2.41e 1	1.21e 1
640	1.03e 1	8.75e 0	1.81e 1	3.14e 1	2.13e 1	2.99e 1	1.06e 2	5.80e 1	5.35e 1
1280	—	—	—	7.97e 1	6.69e 1	3.13e 2	2.40e 2	1.68e 2	4.03e 2
svd	1.19e 1			8.77e 1			6.90e 2		

dom matrix for this example. (ii) The standard SRFT and the modified SRFT have essentially identical errors. (iii) There is almost no difference between the Gaussian random matrix and the random orthonormal matrix in the first three plots, while the fourth plot shows that the random orthonormal matrix performs better. This behavior occurs because, with high probability, a tall Gaussian matrix is well conditioned and a (nearly) square Gaussian matrix is not.

REMARK 7.1. The running times reported in Table 7.1 and in Figure 7.7 depend strongly on both the computer hardware and the coding of the algorithms. The experiments reported here were performed on a standard office desktop with a 3.2 GHz Pentium IV processor and 2 GB of RAM. The algorithms were implemented in Fortran 90 and compiled with the Lahey compiler. The Lahey versions of BLAS and LAPACK were used to accelerate all matrix–matrix multiplications, as well as the SVD computations in Algorithms 5.1 and 5.2. We used the code for the modified SRFT (4.8) provided in the publicly available software package `id_dist` [90].

### Part III: Theory

This part of the paper, §§8–11, provides a detailed analysis of randomized sampling schemes for constructing an approximate basis for the range of a matrix, the task we refer to as Stage A in the framework of §1.2. More precisely, we assess the quality of the basis  $\mathbf{Q}$  that the proto-algorithm of §1.3 produces by establishing rigorous bounds for the approximation error

$$\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^*\mathbf{A}\|, \quad (7.2)$$

where  $\|\cdot\|$  denotes either the spectral norm or the Frobenius norm. The difficulty in developing these bounds is that the matrix  $\mathbf{Q}$  is random, and its distribution is a complicated nonlinear function of the input matrix  $\mathbf{A}$  and the random test matrix  $\mathbf{\Omega}$ . Naturally, any estimate for the approximation error must depend on the properties of the input matrix and the distribution of the test matrix.

To address these challenges, we split the argument into two pieces. The first part exploits techniques from linear algebra to deliver a generic error bound that depends on the interaction between the test matrix  $\mathbf{\Omega}$  and the right singular vectors of the input matrix  $\mathbf{A}$ , as well as the tail singular values of  $\mathbf{A}$ . In the second part of the argument, we take into account the distribution of the random matrix to estimate the error for specific instantiations of the proto-algorithm. This bipartite proof is

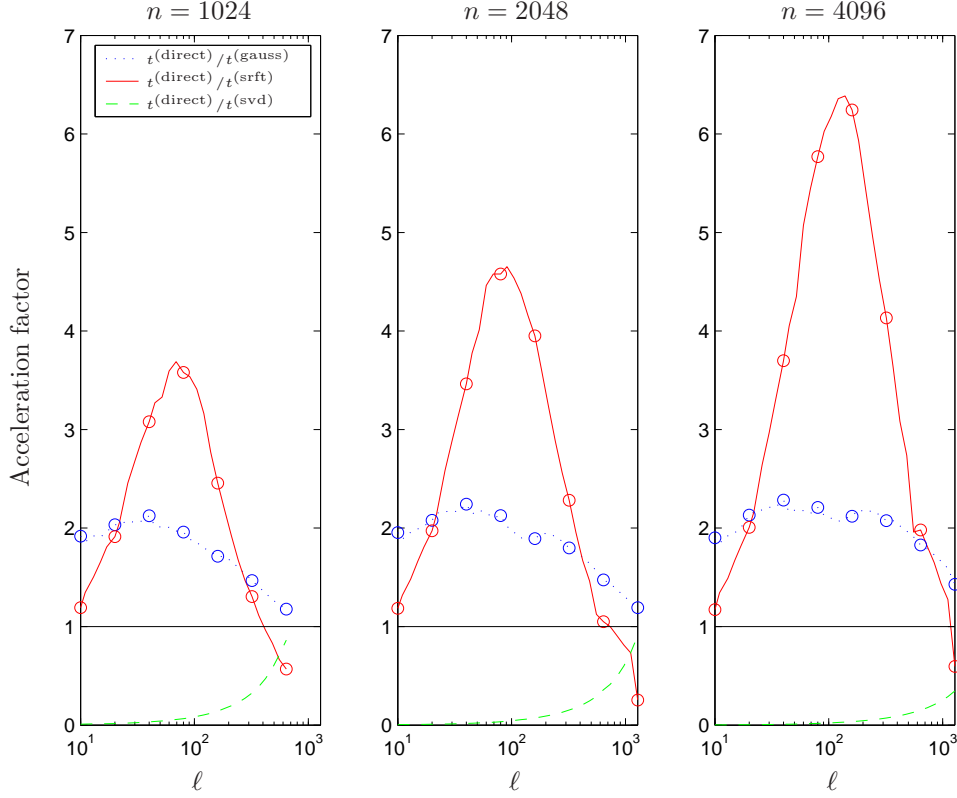


FIG. 7.7. Acceleration factor. *The relative cost of computing an  $\ell$ -term partial SVD of an  $n \times n$  Gaussian matrix using `direct`, a benchmark classical algorithm, versus each of the three competitors described in §7.4. The solid red curve shows the speedup using an SRFT test matrix, and the dotted blue curve shows the speedup with a Gaussian test matrix. The dashed green curve indicates that a full SVD computation using classical methods is substantially slower. Table 7.1 reports the absolute runtimes that yield the circled data points.*

common in the literature on randomized linear algebra, but our argument is most similar in spirit to [17].

Section 8 surveys the basic linear algebraic tools we need. Section 9 uses these methods to derive a generic error bound. Afterward, we specialize these results to the case where the test matrix is Gaussian (§10) and the case where the test matrix is a subsampled random Fourier transform (§11).

**8. Theoretical preliminaries.** We proceed with some additional background from linear algebra. Section 8.1 sets out properties of positive-semidefinite matrices, and §8.2 offers some results for orthogonal projectors. Standard references for this material include [11, 72].

**8.1. Positive semidefinite matrices.** An Hermitian matrix  $M$  is *positive semidefinite* (briefly, *psd*) when  $u^*Mu \geq 0$  for all  $u \neq 0$ . If the inequalities are strict,  $M$  is *positive definite* (briefly, *pd*). The psd matrices form a convex cone, which induces a partial ordering on the linear space of Hermitian matrices:  $M \preceq N$  if and only if  $N - M$  is psd. This ordering allows us to write  $M \succeq 0$  to indicate that the matrix  $M$  is psd.

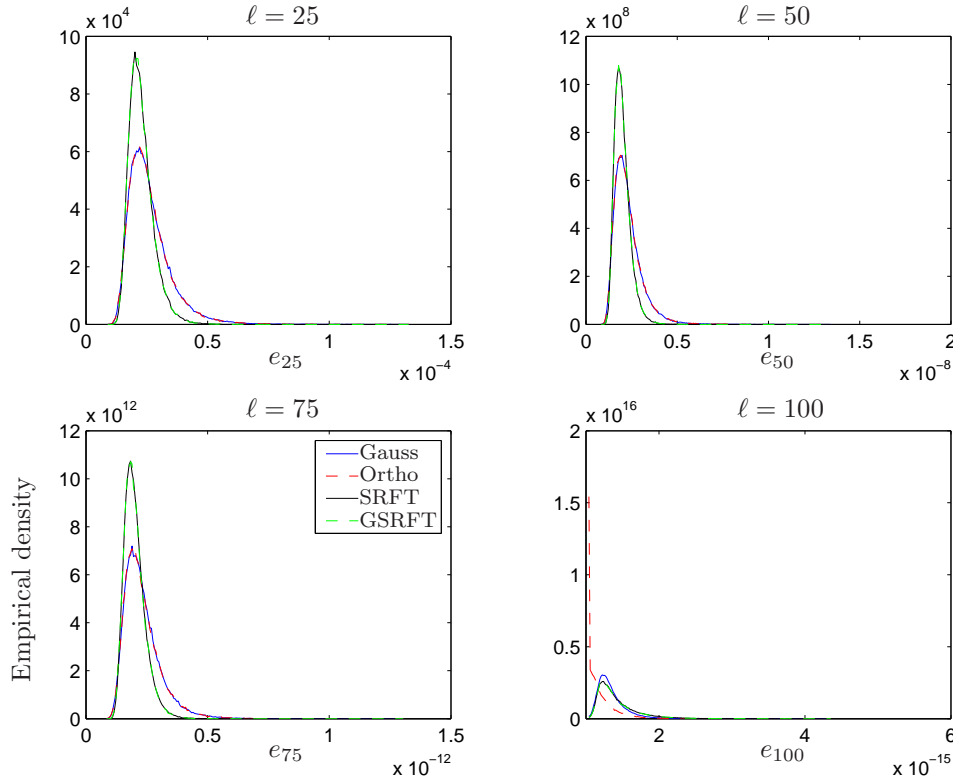


FIG. 7.8. Empirical probability density functions for the error in Algorithm 4.1. As described in §7.4, the algorithm is implemented with four distributions for the random test matrix and used to approximate the  $200 \times 200$  input matrix obtained by discretizing the integral operator (7.1). The four panels capture the empirical error distribution for each version of the algorithm at the moment when  $\ell = 25, 50, 75, 100$  random samples have been drawn.

Alternatively, we can define a psd (resp., pd) matrix as an Hermitian matrix with nonnegative (resp., positive) eigenvalues. In particular, each psd matrix is diagonalizable, and the inverse of a pd matrix is also pd. The spectral norm of a psd matrix  $\mathbf{M}$  has the variational characterization

$$\|\mathbf{M}\| = \max_{\mathbf{u} \neq \mathbf{0}} \frac{\mathbf{u}^* \mathbf{M} \mathbf{u}}{\mathbf{u}^* \mathbf{u}}, \quad (8.1)$$

according to the Rayleigh–Ritz theorem [72, Thm. 4.2.2]. It follows that

$$\mathbf{M} \preceq \mathbf{N} \implies \|\mathbf{M}\| \leq \|\mathbf{N}\|. \quad (8.2)$$

A fundamental fact is that conjugation preserves the psd property.

**PROPOSITION 8.1 (Conjugation Rule).** *Suppose that  $\mathbf{M} \succcurlyeq \mathbf{0}$ . For every  $\mathbf{A}$ , the matrix  $\mathbf{A}^* \mathbf{M} \mathbf{A} \succcurlyeq \mathbf{0}$ . In particular,*

$$\mathbf{M} \preceq \mathbf{N} \implies \mathbf{A}^* \mathbf{M} \mathbf{A} \preceq \mathbf{A}^* \mathbf{N} \mathbf{A}.$$

Our argument invokes the conjugation rule repeatedly. As a first application, we establish a perturbation bound for the matrix inverse near the identity matrix.

PROPOSITION 8.2 (Perturbation of Inverses). *Suppose that  $\mathbf{M} \succcurlyeq \mathbf{0}$ . Then*

$$\mathbf{I} - (\mathbf{I} + \mathbf{M})^{-1} \preccurlyeq \mathbf{M}$$

*Proof.* Define  $\mathbf{R} = \mathbf{M}^{1/2}$ , the psd square root of  $\mathbf{M}$  promised by [72, Thm. 7.2.6]. We have the chain of relations

$$\mathbf{I} - (\mathbf{I} + \mathbf{R}^2)^{-1} = (\mathbf{I} + \mathbf{R}^2)^{-1} \mathbf{R}^2 = \mathbf{R}(\mathbf{I} + \mathbf{R}^2)^{-1} \mathbf{R} \preccurlyeq \mathbf{R}^2.$$

The first equality can be verified algebraically. The second holds because rational functions of a diagonalizable matrix, such as  $\mathbf{R}$ , commute. The last relation follows from the conjugation rule because  $(\mathbf{I} + \mathbf{R}^2)^{-1} \preccurlyeq \mathbf{I}$ .  $\square$

Next, we present a generalization of the fact that the spectral norm of a psd matrix is controlled by its trace.

PROPOSITION 8.3. *We have  $\|\mathbf{M}\| \leq \|\mathbf{A}\| + \|\mathbf{C}\|$  for each partitioned psd matrix*

$$\mathbf{M} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^* & \mathbf{C} \end{bmatrix}.$$

*Proof.* The variational characterization (8.1) of the spectral norm implies that

$$\begin{aligned} \|\mathbf{M}\| &= \sup_{\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 = 1} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}^* \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^* & \mathbf{C} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \\ &\leq \sup_{\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 = 1} (\|\mathbf{A}\| \|\mathbf{x}\|^2 + 2 \|\mathbf{B}\| \|\mathbf{x}\| \|\mathbf{y}\| + \|\mathbf{C}\| \|\mathbf{y}\|^2). \end{aligned}$$

The block generalization of Hadamard's psd criterion [72, Thm. 7.7.7] states that  $\|\mathbf{B}\|^2 \leq \|\mathbf{A}\| \|\mathbf{C}\|$ . Thus,

$$\|\mathbf{M}\| \leq \sup_{\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 = 1} (\|\mathbf{A}\|^{1/2} \|\mathbf{x}\| + \|\mathbf{C}\|^{1/2} \|\mathbf{y}\|)^2 = \|\mathbf{A}\| + \|\mathbf{C}\|.$$

This point completes the argument.  $\square$

**8.2. Orthogonal projectors.** An *orthogonal projector* is an Hermitian matrix  $\mathbf{P}$  that satisfies the polynomial  $\mathbf{P}^2 = \mathbf{P}$ . This identity implies  $\mathbf{0} \preccurlyeq \mathbf{P} \preccurlyeq \mathbf{I}$ . An orthogonal projector is completely determined by its range. For a given matrix  $\mathbf{M}$ , we write  $\mathbf{P}_\mathbf{M}$  for the unique orthogonal projector with  $\text{range}(\mathbf{P}_\mathbf{M}) = \text{range}(\mathbf{M})$ . When  $\mathbf{M}$  has full column rank, we can express this projector explicitly:

$$\mathbf{P}_\mathbf{M} = \mathbf{M}(\mathbf{M}^* \mathbf{M})^{-1} \mathbf{M}^*. \quad (8.3)$$

The orthogonal projector onto the complementary subspace,  $\text{range}(\mathbf{P})^\perp$ , is the matrix  $\mathbf{I} - \mathbf{P}$ . Our argument hinges on several other facts about orthogonal projectors.

PROPOSITION 8.4. *Suppose  $\mathbf{U}$  is unitary. Then  $\mathbf{U}^* \mathbf{P}_\mathbf{M} \mathbf{U} = \mathbf{P}_{\mathbf{U}^* \mathbf{M}}$ .*

*Proof.* Abbreviate  $\mathbf{P} = \mathbf{U}^* \mathbf{P}_\mathbf{M} \mathbf{U}$ . It is clear that  $\mathbf{P}$  is an orthogonal projector since it is Hermitian and  $\mathbf{P}^2 = \mathbf{P}$ . Evidently,

$$\text{range}(\mathbf{P}) = \mathbf{U}^* \text{range}(\mathbf{M}) = \text{range}(\mathbf{U}^* \mathbf{M}).$$

Since the range determines the orthogonal projector, we conclude  $P = P_{U^*M}$ .  $\square$

PROPOSITION 8.5. *Suppose  $\text{range}(N) \subset \text{range}(M)$ . Then, for each matrix  $A$ , it holds that  $\|P_N A\| \leq \|P_M A\|$  and that  $\|(\mathbf{I} - P_M)A\| \leq \|(\mathbf{I} - P_N)A\|$ .*

*Proof.* The projector  $P_N \preceq \mathbf{I}$ , so the conjugation rule yields  $P_M P_N P_M \preceq P_M$ . The hypothesis  $\text{range}(N) \subset \text{range}(M)$  implies that  $P_M P_N = P_N$ , which results in

$$P_M P_N P_M = P_N P_M = (P_M P_N)^* = P_N.$$

In summary,  $P_N \preceq P_M$ . The conjugation rule shows that  $A^* P_N A \preceq A^* P_M A$ . We conclude from (8.2) that

$$\|P_N A\|^2 = \|A^* P_N A\| \leq \|A^* P_M A\| = \|P_M A\|^2.$$

The second statement follows from the first by taking orthogonal complements.  $\square$

Finally, we need a generalization of the scalar inequality  $|px|^q \leq |p| |x|^q$ , which holds when  $|p| \leq 1$  and  $q \geq 1$ .

PROPOSITION 8.6. *Let  $P$  be an orthogonal projector, and let  $M$  be a matrix. For each positive number  $q$ ,*

$$\|PM\| \leq \|P(MM^*)^q M\|^{1/(2q+1)}. \quad (8.4)$$

*Proof.* Suppose that  $R$  is an orthogonal projector,  $D$  is a nonnegative diagonal matrix, and  $t \geq 1$ . We claim that

$$\|RDR\|^t \leq \|RD^t R\|. \quad (8.5)$$

Granted this inequality, we quickly complete the proof. Using an SVD  $M = U\Sigma V^*$ , we compute

$$\begin{aligned} \|PM\|^{2(2q+1)} &= \|PMM^*P\|^{2q+1} = \|(U^*PU) \cdot \Sigma^2 \cdot (U^*PU)\|^{2q+1} \\ &\leq \|(U^*PU) \cdot \Sigma^{2(2q+1)} \cdot (U^*PU)\| = \|P(MM^*)^{2(2q+1)}P\| \\ &= \|P(MM^*)^q M \cdot M^* (MM^*)^q P\| = \|P(MM^*)^q M\|^2. \end{aligned}$$

We have used the unitary invariance of the spectral norm in the second and fourth relations. The inequality (8.5) applies because  $U^*PU$  is an orthogonal projector. Take a square root to finish the argument.

Now, we turn to the claim (8.5). This relation follows immediately from [11, Thm. IX.2.10], but we offer a direct argument based on more elementary considerations. Let  $x$  be a unit vector at which

$$x^*(RDR)x = \|RDR\|.$$

We must have  $Rx = x$ . Otherwise,  $\|Rx\| < 1$  because  $R$  is an orthogonal projector, which implies that the unit vector  $y = Rx / \|Rx\|$  verifies

$$y^*(RDR)y = \frac{(Rx)^*(RDR)(Rx)}{\|Rx\|^2} = \frac{x^*(RDR)x}{\|Rx\|^2} > x^*(RDR)x.$$

Writing  $x_j$  for the entries of  $\mathbf{x}$  and  $d_j$  for the diagonal entries of  $\mathbf{D}$ , we find that

$$\begin{aligned}\|\mathbf{RDR}\|^t &= [\mathbf{x}^*(\mathbf{RDR})\mathbf{x}]^t = [\mathbf{x}^*\mathbf{D}\mathbf{x}]^t = \left[\sum_j d_j x_j^2\right]^t \\ &\leq \left[\sum_j d_j^t x_j^2\right] = \mathbf{x}^*\mathbf{D}^t\mathbf{x} = (\mathbf{R}\mathbf{x})^*\mathbf{D}^t(\mathbf{R}\mathbf{x}) \leq \|\mathbf{RD}^t\mathbf{R}\|.\end{aligned}$$

The inequality is Jensen's, which applies because  $\sum x_j^2 = 1$  and the function  $z \mapsto |z|^t$  is convex for  $t \geq 1$ .  $\square$

**9. Error bounds via linear algebra.** We are now prepared to develop a deterministic error analysis for the proto-algorithm described in §1.3. To begin, we must introduce some notation. Afterward, we establish the key error bound, which strengthens a result from the literature [17, Lem. 4.2]. Finally, we explain why the power method can be used to improve the performance of the proto-algorithm.

**9.1. Setup.** Let  $\mathbf{A}$  be an  $m \times n$  matrix that has a singular value decomposition  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ , as described in Section 3.2.2. Roughly speaking, the proto-algorithm tries to approximate the subspace spanned by the first  $k$  left singular vectors, where  $k$  is now a fixed number. To perform the analysis, it is appropriate to partition the singular value decomposition as follows.

$$\mathbf{A} = \mathbf{U} \begin{bmatrix} k & n-k \\ \mathbf{\Sigma}_1 & \mathbf{\Sigma}_2 \end{bmatrix} \begin{bmatrix} n \\ \mathbf{V}_1^* \\ \mathbf{V}_2^* \end{bmatrix} \begin{matrix} k \\ n-k \end{matrix} \quad (9.1)$$

The matrices  $\mathbf{\Sigma}_1$  and  $\mathbf{\Sigma}_2$  are square. We will see that the left unitary factor  $\mathbf{U}$  does not play a significant role in the analysis.

Let  $\mathbf{\Omega}$  be an  $n \times \ell$  test matrix, where  $\ell$  denotes the number of samples. We assume only that  $\ell \geq k$ . Decompose the test matrix in the coordinate system determined by the right unitary factor of  $\mathbf{A}$ :

$$\mathbf{\Omega}_1 = \mathbf{V}_1^* \mathbf{\Omega} \quad \text{and} \quad \mathbf{\Omega}_2 = \mathbf{V}_2^* \mathbf{\Omega}. \quad (9.2)$$

The error bound for the proto-algorithm depends critically on the properties of the matrices  $\mathbf{\Omega}_1$  and  $\mathbf{\Omega}_2$ . With this notation, the sample matrix  $\mathbf{Y}$  can be expressed as

$$\mathbf{Y} = \mathbf{A}\mathbf{\Omega} = \mathbf{U} \begin{bmatrix} \mathbf{\Sigma}_1 \mathbf{\Omega}_1 \\ \mathbf{\Sigma}_2 \mathbf{\Omega}_2 \end{bmatrix} \begin{matrix} \ell \\ k \\ n-k \end{matrix}$$

It is a useful intuition that the block  $\mathbf{\Sigma}_1 \mathbf{\Omega}_1$  in (9.1) reflects the gross behavior of  $\mathbf{A}$ , while the block  $\mathbf{\Sigma}_2 \mathbf{\Omega}_2$  represents a perturbation.

**9.2. A deterministic error bound for the proto-algorithm.** The proto-algorithm constructs an orthonormal basis  $\mathbf{Q}$  for the range of the sample matrix  $\mathbf{Y}$ , and our goal is to quantify how well this basis captures the action of the input  $\mathbf{A}$ . Since  $\mathbf{Q}\mathbf{Q}^* = \mathbf{P}_Y$ , the challenge is to obtain bounds on the approximation error

$$\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^*\mathbf{A}\| = \|(\mathbf{I} - \mathbf{P}_Y)\mathbf{A}\|.$$

The following theorem shows that the behavior of the proto-algorithm depends on the interaction between the test matrix and the right singular vectors of the input matrix, as well as the singular spectrum of the input matrix.



**THEOREM 9.1** (Deterministic error bound). *Let  $\mathbf{A}$  be an  $m \times n$  matrix with singular value decomposition  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ , and fix  $k \geq 0$ . Choose a test matrix  $\mathbf{\Omega}$ , and construct the sample matrix  $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$ . Partition  $\mathbf{\Sigma}$  as specified in (9.1), and define  $\mathbf{\Omega}_1$  and  $\mathbf{\Omega}_2$  via (9.2). Assuming that  $\mathbf{\Omega}_1$  has full row rank, the approximation error satisfies*

$$\|(\mathbf{I} - \mathbf{P}_{\mathbf{Y}})\mathbf{A}\|^2 \leq \|\mathbf{\Sigma}_2\|^2 + \|\mathbf{\Sigma}_2\mathbf{\Omega}_2\mathbf{\Omega}_1^\dagger\|^2, \quad (9.3)$$

where  $\|\cdot\|$  denotes either the spectral norm or the Frobenius norm.

Theorem 9.1 sharpens the result [17, Lem. 2], which lacks the squares present in (9.3). This refinement yields slightly better error estimates than the earlier bound, and it has consequences for the probabilistic behavior of the error when the test matrix  $\mathbf{\Omega}$  is random. The proof here is different in spirit from the earlier analysis; our argument is inspired by the perturbation theory of orthogonal projectors [125].

*Proof.* We establish the bound for the spectral-norm error. The bound for the Frobenius-norm error follows from an analogous argument that is slightly easier.

Let us begin with some preliminary simplifications. First, we argue that the left unitary factor  $\mathbf{U}$  plays no essential role in the argument. In effect, we execute the proof for an auxiliary input matrix  $\tilde{\mathbf{A}}$  and an associated sample matrix  $\tilde{\mathbf{Y}}$  defined by

$$\tilde{\mathbf{A}} = \mathbf{U}^* \mathbf{A} = \begin{bmatrix} \mathbf{\Sigma}_1 \mathbf{V}_1^* \\ \mathbf{\Sigma}_2 \mathbf{V}_2^* \end{bmatrix} \quad \text{and} \quad \tilde{\mathbf{Y}} = \tilde{\mathbf{A}}\mathbf{\Omega} = \begin{bmatrix} \mathbf{\Sigma}_1 \mathbf{\Omega}_1 \\ \mathbf{\Sigma}_2 \mathbf{\Omega}_2 \end{bmatrix}. \quad (9.4)$$

Owing to the unitary invariance of the spectral norm and to Proposition 8.4, we have the identity

$$\|(\mathbf{I} - \mathbf{P}_{\mathbf{Y}})\mathbf{A}\| = \|\mathbf{U}^*(\mathbf{I} - \mathbf{P}_{\mathbf{Y}})\mathbf{U}\tilde{\mathbf{A}}\| = \|(\mathbf{I} - \mathbf{P}_{\mathbf{U}^*\mathbf{Y}})\tilde{\mathbf{A}}\| = \|(\mathbf{I} - \mathbf{P}_{\tilde{\mathbf{Y}}})\tilde{\mathbf{A}}\|. \quad (9.5)$$

In view of (9.5), it suffices to prove that

$$\|(\mathbf{I} - \mathbf{P}_{\tilde{\mathbf{Y}}})\tilde{\mathbf{A}}\| \leq \|\mathbf{\Sigma}_2\|^2 + \|\mathbf{\Sigma}_2\mathbf{\Omega}_2\mathbf{\Omega}_1^\dagger\|^2. \quad (9.6)$$

Second, we assume that the number  $k$  is chosen so the diagonal entries of  $\mathbf{\Sigma}_1$  are strictly positive. Suppose not. Then  $\mathbf{\Sigma}_2$  is zero because of the ordering of the singular values. As a consequence,

$$\text{range}(\tilde{\mathbf{A}}) = \text{range} \begin{bmatrix} \mathbf{\Sigma}_1 \mathbf{V}_1^* \\ \mathbf{0} \end{bmatrix} = \text{range} \begin{bmatrix} \mathbf{\Sigma}_1 \mathbf{\Omega}_1 \\ \mathbf{0} \end{bmatrix} = \text{range}(\tilde{\mathbf{Y}}).$$

This calculation uses the decompositions presented in (9.4), as well as the fact that both  $\mathbf{V}_1^*$  and  $\mathbf{\Omega}_1$  have full row rank. We conclude that

$$\|(\mathbf{I} - \mathbf{P}_{\tilde{\mathbf{Y}}})\tilde{\mathbf{A}}\| = 0,$$

so the error bound (9.6) holds trivially. (In fact, both sides are zero.)

The main argument is based on ideas from perturbation theory. To illustrate the concept, we start with a matrix related to  $\tilde{\mathbf{Y}}$ :

$$\mathbf{W} = \begin{bmatrix} \overset{\ell}{\mathbf{\Sigma}_1 \mathbf{\Omega}_1} \\ \mathbf{0} \end{bmatrix} \begin{matrix} k \\ n - k \end{matrix}$$

The matrix  $\mathbf{W}$  has the same range as a related matrix formed by “flattening out” the spectrum of the top block. Indeed, since  $\Sigma_1 \Omega_1$  has full row rank,

$$\text{range}(\mathbf{W}) = \text{range} \begin{bmatrix} \mathbf{I} \\ \mathbf{0} \end{bmatrix} \begin{matrix} k \\ n-k \end{matrix}$$

The matrix on the right-hand side has full column rank, so it is legal to apply the formula (8.3) for an orthogonal projector, which immediately yields

$$\mathbf{P}_\mathbf{W} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \quad \text{and} \quad \mathbf{I} - \mathbf{P}_\mathbf{W} = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}. \quad (9.7)$$

In words, the range of  $\mathbf{W}$  aligns with the first  $k$  coordinates, which span the same subspace as the first  $k$  left singular vectors of the auxiliary input matrix  $\tilde{\mathbf{A}}$ . Therefore,  $\text{range}(\mathbf{W})$  captures the action of  $\tilde{\mathbf{A}}$ , which is what we wanted from  $\text{range}(\tilde{\mathbf{Y}})$ .

We treat the auxiliary sample matrix  $\tilde{\mathbf{Y}}$  as a perturbation of  $\mathbf{W}$ , and we hope that their ranges are close to each other. To make the comparison rigorous, let us emulate the arguments outlined in the last paragraph. Referring to the display (9.4), we flatten out the top block of  $\tilde{\mathbf{Y}}$  to obtain the matrix

$$\mathbf{Z} = \tilde{\mathbf{Y}} \cdot \Omega_1^\dagger \Sigma_1^{-1} = \begin{bmatrix} \mathbf{I} \\ \mathbf{F} \end{bmatrix} \quad \text{where} \quad \mathbf{F} = \Sigma_2 \Omega_2 \Omega_1^\dagger \Sigma_1^{-1}. \quad (9.8)$$

Let us return to the error bound (9.6). The construction (9.8) ensures that  $\text{range}(\mathbf{Z}) \subset \text{range}(\tilde{\mathbf{Y}})$ , so Proposition 8.5 implies that the error satisfies

$$\|(\mathbf{I} - \mathbf{P}_{\tilde{\mathbf{Y}}})\tilde{\mathbf{A}}\| \leq \|(\mathbf{I} - \mathbf{P}_\mathbf{Z})\tilde{\mathbf{A}}\|.$$

Squaring this relation, we obtain

$$\|(\mathbf{I} - \mathbf{P}_{\tilde{\mathbf{Y}}})\tilde{\mathbf{A}}\|^2 \leq \|(\mathbf{I} - \mathbf{P}_\mathbf{Z})\tilde{\mathbf{A}}\|^2 = \|\tilde{\mathbf{A}}^*(\mathbf{I} - \mathbf{P}_\mathbf{Z})\tilde{\mathbf{A}}\| = \|\Sigma^*(\mathbf{I} - \mathbf{P}_\mathbf{Z})\Sigma\|. \quad (9.9)$$

The last identity follows from the definition  $\tilde{\mathbf{A}} = \Sigma \mathbf{V}^*$  and the unitary invariance of the spectral norm. Therefore, we can complete the proof of (9.6) by producing a suitable bound for the right-hand side of (9.9).

To continue, we need a detailed representation of the projector  $\mathbf{I} - \mathbf{P}_\mathbf{Z}$ . The construction (9.8) ensures that  $\mathbf{Z}$  has full column rank, so we can apply the formula (8.3) for an orthogonal projector to see that

$$\mathbf{P}_\mathbf{Z} = \mathbf{Z}(\mathbf{Z}^* \mathbf{Z})^{-1} \mathbf{Z}^* = \begin{bmatrix} \mathbf{I} \\ \mathbf{F} \end{bmatrix} (\mathbf{I} + \mathbf{F}^* \mathbf{F})^{-1} \begin{bmatrix} \mathbf{I} \\ \mathbf{F} \end{bmatrix}^*.$$

Expanding this expression, we determine that the complementary projector satisfies

$$\mathbf{I} - \mathbf{P}_\mathbf{Z} = \begin{bmatrix} \mathbf{I} - (\mathbf{I} + \mathbf{F}^* \mathbf{F})^{-1} & -(\mathbf{I} + \mathbf{F}^* \mathbf{F})^{-1} \mathbf{F}^* \\ -\mathbf{F}(\mathbf{I} + \mathbf{F}^* \mathbf{F})^{-1} & \mathbf{I} - \mathbf{F}(\mathbf{I} + \mathbf{F}^* \mathbf{F})^{-1} \mathbf{F} \end{bmatrix}. \quad (9.10)$$

The partitioning here conforms with the partitioning of  $\Sigma$ . When we conjugate the matrix by  $\Sigma$ , copies of  $\Sigma_1^{-1}$ , presently hidden in the top-left block, will cancel to happy effect.

The latter point may not seem obvious, owing to the complicated form of (9.10). In reality, the block matrix is less fearsome than it looks. Proposition 8.2, on the perturbation of inverses, shows that the top-left block verifies

$$\mathbf{I} - (\mathbf{I} + \mathbf{F}^* \mathbf{F})^{-1} \preceq \mathbf{F}^* \mathbf{F}.$$

The bottom-right block satisfies

$$\mathbf{I} - \mathbf{F}(\mathbf{I} + \mathbf{F}^* \mathbf{F})^{-1} \mathbf{F}^* \preceq \mathbf{I}$$

because the conjugation rule guarantees that  $\mathbf{F}(\mathbf{I} + \mathbf{F}^* \mathbf{F})^{-1} \mathbf{F}^* \succcurlyeq \mathbf{0}$ . We abbreviate the off-diagonal blocks with the symbol  $\mathbf{B} = -(\mathbf{I} + \mathbf{F}^* \mathbf{F})^{-1} \mathbf{F}^*$ . In summary,

$$\mathbf{I} - \mathbf{P}_Z \preceq \begin{bmatrix} \mathbf{F}^* \mathbf{F} & \mathbf{B} \\ \mathbf{B}^* & \mathbf{I} \end{bmatrix}.$$

This relation exposes the key structural properties of the projector. Compare this relation with the expression (9.7) for the “ideal” projector  $\mathbf{I} - \mathbf{P}_W$ .

Moving toward the estimate required by (9.9), we conjugate the last relation by  $\Sigma$  to obtain

$$\Sigma^* (\mathbf{I} - \mathbf{P}_Z) \Sigma \preceq \begin{bmatrix} \Sigma_1^* \mathbf{F}^* \mathbf{F} \Sigma_1 & \Sigma_1^* \mathbf{B} \Sigma_2 \\ \Sigma_2^* \mathbf{B}^* \Sigma_1 & \Sigma_2^* \Sigma_2 \end{bmatrix}.$$

The conjugation rule demonstrates that the matrix on the left-hand side is psd, so the matrix on the right-hand side is too. Proposition 8.3 results in the norm bound

$$\|\Sigma^* (\mathbf{I} - \mathbf{P}_Z) \Sigma\| \leq \|\Sigma_1^* \mathbf{F}^* \mathbf{F} \Sigma_1\| + \|\Sigma_2^* \Sigma_2\| = \|\mathbf{F} \Sigma_1\|^2 + \|\Sigma_2\|^2.$$

Recall that  $\mathbf{F} = \Sigma_2 \Omega_2 \Omega_1^\dagger \Sigma_1^{-1}$ , so the factor  $\Sigma_1$  cancels neatly. Therefore,

$$\|\Sigma^* (\mathbf{I} - \mathbf{P}_Z) \Sigma\| \leq \|\Sigma_2 \Omega_2 \Omega_1^\dagger\|^2 + \|\Sigma_2\|^2.$$

Finally, introduce the latter inequality into (9.9) to complete the proof.  $\square$

**9.3. Analysis of the power scheme.** Theorem 9.1 suggests that the performance of the proto-algorithm depends strongly on the relationship between the large singular values of  $\mathbf{A}$  listed in  $\Sigma_1$  and the small singular values listed in  $\Sigma_2$ . When a substantial proportion of the mass of  $\mathbf{A}$  appears in the small singular values, the constructed basis  $\mathbf{Q}$  may have low accuracy. Conversely, when the large singular values dominate, it is much easier to identify a good low-rank basis.

To improve the performance of the proto-algorithm, we can run it with a closely related input matrix whose singular values decay more rapidly [67, 112]. Fix a positive integer  $q$ , and set

$$\mathbf{B} = (\mathbf{A} \mathbf{A}^*)^q \mathbf{A} = \mathbf{U} \Sigma^{2q+1} \mathbf{V}^*.$$

We apply the proto-algorithm to  $\mathbf{B}$ , which generates a sample matrix  $\mathbf{Z} = \mathbf{B} \Omega$  and constructs a basis  $\mathbf{Q}$  for the range of  $\mathbf{Z}$ . Section 4.5 elaborates on the implementation details, and describes a reformulation that sometimes improves the accuracy when the scheme is executed in finite-precision arithmetic. The following result describes how well we can approximate the *original* matrix  $\mathbf{A}$  within the range of  $\mathbf{Z}$ .

**THEOREM 9.2 (Power scheme).** *Let  $\mathbf{A}$  be an  $m \times n$  matrix, and let  $\mathbf{\Omega}$  be an  $n \times \ell$  matrix. Fix a nonnegative integer  $q$ , form  $\mathbf{B} = (\mathbf{A}^* \mathbf{A})^q \mathbf{A}$ , and compute the sample matrix  $\mathbf{Z} = \mathbf{B} \mathbf{\Omega}$ . Then*

$$\|(\mathbf{I} - \mathbf{P}_{\mathbf{Z}}) \mathbf{A}\| \leq \|(\mathbf{I} - \mathbf{P}_{\mathbf{Z}}) \mathbf{B}\|^{1/(2q+1)}.$$

*Proof.* We determine that

$$\|(\mathbf{I} - \mathbf{P}_{\mathbf{Z}}) \mathbf{A}\| \leq \|(\mathbf{I} - \mathbf{P}_{\mathbf{Z}}) (\mathbf{A} \mathbf{A}^*)^q \mathbf{A}\|^{1/(2q+1)} = \|(\mathbf{I} - \mathbf{P}_{\mathbf{Z}}) \mathbf{B}\|^{1/(2q+1)}$$

as a direct consequence of Proposition 8.6.  $\square$

Let us illustrate how the power scheme interacts with the main error bound (9.3). Let  $\sigma_{k+1}$  denote the  $(k+1)$ th singular value of  $\mathbf{A}$ . First, suppose we approximate  $\mathbf{A}$  in the range of the sample matrix  $\mathbf{Y} = \mathbf{A} \mathbf{\Omega}$ . Since  $\|\Sigma_2\| = \sigma_{k+1}$ , Theorem 9.1 implies that

$$\|(\mathbf{I} - \mathbf{P}_{\mathbf{Y}}) \mathbf{A}\| \leq \left(1 + \|\mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger\|^2\right)^{1/2} \sigma_{k+1}. \quad (9.11)$$

Now, define  $\mathbf{B} = (\mathbf{A} \mathbf{A}^*)^q \mathbf{A}$ , and suppose we approximate  $\mathbf{A}$  within the range of the sample matrix  $\mathbf{Z} = \mathbf{B} \mathbf{\Omega}$ . Together, Theorem 9.2 and Theorem 9.1 imply that

$$\|(\mathbf{I} - \mathbf{P}_{\mathbf{Z}}) \mathbf{A}\| \leq \|(\mathbf{I} - \mathbf{P}_{\mathbf{Z}}) \mathbf{B}\|^{1/(2q+1)} \leq \left(1 + \|\mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger\|^2\right)^{1/(4q+2)} \sigma_{k+1}$$

because  $\sigma_{k+1}^{2q+1}$  is the  $(k+1)$ th singular value of  $\mathbf{B}$ . In effect, the power scheme drives down the suboptimality of the bound (9.11) exponentially fast as the power  $q$  increases. In principle, we can make the extra factor as close to one as we like, although this increases the cost of the algorithm.

**9.4. Analysis of truncated SVD.** Finally, let us study the truncated SVD described in Remark 5.1. Suppose that we approximate the input matrix  $\mathbf{A}$  inside the range of the sample matrix  $\mathbf{Z}$ . In essence, the truncation step computes a best rank- $k$  approximation  $\hat{\mathbf{A}}_{(k)}$  of the compressed matrix  $\mathbf{P}_{\mathbf{Z}} \mathbf{A}$ . The next result provides a simple error bound for this method; this argument was proposed by Ming Gu.

**THEOREM 9.3 (Analysis of Truncated SVD).** *Let  $\mathbf{A}$  be an  $m \times n$  matrix with singular values  $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots$ , and let  $\mathbf{Z}$  be an  $m \times \ell$  matrix, where  $\ell \geq k$ . Suppose that  $\hat{\mathbf{A}}_{(k)}$  is a best rank- $k$  approximation of  $\mathbf{P}_{\mathbf{Z}} \mathbf{A}$  with respect to the spectral norm. Then*

$$\|\mathbf{A} - \hat{\mathbf{A}}_{(k)}\| \leq \sigma_{k+1} + \|(\mathbf{I} - \mathbf{P}_{\mathbf{Z}}) \mathbf{A}\|.$$

*Proof.* Apply the triangle inequality to split the error into two components.

$$\|\mathbf{A} - \hat{\mathbf{A}}_{(k)}\| \leq \|\mathbf{A} - \mathbf{P}_{\mathbf{Z}} \mathbf{A}\| + \|\mathbf{P}_{\mathbf{Z}} \mathbf{A} - \hat{\mathbf{A}}_{(k)}\|. \quad (9.12)$$

We have already developed a detailed theory for estimating the first term. To analyze the second term, we introduce a best rank- $k$  approximation  $\mathbf{A}_{(k)}$  of the matrix  $\mathbf{A}$ . Note that

$$\|\mathbf{P}_{\mathbf{Z}} \mathbf{A} - \hat{\mathbf{A}}_{(k)}\| \leq \|\mathbf{P}_{\mathbf{Z}} \mathbf{A} - \mathbf{P}_{\mathbf{Z}} \mathbf{A}_{(k)}\|$$

because  $\hat{\mathbf{A}}_{(k)}$  is a best rank- $k$  approximation to the matrix  $\mathbf{P}_Z \mathbf{A}$ , whereas  $\mathbf{P}_Z \mathbf{A}_{(k)}$  is an undistinguished rank- $k$  matrix. It follows that

$$\|\mathbf{P}_Z \mathbf{A} - \hat{\mathbf{A}}_{(k)}\| \leq \|\mathbf{P}_Z(\mathbf{A} - \mathbf{A}_{(k)})\| \leq \|\mathbf{A} - \mathbf{A}_{(k)}\| = \sigma_{k+1}. \quad (9.13)$$

The second inequality holds because the orthogonal projector is a contraction; the last identity follows from Mirsky's theorem [97]. Combine (9.12) and (9.13) to reach the main result.  $\square$

REMARK 9.1. In the randomized setting, the truncation step appears to be less damaging than the error bound of Theorem 9.3 suggests, but we currently lack a complete theoretical understanding of its behavior.

**10. Gaussian test matrices.** The error bound in Theorem 9.1 shows that the performance of the proto-algorithm depends on the interaction between the test matrix  $\mathbf{\Omega}$  and the right singular vectors of the input matrix  $\mathbf{A}$ . Algorithm 4.1 is a particularly simple version of the proto-algorithm that draws the test matrix according to the standard Gaussian distribution. The literature contains a wealth of information about these matrices, which allows us to perform a very precise error analysis.

We focus on the real case in this section. Analogous results hold in the complex case, where the algorithm even exhibits superior performance.

**10.1. Technical background.** A *standard Gaussian matrix* is a random matrix whose entries are independent standard normal variables. The distribution of a standard Gaussian matrix is rotationally invariant: If  $\mathbf{U}$  and  $\mathbf{V}$  are orthonormal matrices, then  $\mathbf{U}^* \mathbf{G} \mathbf{V}$  also has the standard Gaussian distribution.

Our analysis requires detailed information about the properties of Gaussian matrices. In particular, we must understand how the norm of a Gaussian matrix and its pseudoinverse vary. We summarize the relevant results and citations here, reserving the details for Appendix A.

PROPOSITION 10.1 (Expected norm of a scaled Gaussian matrix). *Fix matrices  $\mathbf{S}, \mathbf{T}$ , and draw a standard Gaussian matrix  $\mathbf{G}$ . Then*

$$\left(\mathbb{E} \|\mathbf{SGT}\|_{\text{F}}^2\right)^{1/2} = \|\mathbf{S}\|_{\text{F}} \|\mathbf{T}\|_{\text{F}} \quad \text{and} \quad (10.1)$$

$$\mathbb{E} \|\mathbf{SGT}\| \leq \|\mathbf{S}\| \|\mathbf{T}\|_{\text{F}} + \|\mathbf{S}\|_{\text{F}} \|\mathbf{T}\|. \quad (10.2)$$

The identity (10.1) follows from a direct calculation. The second bound (10.2) relies on methods developed by Gordon [62, 63]. See Propositions A.1 and A.2.

PROPOSITION 10.2 (Expected norm of a pseudo-inverted Gaussian matrix). *Draw a  $k \times (k + p)$  standard Gaussian matrix  $\mathbf{G}$  with  $k \geq 2$  and  $p \geq 2$ . Then*

$$\left(\mathbb{E} \|\mathbf{G}^\dagger\|_{\text{F}}^2\right)^{1/2} = \sqrt{\frac{k}{p-1}} \quad \text{and} \quad (10.3)$$

$$\mathbb{E} \|\mathbf{G}^\dagger\| \leq \frac{e\sqrt{k+p}}{p}. \quad (10.4)$$

The first identity is a standard result from multivariate statistics [99, p. 96]. The second follows from work of Chen and Dongarra [25]. See Proposition A.4 and A.5.

To study the probability that Algorithm 4.1 produces a large error, we rely on tail bounds for functions of Gaussian matrices. The next proposition rephrases a well-known result on concentration of measure [14, Thm. 4.5.7]. See also [83, §1.1] and [82, §5.1].

**PROPOSITION 10.3** (Concentration for functions of a Gaussian matrix). *Suppose that  $h$  is a Lipschitz function on matrices:*

$$|h(\mathbf{X}) - h(\mathbf{Y})| \leq L \|\mathbf{X} - \mathbf{Y}\|_F \quad \text{for all } \mathbf{X}, \mathbf{Y}.$$

*Draw a standard Gaussian matrix  $\mathbf{G}$ . Then*

$$\mathbb{P}\{h(\mathbf{G}) \geq \mathbb{E}h(\mathbf{G}) + Lt\} \leq e^{-t^2/2}.$$

Finally, we state some large deviation bounds for the norm of a pseudo-inverted Gaussian matrix.

**PROPOSITION 10.4** (Norm bounds for a pseudo-inverted Gaussian matrix). *Let  $\mathbf{G}$  be a  $k \times (k + p)$  Gaussian matrix where  $p \geq 4$ . For all  $t \geq 1$ ,*

$$\mathbb{P}\left\{\|\mathbf{G}^\dagger\|_F \geq \sqrt{\frac{12k}{p}} \cdot t\right\} \leq 4t^{-p} \quad \text{and} \quad (10.5)$$

$$\mathbb{P}\left\{\|\mathbf{G}^\dagger\| \geq \frac{e\sqrt{k+p}}{p+1} \cdot t\right\} \leq t^{-(p+1)}. \quad (10.6)$$

Compare these estimates with Proposition 10.2. It seems that (10.5) is new; we were unable to find a comparable analysis in the random matrix literature. Although the form of (10.5) is not optimal, it allows us to produce more transparent results than a fully detailed estimate. The bound (10.6) essentially appears in the work of Chen and Dongarra [25]. See Propositions A.3 and Theorem A.6 for more information.

**10.2. Average-case analysis of Algorithm 4.1.** We separate our analysis into two pieces. First, we present information about expected values. In the next subsection, we describe bounds on the probability of a large deviation.

We begin with the simplest result, which provides an estimate for the expected approximation error in the Frobenius norm. All proofs are postponed to the end of the section.

**THEOREM 10.5** (Average Frobenius error). *Suppose that  $\mathbf{A}$  is a real  $m \times n$  matrix with singular values  $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots$ . Choose a target rank  $k \geq 2$  and an oversampling parameter  $p \geq 2$ , where  $k + p \leq \min\{m, n\}$ . Draw an  $n \times (k + p)$  standard Gaussian matrix  $\mathbf{\Omega}$ , and construct the sample matrix  $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$ . Then the expected approximation error*

$$\mathbb{E}\|(\mathbf{I} - \mathbf{P}_Y)\mathbf{A}\|_F \leq \left(1 + \frac{k}{p-1}\right)^{1/2} \left(\sum_{j>k} \sigma_j^2\right)^{1/2}.$$

This theorem predicts several intriguing behaviors of Algorithm 4.1. The Eckart–Young theorem [54] shows that  $(\sum_{j>k} \sigma_j^2)^{1/2}$  is the minimal Frobenius-norm error when approximating  $\mathbf{A}$  with a rank- $k$  matrix. This quantity is the appropriate benchmark for the performance of the algorithm. If the small singular values of  $\mathbf{A}$  are very flat, the series may be as large as  $\sigma_{k+1} \sqrt{\min\{m, n\} - k}$ . On the other hand, when the singular values exhibit some decay, the error may be on the same order as  $\sigma_{k+1}$ .

The error bound always exceeds this baseline error, but it may be polynomially larger, depending on the ratio between the target rank  $k$  and the oversampling parameter  $p$ . For  $p$  small (say, less than five), the error is somewhat variable because the small singular values of a nearly square Gaussian matrix are very unstable. As the oversampling increases, the performance improves quickly. When  $p \sim k$ , the error is already within a constant factor of the baseline.

The error bound for the spectral norm is somewhat more complicated, but it reveals some interesting new features.

**THEOREM 10.6** (Average spectral error). *Under the hypotheses of Theorem 10.5,*

$$\mathbb{E} \|(\mathbf{I} - \mathbf{P}_Y)\mathbf{A}\| \leq \left(1 + \sqrt{\frac{k}{p-1}}\right) \sigma_{k+1} + \frac{e\sqrt{k+p}}{p} \left(\sum_{j>k} \sigma_j^2\right)^{1/2}.$$

Mirsky [97] has shown that the quantity  $\sigma_{k+1}$  is the minimum spectral-norm error when approximating  $\mathbf{A}$  with a rank- $k$  matrix, so the first term in Theorem 10.6 is analogous with the error bound in Theorem 10.5. The second term represents a new phenomenon: we also pay for the Frobenius-norm error in approximating  $\mathbf{A}$ . Note that, as the amount  $p$  of oversampling increases, the polynomial factor in the second term declines much more quickly than the factor in the first term. When  $p \sim k$ , the factor on the  $\sigma_{k+1}$  term is constant, while the factor on the series has order  $k^{-1/2}$ .

We also note that the bound in Theorem 10.6 implies

$$\mathbb{E} \|(\mathbf{I} - \mathbf{P}_Y)\mathbf{A}\| \leq \left[1 + \sqrt{\frac{k}{p-1}} + \frac{e\sqrt{k+p}}{p} \cdot \sqrt{\min\{m, n\} - k}\right] \sigma_{k+1},$$

so the average spectral-norm error always lies within a small polynomial factor of the baseline  $\sigma_{k+1}$ .

Let us continue with the proofs of these results.

*Proof.* [Theorem 10.5] Let  $\mathbf{V}$  be the right unitary factor of  $\mathbf{A}$ . Partition  $\mathbf{V} = [\mathbf{V}_1 \mid \mathbf{V}_2]$  into blocks containing, respectively,  $k$  and  $n - k$  columns. Recall that

$$\mathbf{\Omega}_1 = \mathbf{V}_1^* \mathbf{\Omega} \quad \text{and} \quad \mathbf{\Omega}_2 = \mathbf{V}_2^* \mathbf{\Omega}.$$

The Gaussian distribution is rotationally invariant, so  $\mathbf{V}^* \mathbf{\Omega}$  is also a standard Gaussian matrix. Observe that  $\mathbf{\Omega}_1$  and  $\mathbf{\Omega}_2$  are *nonoverlapping* submatrices of  $\mathbf{V}^* \mathbf{\Omega}$ , so these two matrices are not only standard Gaussian but also stochastically independent. Furthermore, the rows of a (fat) Gaussian matrix are almost surely in general position, so the  $k \times (k + p)$  matrix  $\mathbf{\Omega}_1$  has full row rank with probability one.

Hölder's inequality and Theorem 9.1 together imply that

$$\mathbb{E} \|(\mathbf{I} - \mathbf{P}_Y)\mathbf{A}\|_{\text{F}} \leq \left(\mathbb{E} \|(\mathbf{I} - \mathbf{P}_Y)\mathbf{A}\|_{\text{F}}^2\right)^{1/2} \leq \left(\|\mathbf{\Sigma}_2\|_{\text{F}}^2 + \mathbb{E} \|\mathbf{\Sigma}_2 \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger\|_{\text{F}}^2\right)^{1/2}.$$



We compute this expectation by conditioning on the value of  $\mathbf{\Omega}_1$  and applying Proposition 10.1 to the scaled Gaussian matrix  $\mathbf{\Omega}_2$ . Thus,

$$\begin{aligned}\mathbb{E} \|\mathbf{\Sigma}_2 \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger\|_F^2 &= \mathbb{E} \left( \mathbb{E} \left[ \|\mathbf{\Sigma}_2 \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger\|_F^2 \mid \mathbf{\Omega}_1 \right] \right) = \mathbb{E} \left( \|\mathbf{\Sigma}_2\|_F^2 \|\mathbf{\Omega}_1^\dagger\|_F^2 \right) \\ &= \|\mathbf{\Sigma}_2\|_F^2 \cdot \mathbb{E} \|\mathbf{\Omega}_1^\dagger\|_F^2 = \frac{k}{p-1} \cdot \|\mathbf{\Sigma}_2\|_F^2,\end{aligned}$$

where the last expectation follows from relation (10.3) of Proposition 10.2. In summary,

$$\mathbb{E} \|(\mathbf{I} - \mathbf{P}_Y) \mathbf{A}\|_F \leq \left( 1 + \frac{k}{p-1} \right)^{1/2} \|\mathbf{\Sigma}_2\|_F.$$

Observe that  $\|\mathbf{\Sigma}_2\|_F^2 = \sum_{j>k} \sigma_j^2$  to complete the proof.  $\square$

*Proof.* [Theorem 10.6] The argument is similar to the proof of Theorem 10.5. First, Theorem 9.1 implies that

$$\mathbb{E} \|(\mathbf{I} - \mathbf{P}_Y) \mathbf{A}\| \leq \mathbb{E} \left( \|\mathbf{\Sigma}_2\|^2 + \|\mathbf{\Sigma}_2 \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger\|^2 \right)^{1/2} \leq \|\mathbf{\Sigma}_2\| + \mathbb{E} \|\mathbf{\Sigma}_2 \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger\|.$$

We condition on  $\mathbf{\Omega}_1$  and apply Proposition 10.1 to bound the expectation with respect to  $\mathbf{\Omega}_2$ . Thus,

$$\begin{aligned}\mathbb{E} \|\mathbf{\Sigma}_2 \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger\| &\leq \mathbb{E} \left( \|\mathbf{\Sigma}_2\| \|\mathbf{\Omega}_1^\dagger\|_F + \|\mathbf{\Sigma}_2\|_F \|\mathbf{\Omega}_1^\dagger\| \right) \\ &\leq \|\mathbf{\Sigma}_2\| \left( \mathbb{E} \|\mathbf{\Omega}_1^\dagger\|_F^2 \right)^{1/2} + \|\mathbf{\Sigma}_2\|_F \cdot \mathbb{E} \|\mathbf{\Omega}_1^\dagger\|.\end{aligned}$$

where the second relation requires Hölder's inequality. Applying both parts of Proposition 10.2, we obtain

$$\mathbb{E} \|\mathbf{\Sigma}_2 \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger\| \leq \sqrt{\frac{k}{p-1}} \|\mathbf{\Sigma}_2\| + \frac{e\sqrt{k+p}}{p} \|\mathbf{\Sigma}_2\|_F.$$

Note that  $\|\mathbf{\Sigma}_2\| = \sigma_{k+1}$  to wrap up.  $\square$

**10.3. Probabilistic error bounds for Algorithm 4.1.** We can develop tail bounds for the approximation error, which demonstrate that the average performance of the algorithm is representative of the actual performance. We begin with the Frobenius norm because the result is somewhat simpler.

**THEOREM 10.7** (Deviation bounds for the Frobenius error). *Frame the hypotheses of Theorem 10.5. Assume further that  $p \geq 4$ . For all  $u, t \geq 1$ ,*

$$\|(\mathbf{I} - \mathbf{P}_Y) \mathbf{A}\|_F \leq \left( 1 + t \cdot \sqrt{12k/p} \right) \left( \sum_{j>k} \sigma_j^2 \right)^{1/2} + ut \cdot \frac{e\sqrt{k+p}}{p+1} \cdot \sigma_{k+1},$$

*with failure probability at most  $5t^{-p} + 2e^{-u^2/2}$ .*

To parse this theorem, observe that the first term in the error bound corresponds with the expected approximation error in Theorem 10.5. The second term represents a deviation above the mean.

An analogous result holds for the spectral norm.

**THEOREM 10.8** (Deviation bounds for the spectral error). *Frame the hypotheses of Theorem 10.5. Assume further that  $p \geq 4$ . For all  $u, t \geq 1$ ,*

$$\begin{aligned} & \|(\mathbf{I} - \mathbf{P}_Y)\mathbf{A}\| \\ & \leq \left[ \left(1 + t \cdot \sqrt{12k/p}\right) \sigma_{k+1} + t \cdot \frac{e\sqrt{k+p}}{p+1} \left(\sum_{j>k} \sigma_j^2\right)^{1/2} \right] + ut \cdot \frac{e\sqrt{k+p}}{p+1} \sigma_{k+1}, \end{aligned}$$

with failure probability at most  $5t^{-p} + e^{-u^2/2}$ .

The bracket corresponds with the expected spectral-norm error while the remaining term represents a deviation above the mean. Neither the numerical constants nor the precise form of the bound are optimal because of the slackness in Proposition 10.4. Nevertheless, the theorem gives a fairly good picture of what is actually happening.

We acknowledge that the current form of Theorem 10.8 is complicated. To produce more transparent results, we make appropriate selections for the parameters  $u, t$  and bound the numerical constants.

**COROLLARY 10.9** (Simplified deviation bounds for the spectral error). *Frame the hypotheses of Theorem 10.5, and assume further that  $p \geq 4$ . Then*

$$\|(\mathbf{I} - \mathbf{P}_Y)\mathbf{A}\| \leq \left(1 + 17\sqrt{1 + k/p}\right) \sigma_{k+1} + \frac{8\sqrt{k+p}}{p+1} \left(\sum_{j>k} \sigma_j^2\right)^{1/2},$$

with failure probability at most  $6e^{-p}$ . Moreover,

$$\|(\mathbf{I} - \mathbf{P}_Y)\mathbf{A}\| \leq \left(1 + 8\sqrt{(k+p) \cdot p \log p}\right) \sigma_{k+1} + 3\sqrt{k+p} \left(\sum_{j>k} \sigma_j^2\right)^{1/2},$$

with failure probability at most  $6p^{-p}$ .

*Proof.* The first part of the result follows from the choices  $t = e$  and  $u = \sqrt{2p}$ , and the second emerges when  $t = p$  and  $u = \sqrt{2p \log p}$ . Another interesting parameter selection is  $t = p^{c/p}$  and  $u = \sqrt{2c \log p}$ , which yields a failure probability  $6p^{-c}$ .  $\square$

Corollary 10.9 should be compared with [91, Obs. 4.4–4.5]. Although our result contains sharper error estimates, the failure probabilities are usually worse. The error bound (1.9) presented in §1.5 follows after further simplification of the second bound from Corollary 10.9.

We continue with a proof of Theorem 10.8. The same argument can be used to obtain a bound for the Frobenius-norm error, but we omit a detailed account.

*Proof.* [Theorem 10.8] Since  $\mathbf{\Omega}_1$  and  $\mathbf{\Omega}_2$  are independent from each other, we can study how the error depends on the matrix  $\mathbf{\Omega}_2$  by conditioning on the event that  $\mathbf{\Omega}_1$  is not too irregular. To that end, we define a (parameterized) event on which the spectral and Frobenius norms of the matrix  $\mathbf{\Omega}_1^\dagger$  are both controlled. For  $t \geq 1$ , let

$$E_t = \left\{ \mathbf{\Omega}_1 : \|\mathbf{\Omega}_1^\dagger\| \leq \frac{e\sqrt{k+p}}{p+1} \cdot t \quad \text{and} \quad \|\mathbf{\Omega}_1^\dagger\|_F \leq \sqrt{\frac{12k}{p}} \cdot t \right\}.$$

Invoking both parts of Proposition 10.4, we find that

$$\mathbb{P}(E_t^c) \leq t^{-(p+1)} + 4t^{-p} \leq 5t^{-p}.$$

Consider the function  $h(\mathbf{X}) = \|\Sigma_2 \mathbf{X} \Omega_1^\dagger\|$ . We quickly compute its Lipschitz constant  $L$  with the lower triangle inequality and some standard norm estimates:

$$\begin{aligned} |h(\mathbf{X}) - h(\mathbf{Y})| &\leq \|\Sigma_2(\mathbf{X} - \mathbf{Y})\Omega_1^\dagger\| \\ &\leq \|\Sigma_2\| \|\mathbf{X} - \mathbf{Y}\| \|\Omega_1^\dagger\| \leq \|\Sigma_2\| \|\Omega_1^\dagger\| \|\mathbf{X} - \mathbf{Y}\|_F. \end{aligned}$$

Therefore,  $L \leq \|\Sigma_2\| \|\Omega_1^\dagger\|$ . Relation (10.2) of Proposition 10.1 implies that

$$\mathbb{E}[h(\Omega_2) \mid \Omega_1] \leq \|\Sigma_2\| \|\Omega_1^\dagger\|_F + \|\Sigma_2\|_F \|\Omega_1^\dagger\|.$$

Applying the concentration of measure inequality, Proposition 10.3, conditionally to the random variable  $h(\Omega_2) = \|\Sigma_2 \Omega_2 \Omega_1^\dagger\|$  results in

$$\mathbb{P}\left\{\|\Sigma_2 \Omega_2 \Omega_1^\dagger\| > \|\Sigma_2\| \|\Omega_1^\dagger\|_F + \|\Sigma_2\|_F \|\Omega_1^\dagger\| + \|\Sigma_2\| \|\Omega_1^\dagger\| \cdot u \mid E_t\right\} \leq e^{-u^2/2}.$$

Under the event  $E_t$ , we have explicit bounds on the norms of  $\Omega_1^\dagger$ , so

$$\begin{aligned} \mathbb{P}\left\{\|\Sigma_2 \Omega_2 \Omega_1^\dagger\| > \|\Sigma_2\| \sqrt{\frac{12k}{p}} \cdot t + \|\Sigma_2\|_F \frac{e\sqrt{k+p}}{p+1} \cdot t + \|\Sigma_2\| \frac{e\sqrt{k+p}}{p+1} \cdot ut \mid E_t\right\} \\ \leq e^{-u^2/2}. \end{aligned}$$

Use the fact  $\mathbb{P}(E_t^c) \leq 5t^{-p}$  to remove the conditioning. Therefore,

$$\begin{aligned} \mathbb{P}\left\{\|\Sigma_2 \Omega_2 \Omega_1^\dagger\| > \|\Sigma_2\| \sqrt{\frac{12k}{p}} \cdot t + \|\Sigma_2\|_F \frac{e\sqrt{k+p}}{p+1} \cdot t + \|\Sigma_2\| \frac{e\sqrt{k+p}}{p+1} \cdot ut\right\} \\ \leq 5t^{-p} + e^{-u^2/2}. \end{aligned}$$

Insert the expressions for the norms of  $\Sigma_2$  into this result to complete the probability bound. Finally, introduce this estimate into the error bound from Theorem 9.1.  $\square$

**10.4. Analysis of the power scheme.** Theorem 10.6 makes it clear that the performance of the randomized approximation scheme, Algorithm 4.1, depends heavily on the singular spectrum of the input matrix. The power scheme outlined in Algorithm 4.3 addresses this problem by enhancing the decay of spectrum. We can combine our analysis of Algorithm 4.1 with Theorem 9.2 to obtain a detailed report on the behavior of the performance of the power scheme using a Gaussian matrix.

**COROLLARY 10.10** (Average spectral error for the power scheme). *Frame the hypotheses of Theorem 10.5. Define  $\mathbf{B} = (\mathbf{A}\mathbf{A}^*)^q \mathbf{A}$  for a nonnegative integer  $q$ , and construct the sample matrix  $\mathbf{Z} = \mathbf{B}\Omega$ . Then*

$$\mathbb{E}\|(\mathbf{I} - \mathbf{P}_Z)\mathbf{A}\| \leq \left[ \left(1 + \sqrt{\frac{k}{p-1}}\right) \sigma_{k+1}^{2q+1} + \frac{e\sqrt{k+p}}{p} \left(\sum_{j>k} \sigma_j^{2(2q+1)}\right)^{1/2} \right]^{1/(2q+1)}.$$

*Proof.* By Hölder's inequality and Theorem 9.2,

$$\mathbb{E}\|(\mathbf{I} - \mathbf{P}_Z)\mathbf{A}\| \leq \left(\mathbb{E}\|(\mathbf{I} - \mathbf{P}_Z)\mathbf{A}\|^{2q+1}\right)^{1/(2q+1)} \leq \left(\mathbb{E}\|(\mathbf{I} - \mathbf{P}_Z)\mathbf{B}\|\right)^{1/(2q+1)}.$$

Invoke Theorem 10.6 to bound the right-hand side, noting that  $\sigma_j(\mathbf{B}) = \sigma_j^{2q+1}$ .  $\square$

The true message of Corollary 10.10 emerges if we bound the series using its largest term  $\sigma_{k+1}^{4q+2}$  and draw the factor  $\sigma_{k+1}$  out of the bracket:

$$\mathbb{E} \|(\mathbf{I} - \mathbf{P}_Z)\mathbf{A}\| \leq \left[ 1 + \sqrt{\frac{k}{p-1}} + \frac{e\sqrt{k+p}}{p} \cdot \sqrt{\min\{m, n\} - k} \right]^{1/(2q+1)} \sigma_{k+1}.$$

In words, as we increase the exponent  $q$ , the power scheme drives the extra factor in the error to one exponentially fast. By the time  $q \sim \log(\min\{m, n\})$ ,

$$\mathbb{E} \|(\mathbf{I} - \mathbf{P}_Z)\mathbf{A}\| \sim \sigma_{k+1},$$

which is the baseline for the spectral norm.

In most situations, the error bound given by Corollary 10.10 is substantially better than the estimates discussed in the last paragraph. For example, suppose that the tail singular values exhibit the decay profile

$$\sigma_j \lesssim j^{(1+\varepsilon)/(4q+2)} \quad \text{for } j > k \text{ and } \varepsilon > 0.$$

Then the series in Corollary 10.10 is comparable with its largest term, which allows us to remove the dimensional factor  $\min\{m, n\}$  from the error bound.

To obtain large deviation bounds for the performance of the power scheme, simply combine Theorem 9.2 with Theorem 10.8. We omit a detailed statement.

**REMARK 10.1.** We lack an analogous theory for the Frobenius norm because Theorem 9.2 depends on Proposition 8.6, which is not true for the Frobenius norm. It is possible to obtain some results by estimating the Frobenius norm in terms of the spectral norm.

**11. SRFT test matrices.** Another way to implement the proto-algorithm from §1.3 is to use a structured random matrix so that the matrix product in Step 2 can be performed quickly. One type of structured random matrix that has been proposed in the literature is the *subsampled random Fourier transform*, or SRFT, which we discussed in §4.6. In this section, we present bounds on the performance of the proto-algorithm when it is implemented with an SRFT test matrix. In contrast with the results for Gaussian test matrices, the results in this section hold for both real and complex input matrices.

**11.1. Construction and Properties.** Recall from §4.6 that an SRFT is a tall  $n \times \ell$  matrix of the form  $\mathbf{\Omega} = \sqrt{n/\ell} \cdot \mathbf{D}\mathbf{F}\mathbf{R}^*$  where

- $\mathbf{D}$  is a random  $n \times n$  diagonal matrix whose entries are independent and uniformly distributed on the complex unit circle;
- $\mathbf{F}$  is the  $n \times n$  unitary discrete Fourier transform; and
- $\mathbf{R}$  is a random  $\ell \times n$  matrix that restricts an  $n$ -dimensional vector to  $\ell$  coordinates, chosen uniformly at random.

Up to scaling, an SRFT is just a section of a unitary matrix, so it satisfies the norm identity  $\|\mathbf{\Omega}\| = \sqrt{n/\ell}$ . The critical fact is that an appropriately designed SRFT approximately preserves the geometry of an *entire subspace of vectors*.

**THEOREM 11.1** (The SRFT preserves geometry). *Fix an  $n \times k$  orthonormal matrix  $\mathbf{V}$ , and draw an  $n \times \ell$  SRFT matrix  $\mathbf{\Omega}$  where the parameter  $\ell$  satisfies*

$$4 \left[ \sqrt{k} + \sqrt{8 \log(kn)} \right]^2 \log(k) \leq \ell \leq n.$$

*Then*

$$0.40 \leq \sigma_k(\mathbf{V}^* \mathbf{\Omega}) \quad \text{and} \quad \sigma_1(\mathbf{V}^* \mathbf{\Omega}) \leq 1.48$$

*with failure probability at most  $O(k^{-1})$ .*

In words, the kernel of an SRFT of dimension  $\ell \sim k \log(k)$  is unlikely to intersect a fixed  $k$ -dimensional subspace. In contrast with the Gaussian case, the logarithmic factor  $\log(k)$  in the lower bound on  $\ell$  cannot generally be removed (Remark 11.2).

Theorem 11.1 follows from a straightforward variation of the argument in [134], which establishes equivalent bounds for a real analog of the SRFT, called the *subsampling randomized Hadamard transform* (SRHT). We omit further details.

**REMARK 11.1.** For large problems, we can obtain better numerical constants [134, Thm. 3.2]. Fix a small, positive number  $\iota$ . If  $k \gg \log(n)$ , then sampling

$$\ell \geq (1 + \iota) \cdot k \log(k)$$

coordinates is sufficient to ensure that  $\sigma_k(\mathbf{V}^* \mathbf{\Omega}) \geq \iota$  with failure probability at most  $O(k^{-c_\iota})$ . This sampling bound is essentially optimal because  $(1 - \iota) \cdot k \log(k)$  samples are not adequate in the worst case; see Remark 11.2.

**REMARK 11.2.** The logarithmic factor in Theorem 11.1 is *necessary* when the orthonormal matrix  $\mathbf{V}$  is particularly evil. Let us describe an infinite family of worst-case examples. Fix an integer  $k$ , and let  $n = k^2$ . Form an  $n \times k$  orthonormal matrix  $\mathbf{V}$  by regular decimation of the  $n \times n$  identity matrix. More precisely,  $\mathbf{V}$  is the matrix whose  $j$ th row has a unit entry in column  $(j - 1)/k$  when  $j \equiv 1 \pmod{k}$  and is zero otherwise. To see why this type of matrix is nasty, it is helpful to consider the auxiliary matrix  $\mathbf{W} = \mathbf{V}^* \mathbf{D} \mathbf{F}$ . Observe that, up to scaling and modulation of columns,  $\mathbf{W}$  consists of  $k$  copies of a  $k \times k$  DFT concatenated horizontally.

Suppose that we apply the SRFT  $\mathbf{\Omega} = \mathbf{D} \mathbf{F} \mathbf{R}^*$  to the matrix  $\mathbf{V}^*$ . We obtain a matrix of the form  $\mathbf{X} = \mathbf{V}^* \mathbf{\Omega} = \mathbf{W} \mathbf{R}^*$ , which consists of  $\ell$  random columns sampled from  $\mathbf{W}$ . Theorem 11.1 certainly cannot hold unless  $\sigma_k(\mathbf{X}) > 0$ . To ensure the latter event occurs, we must pick at least one copy each of the  $k$  distinct columns of  $\mathbf{W}$ . This is the coupon collector's problem [98, Sec. 3.6] in disguise. To obtain a complete set of  $k$  coupons (i.e., columns) with nonnegligible probability, we must draw at least  $k \log(k)$  columns. The fact that we are sampling without replacement does not improve the analysis appreciably because the matrix has too many columns.

**11.2. Performance guarantees.** We are now prepared to present detailed information on the performance of the proto-algorithm when the test matrix  $\mathbf{\Omega}$  is an SRFT.

**THEOREM 11.2** (Error bounds for SRFT). *Fix an  $m \times n$  matrix  $\mathbf{A}$  with singular values  $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots$ . Draw an  $n \times \ell$  SRFT matrix  $\mathbf{\Omega}$ , where*

$$4 \left[ \sqrt{k} + \sqrt{8 \log(kn)} \right]^2 \log(k) \leq \ell \leq n.$$

Construct the sample matrix  $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$ . Then

$$\begin{aligned} \|(\mathbf{I} - \mathbf{P}_{\mathbf{Y}})\mathbf{A}\| &\leq \sqrt{1 + 7n/\ell} \cdot \sigma_{k+1} \quad \text{and} \\ \|(\mathbf{I} - \mathbf{P}_{\mathbf{Y}})\mathbf{A}\|_{\text{F}} &\leq \sqrt{1 + 7n/\ell} \cdot \left( \sum_{j>k} \sigma_j^2 \right)^{1/2} \end{aligned}$$

with failure probability at most  $O(k^{-1})$ .

As we saw in §10.2, the quantity  $\sigma_{k+1}$  is the minimal spectral-norm error possible when approximating  $\mathbf{A}$  with a rank- $k$  matrix. Similarly, the series in the second bound is the minimal Frobenius-norm error when approximating  $\mathbf{A}$  with a rank- $k$  matrix. We see that both error bounds lie within a polynomial factor of the baseline, and this factor decreases with the number  $\ell$  of samples we retain.

The likelihood of error with an SRFT test matrix is substantially worse than in the Gaussian case. The failure probability here is roughly  $k^{-1}$ , while in the Gaussian case, the failure probability is roughly  $e^{-(\ell-k)}$ . This qualitative difference is not an artifact of the analysis; discrete sampling techniques inherently fail with higher probability.

Matrix approximation schemes based on SRFTs often perform much better in practice than the error analysis here would indicate. While it is not generally possible to guarantee accuracy with a sampling parameter less than  $\ell \sim k \log(k)$ , we have found empirically that the choice  $\ell = k + 20$  is adequate in almost all applications. Indeed, SRFTs sometimes perform even *better* than Gaussian matrices (see, e.g., Figure 7.8).

We complete the section with the proof of Theorem 11.2.

*Proof.* [Theorem 11.2] Let  $\mathbf{V}$  be the right unitary factor of matrix  $\mathbf{A}$ , and partition  $\mathbf{V} = [\mathbf{V}_1 \mid \mathbf{V}_2]$  into blocks containing, respectively,  $k$  and  $n - k$  columns. Recall that

$$\mathbf{\Omega}_1 = \mathbf{V}_1^* \mathbf{\Omega} \quad \text{and} \quad \mathbf{\Omega}_2 = \mathbf{V}_2^* \mathbf{\Omega}.$$

where  $\mathbf{\Omega}$  is the conjugate transpose of an SRFT. Theorem 11.1 ensures that the submatrix  $\mathbf{\Omega}_1$  has full row rank, with failure probability at most  $O(k^{-1})$ . Therefore, Theorem 9.1 implies that

$$\|(\mathbf{I} - \mathbf{P}_{\mathbf{Y}})\mathbf{A}\| \leq \|\mathbf{\Sigma}_2\| \left[ 1 + \|\mathbf{\Omega}_1^\dagger\|^2 \cdot \|\mathbf{\Omega}_2\|^2 \right]^{1/2},$$

where  $\|\cdot\|$  denotes either the spectral norm or the Frobenius norm. Our application of Theorem 11.1 also ensures that the spectral norm of  $\mathbf{\Omega}_1^\dagger$  is under control.

$$\|\mathbf{\Omega}_1^\dagger\|^2 \leq \frac{1}{0.40^2} < 7.$$

We may bound the spectral norm of  $\mathbf{\Omega}_2$  deterministically.

$$\|\mathbf{\Omega}_2\| = \|\mathbf{V}_2^* \mathbf{\Omega}\| \leq \|\mathbf{V}_2^*\| \|\mathbf{\Omega}\| = \sqrt{n/\ell}$$

since  $\mathbf{V}_2$  and  $\sqrt{\ell/n} \cdot \mathbf{\Omega}$  are both orthonormal matrices. Combine these estimates to complete the proof.  $\square$

**Acknowledgments.** The authors have benefited from valuable discussions with many researchers, among them Inderjit Dhillon, Petros Drineas, Ming Gu, Edo Liberty, Michael Mahoney, Vladimir Rokhlin, Yoel Shkolnisky, and Arthur Szlam. In

particular, we would like to thank Mark Tygert for his insightful remarks on early drafts of this paper. The example in Section 7.2 was provided by François Meyer of the University of Colorado at Boulder. The example in Section 7.3 comes from the FERET database of facial images collected under the FERET program, sponsored by the DoD Counterdrug Technology Development Program Office. The work reported was initiated during the program *Mathematics of Knowledge and Search Engines* held at IPAM in the fall of 2007. Finally, we would like to thank the anonymous referees, whose thoughtful remarks have helped us to improve the manuscript dramatically.

**Appendix A. On Gaussian matrices.** This appendix collects some of the properties of Gaussian matrices that we use in our analysis. Most of the results follow quickly from material that is already available in the literature. One fact, however, requires a surprisingly difficult new argument. We focus on the real case here; the complex case is similar but actually yields better results.

**A.1. Expectation of norms.** We begin with the expected Frobenius norm of a scaled Gaussian matrix, which follows from an easy calculation.

PROPOSITION A.1. *Fix real matrices  $\mathbf{S}, \mathbf{T}$ , and draw a standard Gaussian matrix  $\mathbf{G}$ . Then*

$$\left(\mathbb{E} \|\mathbf{SGT}\|_{\text{F}}^2\right)^{1/2} = \|\mathbf{S}\|_{\text{F}} \|\mathbf{T}\|_{\text{F}}.$$

*Proof.* The distribution of a Gaussian matrix is invariant under orthogonal transformations, and the Frobenius norm is also unitarily invariant. As a result, it represents no loss of generality to assume that  $\mathbf{S}$  and  $\mathbf{T}$  are diagonal. Therefore,

$$\mathbb{E} \|\mathbf{SGT}\|_{\text{F}}^2 = \mathbb{E} \left[ \sum_{jk} |s_{jj} g_{jk} t_{kk}|^2 \right] = \sum_{jk} |s_{jj}|^2 |t_{kk}|^2 = \|\mathbf{S}\|_{\text{F}}^2 \|\mathbf{T}\|_{\text{F}}^2.$$

Since the right-hand side is unitarily invariant, we have also identified the value of the expectation for general matrices  $\mathbf{S}$  and  $\mathbf{T}$ .  $\square$

The literature contains an excellent bound for the expected spectral norm of a scaled Gaussian matrix. The result is due to Gordon [62, 63], who established the bound using a sharp version of Slepian's lemma. See [83, §3.3] and [34, §2.3] for additional discussion.

PROPOSITION A.2. *Fix real matrices  $\mathbf{S}, \mathbf{T}$ , and draw a standard Gaussian matrix  $\mathbf{G}$ . Then*

$$\mathbb{E} \|\mathbf{SGT}\| \leq \|\mathbf{S}\| \|\mathbf{T}\|_{\text{F}} + \|\mathbf{S}\|_{\text{F}} \|\mathbf{T}\|.$$

**A.2. Spectral norm of pseudoinverse.** Now, we turn to the pseudoinverse of a Gaussian matrix. Recently, Chen and Dongarra developed a good bound on the probability that its spectral norm is large. The statement here follows from [25, Lem. 4.1] after an application of Stirling's approximation. See also [91, Lem. 2.14]

PROPOSITION A.3. *Let  $\mathbf{G}$  be an  $m \times n$  standard Gaussian matrix with  $n \geq m \geq 2$ . For each  $t > 0$ ,*

$$\mathbb{P} \{ \|\mathbf{G}^\dagger\| > t \} \leq \frac{1}{\sqrt{2\pi(n-m+1)}} \left[ \frac{e\sqrt{n}}{n-m+1} \right]^{n-m+1} t^{-(n-m+1)}.$$



We can use Proposition A.3 to bound the expected spectral norm of a pseudo-inverted Gaussian matrix.

PROPOSITION A.4. *Let  $\mathbf{G}$  be a  $m \times n$  standard Gaussian matrix with  $n - m \geq 1$  and  $m \geq 2$ . Then*

$$\mathbb{E} \|\mathbf{G}^\dagger\| < \frac{e\sqrt{n}}{n - m}$$

*Proof.* Let us make the abbreviations  $p = n - m$  and

$$C = \frac{1}{\sqrt{2\pi(p+1)}} \left[ \frac{e\sqrt{n}}{p+1} \right]^{p+1}.$$

We compute the expectation by way of a standard argument. The integral formula for the mean of a nonnegative random variable implies that, for all  $E > 0$ ,

$$\begin{aligned} \mathbb{E} \|\mathbf{G}^\dagger\| &= \int_0^\infty \mathbb{P} \{ \|\mathbf{G}^\dagger\| > t \} dt \leq E + \int_E^\infty \mathbb{P} \{ \|\mathbf{G}^\dagger\| > t \} dt \\ &\leq E + C \int_E^\infty t^{-(p+1)} dt = E + \frac{1}{p} C E^{-p}, \end{aligned}$$

where the second inequality follows from Proposition A.3. The right-hand side is minimized when  $E = C^{1/(p+1)}$ . Substitute and simplify.  $\square$

**A.3. Frobenius norm of pseudoinverse.** The squared Frobenius norm of a pseudo-inverted Gaussian matrix is closely connected with the trace of an inverted Wishart matrix. This observation leads to an exact expression for the expectation.

PROPOSITION A.5. *Let  $\mathbf{G}$  be an  $m \times n$  standard Gaussian matrix with  $n - m \geq 2$ . Then*

$$\mathbb{E} \|\mathbf{G}^\dagger\|_F^2 = \frac{m}{n - m - 1}.$$

*Proof.* Observe that

$$\|\mathbf{G}^\dagger\|_F^2 = \text{trace} [(\mathbf{G}^\dagger)^* \mathbf{G}^\dagger] = \text{trace} [(\mathbf{G}\mathbf{G}^*)^{-1}].$$

The second identity holds almost surely because the Wishart matrix  $\mathbf{G}\mathbf{G}^*$  is invertible with probability one. The random matrix  $(\mathbf{G}\mathbf{G}^*)^{-1}$  follows the inverted Wishart distribution, so we can compute its expected trace explicitly using a formula from [99, p. 97].  $\square$

On the other hand, very little seems to be known about the tail behavior of the Frobenius norm of a pseudo-inverted Gaussian matrix. The following theorem, which is new, provides an adequate bound on the probability of a large deviation.

THEOREM A.6. *Let  $\mathbf{G}$  be an  $m \times n$  standard Gaussian matrix with  $n - m \geq 4$ . For each  $t \geq 1$ ,*

$$\mathbb{P} \left\{ \|\mathbf{G}^\dagger\|_F^2 > \frac{12m}{n - m} \cdot t \right\} \leq 4t^{-(n-m)/2}.$$

Neither the precise form of Theorem A.6 nor the constants are ideal; we have focused instead on establishing a useful bound with minimal fuss. The rest of the section is devoted to the rather lengthy proof. Unfortunately, most of the standard methods for producing tail bounds fail for random variables that do not exhibit normal or exponential concentration. Our argument relies on special properties of Gaussian matrices and a dose of brute force.

**A.3.1. Technical background.** We begin with a piece of notation. For any number  $q \geq 1$ , we define the  $L_q$  norm of a random variable  $Z$  by

$$\mathbb{E}^q(Z) = (\mathbb{E}|Z|^q)^{1/q}.$$

In particular, the  $L_q$  norm satisfies the triangle inequality.

We continue with a collection of technical results. First, we present a striking fact about the structure of Gaussian matrices [55, §3.5].

**PROPOSITION A.7.** *For  $n \geq m$ , an  $m \times n$  standard Gaussian matrix is orthogonally equivalent with a random bidiagonal matrix*

$$\mathbf{L} = \begin{bmatrix} X_n & & & & \\ Y_{m-1} & X_{n-1} & & & \\ & Y_{m-2} & X_{n-2} & & \\ & & \ddots & \ddots & \\ & & & Y_1 & X_{n-(m-1)} \end{bmatrix}_{m \times n}, \quad (\text{A.1})$$

where, for each  $j$ , the random variables  $X_j^2$  and  $Y_j^2$  follow the  $\chi^2$  distribution with  $j$  degrees of freedom. Furthermore, these variates are mutually independent.

We also require the moments of a chi-square variate, which are expressed in terms of special functions.

**PROPOSITION A.8.** *Let  $\Xi$  be a  $\chi^2$  variate with  $k$  degrees of freedom. When  $0 \leq q < k/2$ ,*

$$\mathbb{E}(\Xi^q) = \frac{2^q \Gamma(k/2 + q)}{\Gamma(k/2)} \quad \text{and} \quad \mathbb{E}(\Xi^{-q}) = \frac{\Gamma(k/2 - q)}{2^q \Gamma(k/2)}.$$

*Proof.* Recall that a  $\chi^2$  variate with  $k$  degrees of freedom has the probability density function

$$f(t) = \frac{1}{2^{k/2} \Gamma(k/2)} t^{k/2-1} e^{-t/2}, \quad \text{for } t \geq 0.$$

By the integral formula for expectation,

$$\mathbb{E}(\Xi^q) = \int_0^\infty t^q f(t) dt = \frac{2^q \Gamma(k/2 + q)}{\Gamma(k/2)},$$

where the second equality follows from Euler's integral expression for the gamma function. The other calculation is similar.  $\square$

To streamline the proof, we eliminate the gamma functions from Proposition A.8. The next result bounds the positive moments of a chi-square variate.

LEMMA A.9. *Let  $\Xi$  be a  $\chi^2$  variate with  $k$  degrees of freedom. For  $q \geq 1$ ,*

$$\mathbb{E}^q(\Xi) \leq k + q.$$

*Proof.* Write  $q = r + \theta$ , where  $r = \lfloor q \rfloor$ . Repeated application of the functional equation  $z\Gamma(z) = \Gamma(z+1)$  yields

$$\mathbb{E}^q(\Xi) = \left[ \frac{2^\theta \Gamma(k/2 + \theta)}{\Gamma(k/2)} \cdot \prod_{j=1}^r (k + 2(q - j)) \right]^{1/q}.$$

The gamma function is logarithmically convex, so

$$\frac{2^\theta \Gamma(k/2 + \theta)}{\Gamma(k/2)} \leq \frac{2^\theta \cdot \Gamma(k/2)^{1-\theta} \cdot \Gamma(k/2 + 1)^\theta}{\Gamma(k/2)} = k^\theta \leq \left[ \prod_{j=1}^r (k + 2(q - j)) \right]^{\theta/r}.$$

The second inequality holds because  $k$  is smaller than each term in the product, hence is smaller than their geometric mean. As a consequence,

$$\mathbb{E}^q(\Xi) \leq \left[ \prod_{j=1}^r (k + 2(q - j)) \right]^{1/r} \leq \frac{1}{r} \sum_{j=1}^r (k + 2(q - j)) \leq k + q.$$

The second relation is the inequality between the geometric and arithmetic mean.  $\square$

Finally, we develop a bound for the negative moments of a chi-square variate.

LEMMA A.10. *Let  $\Xi$  be a  $\chi^2$  variate with  $k$  degrees of freedom, where  $k \geq 5$ . When  $2 \leq q \leq (k-1)/2$ ,*

$$\mathbb{E}^q(\Xi^{-1}) < \frac{3}{k}.$$

*Proof.* We establish the bound for  $q = (k-1)/2$ . For smaller values of  $q$ , the result follows from Hölder's inequality. Proposition A.8 shows that

$$\mathbb{E}^q(\Xi^{-1}) = \left[ \frac{\Gamma(k/2 - q)}{2^q \Gamma(k/2)} \right]^{1/q} = \left[ \frac{\Gamma(1/2)}{2^q \Gamma(k/2)} \right]^{1/q}.$$

Stirling's approximation ensures that  $\Gamma(k/2) \geq \sqrt{2\pi} \cdot (k/2)^{(k-1)/2} \cdot e^{-k/2}$ . Since the value  $\Gamma(1/2) = \sqrt{\pi}$ ,

$$\mathbb{E}^q(\Xi^{-1}) \leq \left[ \frac{\sqrt{\pi}}{2^q \sqrt{2\pi} \cdot (k/2)^q \cdot e^{-q-1/2}} \right]^{1/q} = \frac{e}{k} \left[ \frac{e}{2} \right]^{1/2q} < \frac{3}{k},$$

where we used the assumption  $q \geq 2$  to complete the numerical estimate.  $\square$

**A.3.2. Proof of Theorem A.6.** Let  $\mathbf{G}$  be an  $m \times n$  Gaussian matrix, where we assume that  $n - m \geq 4$ . Define the random variable

$$Z = \|\mathbf{G}^\dagger\|_{\text{F}}^2.$$

Our goal is to develop a tail bound for  $Z$ . The argument is inspired by work of Szarek [130, §6] for square Gaussian matrices.

The first step is to find an explicit, tractable representation for the random variable. According to Proposition A.7, a Gaussian matrix  $\mathbf{G}$  is orthogonally equivalent with a bidiagonal matrix  $\mathbf{L}$  of the form (A.1). Making an analogy with the inversion formula for a triangular matrix, we realize that the pseudoinverse of  $\mathbf{L}$  is given by

$$\mathbf{L}^\dagger = \begin{bmatrix} X_n^{-1} & & & & \\ \frac{-Y_{m-1}}{X_n X_{n-1}} & X_{n-1}^{-1} & & & \\ & \frac{-Y_{m-2}}{X_{n-1} X_{n-2}} & X_{n-2}^{-1} & & \\ & & \ddots & \ddots & \\ & & & \frac{-Y_1}{X_{n-(m-2)} X_{n-(m-1)}} & X_{n-(m-1)}^{-1} \end{bmatrix}_{n \times m}.$$

Because  $\mathbf{L}^\dagger$  is orthogonally equivalent with  $\mathbf{G}^\dagger$  and the Frobenius norm is unitarily invariant, we have the relations

$$Z = \|\mathbf{G}^\dagger\|_{\text{F}}^2 = \|\mathbf{L}^\dagger\|_{\text{F}}^2 \leq \sum_{j=0}^{m-1} \frac{1}{X_{n-j}^2} \left( 1 + \frac{Y_{m-j}^2}{X_{n-j+1}^2} \right),$$

where we have added an extra subdiagonal term (corresponding with  $j = 0$ ) so that we can avoid exceptional cases later. We abbreviate the summands as

$$W_j = \frac{1}{X_{n-j}^2} \left( 1 + \frac{Y_{m-j}^2}{X_{n-j+1}^2} \right), \quad j = 0, 1, 2, \dots, m-1.$$

Next, we develop a large deviation bound for each summand by computing a moment and invoking Markov's inequality. For the exponent  $q = (n-m)/2$ , Lemmas A.9 and A.10 yield

$$\begin{aligned} \mathbb{E}^q(W_j) &= \mathbb{E}^q(X_{n-j}^{-2}) \cdot \mathbb{E}^q \left[ 1 + \frac{Y_{m-j}^2}{X_{n-j+1}^2} \right] \\ &\leq \mathbb{E}^q(X_{n-j}^{-2}) \left[ 1 + \mathbb{E}^q(Y_{m-j}^2) \cdot \mathbb{E}^q(X_{n-j+1}^{-2}) \right] \\ &\leq \frac{3}{n-j} \left[ 1 + \frac{3(m-j+q)}{n-j+1} \right] \\ &= \frac{3}{n-j} \left[ 1 + 3 - \frac{3(n-m+1-q)}{n-j+1} \right] \end{aligned}$$

Note that the first two relations require the independence of the variates and the triangle inequality for the  $L_q$  norm. The maximum value of the bracket evidently occurs when  $j = 0$ , so

$$\mathbb{E}^q(W_j) < \frac{12}{n-j}, \quad j = 0, 1, 2, \dots, m-1.$$

Markov's inequality results in

$$\mathbb{P} \left\{ W_j \geq \frac{12}{n-j} \cdot u \right\} \leq u^{-q}.$$

Select  $u = t \cdot (n-j)/(n-m)$  to reach

$$\mathbb{P} \left\{ W_j \geq \frac{12}{n-m} \cdot t \right\} \leq \left[ \frac{n-m}{n-j} \right]^q t^{-q}.$$

To complete the argument, we combine these estimates by means of the union bound and clean up the resulting mess. Since  $Z \leq \sum_{j=0}^{m-1} W_j$ ,

$$\mathbb{P} \left\{ Z \geq \frac{12m}{n-m} \cdot t \right\} \leq t^{-q} \sum_{j=0}^{m-1} \left[ \frac{n-m}{n-j} \right]^q.$$

To control the sum on the right-hand side, observe that

$$\begin{aligned} \sum_{j=0}^{m-1} \left[ \frac{n-m}{n-j} \right]^q &< (n-m)^q \int_0^m (n-x)^{-q} dx \\ &< \frac{(n-m)^q}{q-1} (n-m)^{-q+1} = \frac{2(n-m)}{n-m-2} \leq 4, \end{aligned}$$

where the last inequality follows from the hypothesis  $n-m \geq 4$ . Together, the estimates in this paragraph produce the advertised bound.

REMARK A.1. It would be very interesting to find more conceptual and extensible argument that yields accurate concentration results for inverse spectral functions of a Gaussian matrix.

#### REFERENCES

- [1] D. ACHLIOPTAS, *Database-friendly random projections: Johnson–Lindenstrauss with binary coins*, J. Comput. System Sci., 66 (2003), pp. 671–687.
- [2] D. ACHLIOPTAS AND F. MCSHERRY, *Fast computation of low-rank matrix approximations*, J. Assoc. Comput. Mach., 54 (2007), pp. Art. 9, 19 pp. (electronic).
- [3] N. AILON AND B. CHAZELLE, *Approximate nearest neighbors and the fast Johnson–Lindenstrauss transform*, in STOC '06: Proc. 38th Ann. ACM Symp. Theory of Computing, 2006, pp. 557–563.
- [4] N. AILON AND E. LIBERTY, *Fast dimension reduction using Rademacher series on dual BCH codes*, in STOC '08: Proc. 40th Ann. ACM Symp. Theory of Computing, 2008.
- [5] N. ALON, P. GIBBONS, Y. MATIAS, AND M. SZEGEDY, *Tracking join and self-join sizes in limited storage*, in Proc. 18th ACM Symp. Principles of Database Systems (PODS), 1999, pp. 10–20.
- [6] N. ALON, Y. MATIAS, AND M. SZEGEDY, *The space complexity of approximating frequency moments*, in STOC '96: Proc. 28th Ann. ACM Symp. Theory of Algorithms, 1996, pp. 20–29.
- [7] S. ARORA, E. HAZAN, AND S. KALE, *A fast random sampling algorithm for sparsifying matrices*, in Approximation, randomization and combinatorial optimization: Algorithms and Techniques, Springer, Berlin, 2006, pp. 272–279.
- [8] A. R. BARRON, *Universal approximation bounds for superpositions of a sigmoidal function*, IEEE Trans. Inform. Theory, 39 (1993), pp. 930–945.
- [9] I. BEICHL, *The Metropolis algorithm*, Comput. Sci. Eng., (2000), pp. 65–69.
- [10] M.-A. BELLABAS AND P. J. WOLFE, *On sparse representations of linear operators and the approximation of matrix products*, in Proc. 42nd Ann. Conf. Information Sciences and Systems (CISS), 2008, pp. 258–263.

- [11] R. BHATIA, *Matrix Analysis*, no. 169 in GTM, Springer, Berlin, 1997.
- [12] Å. BJÖRCK, *Numerics of Gram–Schmidt orthogonalization*, Linear Algebra Appl., 197–198 (1994), pp. 297–316.
- [13] ———, *Numerical Methods for Least Squares Problems*, SIAM, Philadelphia, PA, 1996.
- [14] V. BOGDANOV, *Gaussian Measures*, American Mathematical Society, Providence, RI, 1998.
- [15] J. BOURGAIN, *On Lipschitz embedding of finite metric spaces in Hilbert space*, Israel J. Math., 52 (1985), pp. 46–52.
- [16] C. BOUTSIDIS AND P. DRINEAS, *Random projections for nonnegative least squares*, Linear Algebra Appl., 431 (2009), pp. 760–771.
- [17] C. BOUTSIDIS, P. DRINEAS, AND M. W. MAHONEY, *An improved approximation algorithm for the column subset selection problem*, in Proc. 20th Ann. ACM–SIAM Symp. Discrete Algorithms (SODA), 2009.
- [18] C. BOUTSIDIS, M. W. MAHONEY, AND P. DRINEAS, *Unsupervised feature selection for principal components analysis*, in Proc. ACM SIGKDD Intl. Conf. Knowledge Discovery and Data Mining (KDD), Aug. 2008.
- [19] E. CANDÈS AND J. K. ROMBERG, *Sparsity and incoherence in compressive sampling*, Inverse Problems, 23 (2007), pp. 969–985.
- [20] E. CANDÈS, J. K. ROMBERG, AND T. TAO, *Robust uncertainty principles: Exact signal reconstruction from highly incomplete Fourier information*, IEEE Trans. Inform. Theory, 52 (2006), pp. 489–509.
- [21] E. J. CANDÈS, *Compressive sampling*, in Proc. 2006 Intl. Cong. Mathematicians, Madrid, 2006.
- [22] E. J. CANDÈS AND B. RECHT, *Exact matrix completion via convex optimization*, Found. Comput. Math., 9 (2009), pp. 717–772.
- [23] E. J. CANDÈS AND T. TAO, *The power of convex relaxation: Near-optimal matrix completion*, IEEE Trans. Inform. Theory, 56 (2010), pp. 2053–2080.
- [24] B. CARL, *Inequalities of Bernstein–Jackson-type and the degree of compactness in Banach spaces*, Ann. Inst. Fourier (Grenoble), 35 (1985), pp. 79–118.
- [25] Z. CHEN AND J. J. DONGARRA, *Condition numbers of Gaussian random matrices*, SIAM J. Matrix Anal. Appl., 27 (2005), pp. 603–620.
- [26] H. CHENG, Z. GIMBUTAS, P.-G. MARTINSSON, AND V. ROKHLIN, *On the compression of low rank matrices*, SIAM J. Sci. Comput., 26 (2005), pp. 1389–1404 (electronic).
- [27] 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.
- [28] K. L. CLARKSON, *Subgradient and sampling algorithms for  $\ell_1$  regression*, in Proc. 16th Ann. ACM–SIAM Symp. Discrete Algorithms (SODA), 2005, pp. 257–266.
- [29] K. L. CLARKSON AND D. P. WOODRUFF, *Numerical linear algebra in the streaming model*, in STOC ’09: Proc. 41st Ann. ACM Symp. Theory of Computing, 2009.
- [30] R. R. COIFMAN, S. LAFON, A. B. LEE, M. MAGGIONI, B. NADLER, F. WARNER, AND S. W. ZUCKER, *Geometric diffusions as a tool for harmonic analysis and structure definition of data: Diffusion maps*, Proc. Natl. Acad. Sci. USA, 102 (2005), pp. 7426–7431.
- [31] A. DASGUPTA, P. DRINEAS, B. HARB, R. KUMAR, AND M. W. MAHONEY, *Sampling algorithms and coresets for  $\ell_p$  regression*, SIAM J. Comput., 38 (2009), pp. 2060–2078.
- [32] S. DASGUPTA AND A. GUPTA, *An elementary proof of the Johnson–Lindenstrauss lemma*, Computer Science Dept. Tech. Report 99-006, Univ. California at Berkeley, Mar. 1999.
- [33] A. D’ASPREMONT, *Subsampling algorithms for semidefinite programming*. Available at [arXiv:0803.1990](https://arxiv.org/abs/0803.1990), Apr. 2009.
- [34] K. R. DAVIDSON AND S. J. SZAREK, *Local operator theory, random matrices, and Banach spaces*, in Handbook of Banach Space Geometry, W. B. Johnson and J. Lindenstrauss, eds., Elsevier, 2002, pp. 317–366.
- [35] J. DEMMEL, I. DUMITRIU, AND O. HOLTZ, *Fast linear algebra is stable*, Numer. Math., 108 (2007), pp. 59–91.
- [36] A. DESHPANDE AND L. RADEMACHER, *Efficient volume sampling for row/column subset selection*. Available at [arXiv:1004.4057](https://arxiv.org/abs/1004.4057), Apr. 2010.
- [37] A. DESHPANDE, L. RADEMACHER, S. VEMPALA, AND G. WANG, *Matrix approximation and projective clustering via volume sampling*, in Proc. 17th Ann. ACM–SIAM Symp. Discrete Algorithms (SODA), 2006, pp. 1117–1126.
- [38] A. DESHPANDE AND S. VEMPALA, *Adaptive sampling and fast low-rank matrix approximation*, in Approximation, randomization and combinatorial optimization, vol. 4110 of LNCS, Springer, Berlin, 2006, pp. 292–303.
- [39] J. D. DIXON, *Estimating extremal eigenvalues and condition numbers of matrices*, SIAM J. Numer. Anal., 20 (1983), pp. 812–814.

- [40] J. DONGARRA AND F. SULLIVAN, *The top 10 algorithms*, Comput. Sci. Eng., (2000), pp. 22–23.
- [41] D. L. DONOHO, *Compressed sensing*, IEEE Trans. Inform. Theory, 52 (2006), pp. 1289–1306.
- [42] D. L. DONOHO, M. VETTERLI, R. A. DEVORE, AND I. DAUBECHIES, *Data compression and harmonic analysis*, IEEE Trans. Inform. Theory, 44 (1998), pp. 2433–2452.
- [43] P. DRINEAS, A. FRIEZE, R. KANNAN, S. VEMPALA, AND V. VINAY, *Clustering of large graphs via the singular value decomposition*, Machine Learning, 56 (2004), pp. 9–33.
- [44] P. DRINEAS, A. FRIEZE, R. KANNAN, S. VEMPALA, AND V. VINAY, *Clustering in large graphs and matrices*, in Proc. 10th Ann. ACM Symp. on Discrete Algorithms (SODA), 1999, pp. 291–299.
- [45] P. DRINEAS, R. KANNAN, AND M. W. MAHONEY, *Fast Monte Carlo algorithms for matrices. I. Approximating matrix multiplication*, SIAM J. Comput., 36 (2006), pp. 132–157.
- [46] ———, *Fast Monte Carlo algorithms for matrices. II. Computing a low-rank approximation to a matrix*, SIAM J. Comput., 36 (2006), pp. 158–183 (electronic).
- [47] ———, *Fast Monte Carlo algorithms for matrices. III. Computing a compressed approximate matrix decomposition*, SIAM J. Comput., 36 (2006), pp. 184–206.
- [48] P. DRINEAS AND M. W. MAHONEY, *On the Nystrom method for approximating a Gram matrix for improved kernel-based learning*, J. Mach. Learn. Res., 6 (2005), pp. 2153–2175.
- [49] ———, *A randomized algorithm for a tensor-based generalization of the singular value decomposition*, Linear Algebra Appl., 420 (2007), pp. 553–571.
- [50] P. DRINEAS, M. W. MAHONEY, AND S. MUTHUKRISHNAN, *Subspace sampling and relative-error matrix approximation: Column-based methods*, in APPROX and RANDOM 2006, J. Diaz et al., ed., no. 4110 in LNCS, Springer, Berlin, 2006, pp. 321–326.
- [51] ———, *Relative-error CUR matrix decompositions*, SIAM J. Matrix Anal. Appl., 30 (2008), pp. 844–881.
- [52] P. DRINEAS, M. W. MAHONEY, S. MUTHUKRISHNAN, AND T. SARLÓS, *Faster least squares approximation*, Num. Math., (2009). To appear. Available at [arXiv:0710.1435](https://arxiv.org/abs/0710.1435).
- [53] A. DVORETSKY, *Some results on convex bodies and Banach spaces*, in Proc. Intl. Symp. Linear Spaces, Jerusalem, 1961, pp. 123–160.
- [54] C. ECKART AND G. YOUNG, *The approximation of one matrix by another of lower rank*, Psychometrika, 1 (1936), pp. 211–218.
- [55] A. EDELMAN, *Eigenvalues and condition numbers of random matrices*, Ph.D. thesis, Mathematics Dept., Massachusetts Inst. Tech., Boston, MA, May 1989.
- [56] B. ENGQUIST AND O. RUNBORG, *Wavelet-based numerical homogenization with applications*, in Multiscale and Multiresolution Methods: Theory and Applications, T. J. Barth et al., ed., vol. 20 of LNCSE, Springer, Berlin, 2001, pp. 97–148.
- [57] A. FRIEZE, R. KANNAN, AND S. VEMPALA, *Fast Monte Carlo algorithms for finding low-rank approximations*, in Proc. 39th Ann. IEEE Symp. Foundations of Computer Science (FOCS), 1998, pp. 370–378.
- [58] ———, *Fast Monte Carlo algorithms for finding low-rank approximations*, J. Assoc. Comput. Mach., 51 (2004), pp. 1025–1041. (electronic).
- [59] A. YU. GARNAEV AND E. D. GLUSKIN, *The widths of a Euclidean ball*, Dokl. Akad. Nauk. SSSR, 277 (1984), pp. 1048–1052. In Russian.
- [60] A. GITTENS AND J. A. TROPP, *Error bounds for random matrix approximation schemes*. Available at [arXiv:0911.4108](https://arxiv.org/abs/0911.4108).
- [61] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, Johns Hopkins Studies in the Mathematical Sciences, Johns Hopkins Univ. Press, Baltimore, MD, 3rd ed., 1996.
- [62] Y. GORDON, *Some inequalities for Gaussian processes and applications*, Israel J. Math., 50 (1985), pp. 265–289.
- [63] ———, *Gaussian processes and almost spherical sections of convex bodies*, Ann. Probab., 16 (1988), pp. 180–188.
- [64] S. A. GOREINOV, E. E. TYRTYSHNIKOV, AND N. L. ZAMARASHKIN, *Theory of pseudo-skeleton matrix approximations*, Linear Algebra Appl., 261 (1997), pp. 1–21.
- [65] L. GRASEDYCK AND W. HACKBUSCH, *Construction and arithmetics of  $\mathcal{H}$ -matrices*, Computing, 70 (2003), pp. 295–334.
- [66] L. GREENGARD AND V. ROKHLIN, *A new version of the fast multipole method for the Laplace equation in three dimensions*, Acta Numer., 17 (1997), pp. 229–269.
- [67] M. GU, 2007. Personal communication.
- [68] 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.
- [69] N. HALKO, P.-G. MARTINSSON, Y. SHKOLNISKY, AND M. TYGERT, *An algorithm for the principal component analysis of large data sets*, 2010.
- [70] S. HAR-PELED, *Matrix approximation in linear time*. Manuscript. Available at <http://valis>.



- cs.uiuc.edu/~sariel/research/papers/05/lrank/, 2006.
- [71] T. HASTIE, R. TIBSHIRANI, AND J. FRIEDMAN, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, Springer, Berlin, 2nd ed., 2008.
  - [72] R. A. HORN AND C. R. JOHNSON, *Matrix Analysis*, Cambridge Univ. Press, Cambridge, 1985.
  - [73] P. INDYK AND R. MOTWANI, *Approximate nearest neighbors: Toward removing the curse of dimensionality*, in STOC '98: Proc. 30th Ann. ACM Symp. Theory of Computing, 1998, pp. 604–613.
  - [74] W. B. JOHNSON AND J. LINDENSTRAUSS, *Extensions of Lipschitz mappings into a Hilbert space*, Contemp. Math., 26 (1984), pp. 189–206.
  - [75] D. R. KARGER, *Random sampling in cut, flow, and network design problems*, Math. Oper. Res., 24 (1999), pp. 383–413.
  - [76] ———, *Minimum cuts in near-linear time*, J. ACM, 47 (2000), pp. 46–76.
  - [77] B. S. KAŠIN, *On the widths of certain finite-dimensional sets and classes of smooth functions*, Izv. Akad. Nauk. SSSR Ser. Mat., 41 (1977), pp. 334–351, 478. In Russian.
  - [78] J. KLEINBERG, *Two algorithms for nearest neighbor search in high dimensions*, in STOC '97: Proc. 29th ACM Symp. Theory of Computing, 1997, pp. 599–608.
  - [79] J. KUCZYŃSKI AND H. WOŹNIAKOWSKI, *Estimating the largest eigenvalue by the power and Lanczos algorithms with a random start*, SIAM J. Matrix Anal. Appl., 13 (1992), pp. 1094–1122.
  - [80] E. KUSHILEVITZ, R. OSTROVSKI, AND Y. RABANI, *Efficient search for approximate nearest neighbor in high-dimensional spaces*, SIAM J. Comput., 30 (2000), pp. 457–474.
  - [81] D. LE AND D. S. PARKER, *Using randomization to make recursive matrix algorithms practical*, J. Funct. Programming, 9 (1999), pp. 605–624.
  - [82] M. LEDOUX, *The Concentration of Measure Phenomenon*, no. 89 in MSM, American Mathematical Society, Providence, RI, 2001.
  - [83] M. LEDOUX AND M. TALAGRAND, *Probability in Banach Spaces: Isoperimetry and Processes*, Springer, Berlin, 1991.
  - [84] W. S. LEE, P. L. BARTLETT, AND R. C. WILLIAMSON, *Efficient agnostic learning of neural networks with bounded fan-in*, IEEE Trans. Inform. Theory, 42 (1996), pp. 2118–2132.
  - [85] Z. LEYK AND H. WOŹNIAKOWSKI, *Estimating the largest eigenvector by Lanczos and polynomial algorithms with a random start*, Num. Linear Algebra Appl., 5 (1998), pp. 147–164.
  - [86] E. LIBERTY, *Accelerated dense random projections*, Ph.D. thesis, Computer Science Dept., Yale University, New Haven, CT, 2009.
  - [87] E. LIBERTY, N. AILON, AND A. SINGER, *Dense fast random projections and lean Walsh transforms*, in APPROX and RANDOM 2008, A. Goel et al., ed., no. 5171 in LNCS, Springer, Berlin, 2008, pp. 512–522.
  - [88] E. LIBERTY, F. F. WOOLFE, P.-G. MARTINSSON, V. ROKHLIN, AND M. TYGERT, *Randomized algorithms for the low-rank approximation of matrices*, Proc. Natl. Acad. Sci. USA, 104 (2007), pp. 20167–20172.
  - [89] M. W. MAHONEY AND P. DRINEAS, *CUR matrix decompositions for improved data analysis*, Proc. Natl. Acad. Sci. USA, 106 (2009), pp. 697–702.
  - [90] P.-G. MARTINSSON, V. ROKHLIN, Y. SHKOLNISKY, AND M. TYGERT, *ID: A software package for low-rank approximation of matrices via interpolative decompositions, version 0.2*, 2008.
  - [91] P.-G. MARTINSSON, V. ROKHLIN, AND M. TYGERT, *A randomized algorithm for the approximation of matrices*, Computer Science Dept. Tech. Report 1361, Yale Univ., New Haven, CT, 2006.
  - [92] P.-G. MARTINSSON, A. SZLAM, AND M. TYGERT, *Normalized power iterations for the computation of SVD*. Manuscript., Nov. 2010.
  - [93] J. MATOUŠEK, *Lectures on Discrete Geometry*, Springer, Berlin, 2002.
  - [94] F. MCSHERRY, *Spectral Methods in Data Analysis*, Ph.D. thesis, Computer Science Dept., Univ. Washington, Seattle, WA, 2004.
  - [95] N. METROPOLIS AND S. ULAM, *The Monte Carlo method*, J. Amer. Statist. Assoc., 44 (1949), pp. 335–341.
  - [96] V. D. MILMAN, *A new proof of A. Dvoretzky's theorem on cross-sections of convex bodies*, Funkcional. Anal. i Priložen, 5 (1971), pp. 28–37.
  - [97] L. MIRSKY, *Symmetric gauge functions and unitarily invariant norms*, Quart. J. Math. Oxford Ser. (2), 11 (1960), pp. 50–59.
  - [98] R. MOTWANI AND P. RAGHAVAN, *Randomized Algorithms*, Cambridge Univ. Press, Cambridge, 1995.
  - [99] R. J. MUIRHEAD, *Aspects of Multivariate Statistical Theory*, Wiley, New York, NY, 1982.
  - [100] S. MUTHUKRISHNAN, *Data Streams: Algorithms and Applications*, Now Publ., Boston, MA,

- 2005.
- [101] D. NEEDELL, *Randomized Kaczmarz solver for noisy linear systems*, BIT, 50 (2010), pp. 395–403.
  - [102] N. H. NGUYEN, T. T. DO, AND T. D. TRAN, *A fast and efficient algorithm for low-rank approximation of a matrix*, in STOC '09: Proc. 41st Ann. ACM Symp. Theory of Computing, 2009.
  - [103] C.-T. PAN, *On the existence and computation of rank-revealing LU factorizations*, Linear Algebra Appl., 316 (2000), pp. 199–222.
  - [104] C. H. PAPADIMITRIOU, P. RAGHAVAN, H. TAMAKI, AND S. VEMPALA, *Latent semantic indexing: A probabilistic analysis*, in Proc. 17th ACM Symp. Principles of Database Systems (PODS), 1998.
  - [105] ———, *Latent semantic indexing: A probabilistic analysis*, J. Comput. System Sci., 61 (2000), pp. 217–235.
  - [106] D. S. PARKER AND B. PIERCE, *The randomizing FFT: An alternative to pivoting in Gaussian elimination*, Computer Science Dept. Tech. Report CSD 950037, Univ. California at Los Angeles, 1995.
  - [107] P. J. PHILLIPS, H. MOON, S.A. RIZVI, AND P.J. RAUSS, *The FERET evaluation methodology for face recognition algorithms*, IEEE Trans. Pattern Anal. Mach. Intelligence, 22 (2000), pp. 1090–1104.
  - [108] P. J. PHILLIPS, H. WECHSLER, J. HUANG, AND P. RAUSS, *The FERET database and evaluation procedure for face recognition algorithms*, Image Vision Comput., 16 (1998), pp. 295–306.
  - [109] W. H. PRESS, S. A. TEUKOLSKY, W. T. VETTERLING, AND B. P. FLANNERY, *Numerical Recipes: The Art of Scientific Computing*, Cambridge University Press, Cambridge, 3rd ed., 2007.
  - [110] A. RAHIMI AND B. RECHT, *Random features for large-scale kernel machines*, in Proc. 21st Ann. Conf. Advances in Neural Information Processing Systems (NIPS), 2007.
  - [111] B. RECHT, M. FAZEL, AND P. PARILLO, *Guaranteed minimum-rank solutions of matrix equations via nuclear-norm minimization*, SIAM Rev., 52 (2010), pp. 471–501.
  - [112] V. ROKHLIN, A. SZLAM, AND M. TYGERT, *A randomized algorithm for principal component analysis*, SIAM J. Matrix Anal. Appl., 31 (2009), pp. 1100–1124.
  - [113] V. ROKHLIN AND M. TYGERT, *A fast randomized algorithm for overdetermined linear least-squares regression*, Proc. Natl. Acad. Sci. USA, 105 (2008), pp. 13212–13217.
  - [114] S. ROWEIS, *EM algorithms for PCA and SPCA*, in Proc. 10th Ann. Conf. Advances in Neural Information Processing Systems (NIPS), MIT Press, 1997, pp. 626–632.
  - [115] M. RUDELSON, *Random vectors in the isotropic position*, J. Funct. Anal., 164 (1999), pp. 60–72.
  - [116] M. RUDELSON AND R. VERSHYNIN, *Sampling from large matrices: An approach through geometric functional analysis*, J. Assoc. Comput. Mach., 54 (2007), pp. Art. 21, 19 pp. (electronic).
  - [117] A. F. RUSTON, *Auerbach's theorem*, Math. Proc. Cambridge Philos. Soc., 56 (1964), pp. 476–480.
  - [118] T. SARLÓS, *Improved approximation algorithms for large matrices via random projections*, in Proc. 47th Ann. IEEE Symp. Foundations of Computer Science (FOCS), 2006, pp. 143–152.
  - [119] S. SHALEV-SHWARTZ AND N. SREBRO, *Low  $\ell_1$ -norm and guarantees on sparsifiability*, in ICML/COLT/UAI Sparse Optimization and Variable Selection Workshop, July 2008.
  - [120] X. SHEN AND F. G. MEYER, *Low-dimensional embedding of fMRI datasets*, Neuroimage, 41 (2008), pp. 886–902.
  - [121] N. D. SHYAMALKUMAR AND K. VARADARAJARAN, *Efficient subspace approximation algorithms*, in Proc. 18th Ann. ACM–SIAM Symp. Discrete Algorithms (SODA), 2007, pp. 532–540.
  - [122] L. SIROVICH AND M. KIRBY, *Low-dimensional procedure for the characterization of human faces*, J. Optical Soc. Amer. A, 4 (1987), pp. 519–524.
  - [123] D. SPIELMAN AND N. SRIVASTASA, *Graph sparsification by effective resistances*, in STOC '08: Proc. 40th Ann. ACM Symp. Theory of Computing, 2008.
  - [124] G. STEWART, *Accelerating the orthogonal iteration for the eigenvectors of a Hermitian matrix*, Num. Math., 13 (1969), pp. 362–376. 10.1007/BF02165413.
  - [125] G. W. STEWART, *Perturbation of pseudo-inverses, projections, and linear least squares problems*, SIAM Rev., 19 (1977), pp. 634–662.
  - [126] ———, *Four algorithms for the efficient computation of truncated pivoted QR approximations to a sparse matrix*, Numer. Math., 83 (1999), pp. 313–323.
  - [127] ———, *The decompositional approach to matrix computation*, Comput. Sci. Eng., (2000), pp. 50–59.

- [128] T. STROHMER AND R. VERSHYNIN, *A randomized Kaczmarz algorithm with exponential convergence*, J. Fourier Anal. Appl., 15 (2009), pp. 262–278.
- [129] J. SUN, Y. XIE, H. ZHANG, AND C. FALOUTSOS, *Less is more: Compact matrix decomposition for large sparse graphs*, Stat. Anal. Data Min., 1 (2008), pp. 6–22.
- [130] S. J. SZAREK, *Spaces with large distance from  $\ell_\infty^n$  and random matrices*, Amer. J. Math., 112 (1990), pp. 899–942.
- [131] A. SZLAM, M. MAGGIONI, AND R. R. COIFMAN, *Regularization on graphs with function-adapted diffusion processes*, J. Mach. Learn. Res., 9 (2008), pp. 1711–1739.
- [132] L. N. TREFETHEN AND D. BAU III, *Numerical Linear Algebra*, SIAM, Philadelphia, PA, 1997.
- [133] J. A. TROPP, *On the conditioning of random subdictionaries*, Appl. Comput. Harmon. Anal., 25 (2008), pp. 1–24.
- [134] ———, *Improved analysis of the subsampled randomized Hadamard transform*, Adv. Adaptive Data Anal., 3 (2011). To appear. Available at [arXiv:1011.1595](https://arxiv.org/abs/1011.1595).
- [135] J. VON NEUMANN AND H. H. GOLDSTINE, *Numerical inverting of matrices of high order*, Bull. Amer. Math. Soc., 53 (1947), pp. 1021–1099.
- [136] ———, *Numerical inverting of matrices of high order. II*, Proc. Amer. Math. Soc., 2 (1952), pp. 188–202.
- [137] F. WOOLFE, E. LIBERTY, V. ROKHLIN, AND M. TYGERT, *A fast randomized algorithm for the approximation of matrices*, Appl. Comp. Harmon. Anal., 25 (2008), pp. 335–366.