

## EFFICIENT ERROR AND VARIANCE ESTIMATION FOR RANDOMIZED MATRIX COMPUTATIONS\*

ETHAN N. EPPERLY<sup>†</sup> AND JOEL A. TROPP<sup>†</sup>

**Abstract.** Randomized matrix algorithms have become workhorse tools in scientific computing and machine learning. To use these algorithms safely in applications, they should be coupled with posterior error estimates to assess the quality of the output. To meet this need, this paper proposes two diagnostics: a leave-one-out error estimator for randomized low-rank approximations and a jackknife resampling method to estimate the variance of the output of a randomized matrix computation. Both of these diagnostics are rapid to compute for randomized low-rank approximation algorithms such as the randomized SVD and randomized Nyström approximation, and they provide useful information that can be used to assess the quality of the computed output and guide algorithmic parameter choices.

**Key words.** jackknife resampling, low-rank approximation, error estimation, randomized algorithms

**MSC codes.** 62F40, 65F55, 68W20

**DOI.** 10.1137/23M1558537

**1. Introduction.** In recent years, randomness has become an essential tool in the design of matrix algorithms [2, 7, 12, 13, 29], with randomized algorithms proving especially effective for matrix low-rank approximation. To use these algorithms safely in practice, they should be supported by *posterior error estimates* and other quality metrics that inform the user about the accuracy of the computational output.

This paper presents two diagnostic tools for randomized matrix computations:

- First, we provide a *leave-one-out posterior estimate* for the error  $\|\mathbf{A} - \mathbf{X}\|_F$  for a low-rank approximation  $\mathbf{X}$  to a matrix  $\mathbf{A}$  produced by randomized algorithms such as the randomized SVD or randomized Nyström approximation.
- Second, we present a jackknife method for estimating the *variance*  $\text{Var}(\mathbf{X}) := \mathbb{E} \|\mathbf{X} - \mathbb{E} \mathbf{X}\|_F^2$  of the matrix output  $\mathbf{X}$  of a randomized algorithm. The variance is a useful diagnostic: If the computation is sensitive to the randomness used by the algorithm, the computed output should be treated with suspicion.

By using novel downdating formulas (see (4.1) and (4.3) below), we can rapidly compute both of these estimators for widely used low-rank approximation methods such as the randomized SVD and randomized Nyström approximation. Our diagnostics are also *sample-efficient* in the sense that they require no additional matrix-vector products or other information beyond that used in the original algorithm. The speed and efficiency of these diagnostics make them compelling additions to workflows involving randomized matrix computation.

---

\*Submitted to the journal's Numerical Algorithms for Scientific Computing section March 13, 2023; accepted for publication (in revised form) October 30, 2023; published electronically February 8, 2024.

<https://doi.org/10.1137/23M1558537>

**Funding:** This material is based upon work supported by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research, Department of Energy Computational Science Graduate Fellowship under award DE-SC0021110. The second author was supported in part by ONR awards N00014-17-1-2146 and N00014-18-1-2363 and NSF FRG award 1952777.

<sup>†</sup>Department of Computing and Mathematical Sciences, California Institute of Technology, Pasadena, CA 91125 USA (epperly@caltech.edu, jtropp@cms.caltech.edu).

**1.1. Leave-one-out error estimation.** We begin by motivating our first diagnostic, a leave-one-out estimator for the error of a low-rank matrix approximation. For concreteness, we introduce this estimate in the context of Nyström approximation of positive-semidefinite (psd) matrices; see section 2 for a more general setting.

Suppose we want to approximate a psd matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$ , which we can only access by the matrix-vector product operation  $\omega \mapsto \mathbf{A}\omega$ . Using the matrix-vector product operation, we can compute the matrix-matrix product  $\mathbf{A}\Phi$  with a (random) matrix  $\Phi \in \mathbb{R}^{d \times s}$  with  $s$  columns and form the Nyström approximation [6, 24, 28]

$$(1.1) \quad \mathbf{A}\langle\Phi\rangle := \mathbf{A}\Phi(\Phi^* \mathbf{A}\Phi)^\dagger (\mathbf{A}\Phi)^*.$$

Here,  $*$  denotes the transpose and  $\dagger$  denotes the Moore–Penrose pseudoinverse. The result is a psd, rank- $s$  approximation  $\mathbf{A}\langle\Phi\rangle$  to the matrix  $\mathbf{A}$ . We will generate  $\Phi$  by applying  $q \geq 0$  steps of subspace iteration to a random test matrix  $\Omega$

$$(1.2) \quad \Phi = \mathbf{A}^q \Omega,$$

where  $\Omega$  is populated with statistically independent standard Gaussian entries. The quality of the Nyström approximation improves with a higher approximation rank  $s$  or number of subspace iteration steps  $q$ . Using the forthcoming Algorithm 4.1, we can compute  $\mathbf{A}\langle\Phi\rangle$  in the form of a compact eigenvalue decomposition:

$$(1.3) \quad \mathbf{A}\langle\Phi\rangle = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^*,$$

where  $\mathbf{V} \in \mathbb{R}^{d \times s}$  has orthonormal columns and  $\mathbf{\Lambda} \in \mathbb{R}_+^{s \times s}$  is diagonal. The cost of this procedure is  $\mathcal{O}(qs)$  matrix-vector products with  $\mathbf{A}$  and  $\mathcal{O}(ds^2)$  additional operations.

To use the Nyström approximation with confidence in applications and to guide the choice of parameters  $s$  and  $q$ , we need to understand the accuracy of the approximation  $\mathbf{A} \approx \mathbf{A}\langle\Phi\rangle$ . This motivates our question:

What is the most efficient way to estimate the error  $\|\mathbf{A} - \mathbf{A}\langle\Phi\rangle\|_F$ ?

Inspired by this question, this article proposes the *leave-one-out error estimator*, which provides an estimate of  $\|\mathbf{A} - \mathbf{A}\langle\Phi\rangle\|_F$  using only the information already collected from  $\mathbf{A}$  to form the Nyström approximation.

The leave-one-out estimator is built by recomputing the Nyström approximation using subsamples of the columns of the matrix  $\Omega$ . We can regard the Nyström approximation as a function of the test matrix  $\Omega$ :

$$\Omega \mapsto \mathbf{X} = \mathbf{X}(\Omega) := \mathbf{A}\langle\mathbf{A}^q \Omega\rangle.$$

Let  $\Omega^{(j)}$  denote  $\Omega$  without its  $j$ th column. Introduce *replicates*  $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(s)}$  by recomputing  $\mathbf{X}$  with each column of  $\Omega$  left out in turn:

$$\mathbf{X}^{(j)} = \mathbf{X}(\Omega^{(j)}) \quad \text{for } j = 1, 2, \dots, s.$$

Letting  $\omega_j$  denote the  $j$ th column of  $\Omega$  and  $\|\cdot\|$  denote the Euclidean norm, we define the *leave-one-out error estimator*

$$\widehat{\text{Err}}^2(\mathbf{X}, \mathbf{A}) := \frac{1}{s} \sum_{i=1}^s \|\mathbf{A} - \mathbf{X}^{(i)}\|_{\omega_i}^2.$$

As we show in Theorem 2.1, this error estimator is an *unbiased estimator for the mean-square error of the rank- $(s - 1)$  Nyström approximation*:

$$\mathbb{E} \widehat{\text{Err}}^2(\mathbf{X}, \mathbf{A}) = \mathbb{E} \|\mathbf{A} - \mathbf{X}(\boldsymbol{\Omega}^{(s)})\|_{\mathbb{F}}^2.$$

Once the Nyström approximation has been computed,  $\widehat{\text{Err}}(\mathbf{X}, \mathbf{A})$  is cheap to form, requiring at most  $\mathcal{O}(ds^2)$  additional operations (and only  $\mathcal{O}(s^3)$  operations if  $q = 0$ ). See subsection 4.1 for details. The error estimate  $\widehat{\text{Err}}(\mathbf{X}, \mathbf{A})$  requires no information about  $\mathbf{A}$  beyond what is required to form the approximation  $\mathbf{X}$ .

A good point of comparison for the leave-one-out error estimator is provided by the Girard–Hutchinson norm estimate [12, section 4.8]

$$(1.4) \quad \widehat{\text{Err}}_{\text{GH}}^2(\mathbf{X}, \mathbf{A}) = \frac{1}{t} \sum_{i=1}^t \|(\mathbf{A} - \mathbf{X})\boldsymbol{\nu}_i\|^2 \approx \|\mathbf{A} - \mathbf{X}\|_{\mathbb{F}}^2.$$

Here,  $\boldsymbol{\nu}_1, \dots, \boldsymbol{\nu}_t$  are independent standard Gaussian test vectors. This estimator requires  $t$  matrix-vector products with  $\mathbf{A}$  (beyond those used to form the approximation  $\mathbf{X}$ ), and the norm estimate becomes more accurate with larger  $t$ . The Girard–Hutchinson norm estimator is equivalent to the Girard–Hutchinson trace estimator [12, section 4.2] applied to  $(\mathbf{A} - \mathbf{X})^*(\mathbf{A} - \mathbf{X})$  and served as our inspiration for the leave-one-out error estimator. The leave-one-out estimator improves on the Girard–Hutchinson estimator as it requires no additional matrix-vector products with  $\mathbf{A}$  to compute. In addition, the quality of the leave-one-out estimator automatically improves when the approximation rank  $s$  increases, whereas the Girard–Hutchinson estimate only improves by using more matrix-vector products.

Figure 1 demonstrates the leave-one-out error estimator. In this figure, we apply single-pass Nyström approximation ( $q = 0$ ) to approximate psd kernel matrix  $\mathbf{A} \in \mathbb{R}^{10^4 \times 10^4}$  formed from a random subsample of  $10^4$  points from the QM9 dataset [18, 19] using approximation ranks  $5 \leq s \leq 150$ . In the left panel, we plot the mean error

$$(1.5) \quad \text{Err}(\mathbf{X}, \mathbf{A}) := \mathbb{E} \|\mathbf{A} - \mathbf{X}\|_{\mathbb{F}}$$

and the leave-one-out error estimate  $\widehat{\text{Err}}(\mathbf{X}, \mathbf{A})$ , estimated using 1000 trials. We see that the estimate  $\widehat{\text{Err}}(\mathbf{X}, \mathbf{A})$  closely tracks the true error. Moreover, the error estimate  $\widehat{\text{Err}}(\mathbf{X}, \mathbf{A})$  is fast to compute, with the error estimate taking less than 1% of the total runtime to form. The right panel compares the mean relative error

$$\text{mean relative error} = \mathbb{E} \left[ \frac{|\|\mathbf{A} - \mathbf{X}\|_{\mathbb{F}} - \text{Est}|}{\|\mathbf{A} - \mathbf{X}\|_{\mathbb{F}}} \right], \quad \text{Est} \in \left\{ \widehat{\text{Err}}(\mathbf{X}, \mathbf{A}), \widehat{\text{Err}}_{\text{GH}}(\mathbf{X}, \mathbf{A}) \right\}$$

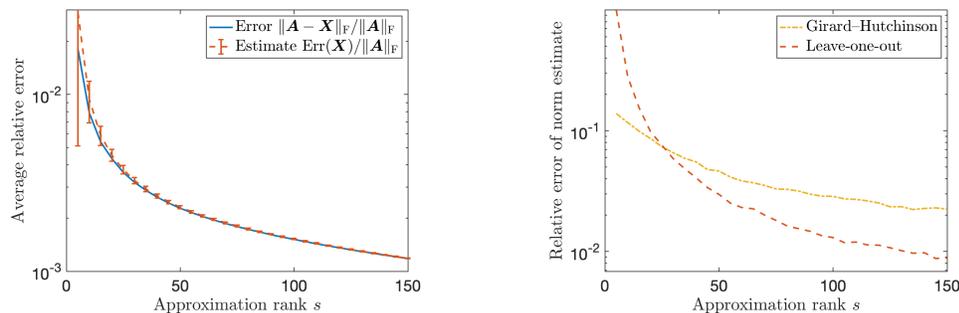


FIG. 1. *Leave-one-out error estimator.* Approximation error  $\text{Err}(\mathbf{X}, \mathbf{A})/\|\mathbf{A}\|_{\mathbb{F}}$  and normalized error estimate  $\widehat{\text{Err}}(\mathbf{X}, \mathbf{A})/\|\mathbf{A}\|_{\mathbb{F}}$  (left) and relative error for the leave-one-out and Girard–Hutchinson error estimators (right) for Nyström approximation  $\mathbf{X}$  to a psd kernel matrix  $\mathbf{A} \in \mathbb{R}^{10^4 \times 10^4}$ . Error and error estimate were computed by an average of 1000 trials. Error bars show one standard deviation.

for both the leave-one-out and Girard–Hutchinson error estimators. Following [25, section 7.9], we use  $t = 10$  matrix-vector products for the Girard–Hutchinson estimator. For  $s \geq 25$ , the leave-one-out estimator is more accurate than Girard–Hutchinson and, for all values of  $s$ , the leave-one-out estimator is cheaper to form than Girard–Hutchinson, requiring just  $\mathcal{O}(s^3)$  operations and no additional matrix-vector products.

**1.2. Matrix jackknife motivating example: Spectral clustering.** Randomized low-rank approximations can also be used for *spectral computations* (i.e., to approximate eigenvalues, eigenvectors, singular values, etc.). In this case, the low-rank approximation error  $\|\mathbf{A} - \mathbf{X}\|_{\mathbb{F}}$  may only provide indirect information about the accuracy of the computation. For situations such as this, we propose a *matrix jackknife variance estimate* as a diagnostic tool. In this section, we illustrate the value of this matrix jackknife approach in a spectral clustering application, before introducing the method in generality in section 3.

**1.2.1. Nyström-accelerated spectral clustering.** Spectral clustering [27] is an algorithm that uses eigenvectors to assign data points, say,  $\mathbf{c}_1, \dots, \mathbf{c}_d \in \mathbb{R}^m$ , into groups. To measure similarity between points, we employ a nonnegative, positive definite kernel function  $\kappa : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}_+$ . One popular choice is the square-exponential kernel

$$(1.6) \quad \kappa(\mathbf{c}, \mathbf{c}') = \exp\left(-\frac{\|\mathbf{c} - \mathbf{c}'\|^2}{2\sigma^2}\right).$$

To cluster the data points into groups, we perform the following steps:

1. Form the kernel matrix  $\mathbf{K}$  with entries  $k_{ij} = \kappa(\mathbf{c}_i, \mathbf{c}_j)$ .
2. Assemble the diagonal matrix  $\mathbf{D} = \text{diag}(\sum_{j=1}^d k_{ij} : i = 1, \dots, d)$ .
3. Compute the  $n_{\text{dim}}$  dominant eigenvectors  $\mathbf{U}$  of  $\mathbf{A} := \mathbf{D}^{-1/2} \mathbf{K} \mathbf{D}^{-1/2}$ .
4. Set  $\mathbf{W} := \mathbf{D}^{-1/2} \mathbf{U}$ .
5. Apply a general-purpose clustering algorithm, such as k-means [1] with  $n_{\text{cen}}$  centers, to the rows of  $\mathbf{W}$ .

Parameters  $n_{\text{dim}}$  and  $n_{\text{cen}}$  set the clustering space dimension and number of clusters.

If one uses direct methods for the eigenvalue problem, the cost of spectral clustering is dominated by the  $\mathcal{O}(d^3)$  cost for the eigenvector calculation in step 3. We can accelerate spectral clustering by using Nyström approximation (1.1). The modification is simple: Use the  $n_{\text{dim}}$  dominant eigenvectors of the Nyström approximation, accessible from the eigendecomposition (1.3), in place of the eigenvectors of  $\mathbf{A}$ .

**1.2.2. Variance estimation for spectral clustering.** As we refine the approximation by increasing  $s$ , the approximate eigenvectors  $\hat{\mathbf{U}}$  will converge to the true eigenvectors  $\mathbf{U}$  (provided  $\mathbf{U}$  is unique). But how do we know when we have taken  $s$  large enough? Our guiding principle is:

In order to trust the answer provided by a randomized algorithm, the output should be insensitive to the randomness used by the algorithm.

The *variance* of the matrix output  $\mathbf{X}$  of a randomized algorithm, defined as

$$(1.7) \quad \text{Var}(\mathbf{X}) := \mathbb{E} \|\mathbf{X} - \mathbb{E} \mathbf{X}\|_{\mathbb{F}}^2,$$

provides a quantitative measurement of the sensitivity of the algorithmic output to randomness used by the algorithm. In the context of spectral clustering, we can use a variance estimate to guide our choice of the rank  $s$ .

To understand the sensitivity of Nyström-accelerated spectral clustering to randomness in the algorithm, we need to specify a target matrix  $\mathbf{X}$  for variance estimation. The input to k-means clustering are the *coordinates*

$$\widehat{\mathbf{W}} := \mathbf{D}^{-1/2} \widehat{\mathbf{U}}.$$

To respect the invariance of k-means to scaling and rotation of the coordinates, our target for variance estimation will be

$$\mathbf{X} = \widehat{\mathbf{W}} \widehat{\mathbf{W}}^* / \|\widehat{\mathbf{W}} \widehat{\mathbf{W}}^*\|_{\mathbb{F}}.$$

**1.2.3. Jackknife variance estimation.** Jackknife variance estimation is similar to the leave-one-out error estimator in that we use replicates  $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(s)}$  recomputed by successively leaving out columns of the random test matrix  $\mathbf{\Omega}$ . As before, we view the target  $\mathbf{X}$  as a function

$$\mathbf{X} = \mathbf{X}(\mathbf{\Omega})$$

of the test matrix  $\mathbf{\Omega}$  defining the Nyström approximation by (1.1) and (1.2). Form *jackknife replicates*  $\mathbf{X}^{(j)}$  and their average  $\mathbf{X}^{(\cdot)}$  via

$$\mathbf{X}^{(j)} = \mathbf{X}(\mathbf{\Omega}^{(j)}) \quad \text{for } j = 1, \dots, s \quad \text{and} \quad \mathbf{X}^{(\cdot)} := \frac{1}{s} \sum_{j=1}^s \mathbf{X}^{(j)},$$

where  $\mathbf{\Omega}^{(j)}$  again denotes  $\mathbf{\Omega}$  without its  $j$ th column. The jackknife estimate for  $\text{Var}(\mathbf{X})$  is

$$\text{Jack}^2(\mathbf{X}) := \sum_{j=1}^s \left\| \mathbf{X}^{(j)} - \mathbf{X}^{(\cdot)} \right\|_{\mathbb{F}}^2.$$

Guarantees for this estimator are provided in Theorem 3.2. Algorithm SM1.1 and Program SM4 provide pseudocode and a MATLAB implementation of Nyström-accelerated spectral clustering with the jackknife estimate  $\text{Jack}(\mathbf{X})$ . The cost of forming the jackknife estimate  $\text{Jack}(\mathbf{X})$  is  $\mathcal{O}(s^2 d)$ , much faster than the  $\mathcal{O}(qsd^2)$  total cost of Nyström-accelerated spectral clustering.

**1.2.4. Numerical example.** To demonstrate the effectiveness of jackknife variance estimation for spectral clustering, we use the following experimental setup. Consider the task of separating the four letters JACK from a point cloud  $\mathbf{c}_1, \dots, \mathbf{c}_{9426} \in \mathbb{R}^2$ . For spectral clustering, use the square-exponential kernel (1.6) and parameters  $n_{\text{dim}} = n_{\text{cen}} = 4$ . For Nyström, use  $q = 3$  steps of subspace iteration and test a range of approximation ranks  $50 \leq s \leq 150$ . For each value of  $s$ , we run 1000 trials and report the natural Monte Carlo estimate of the standard deviation

$$(1.8) \quad \text{Std}(\mathbf{X}) = \left( \mathbb{E} \|\mathbf{X} - \mathbb{E} \mathbf{X}\|_{\mathbb{F}}^2 \right)^{1/2},$$

the mean and standard deviation for the *standard deviation* estimate  $\text{Jack}(\mathbf{X})$ , the relative error  $\text{Err}(\widehat{\mathbf{A}}, \mathbf{A}) / \|\mathbf{A}\|_{\mathbb{F}}$  for the Nyström approximation  $\widehat{\mathbf{A}}$ , and the empirical success probability for spectral clustering.

# JACK JACK

Correct clustering

Incorrect clustering

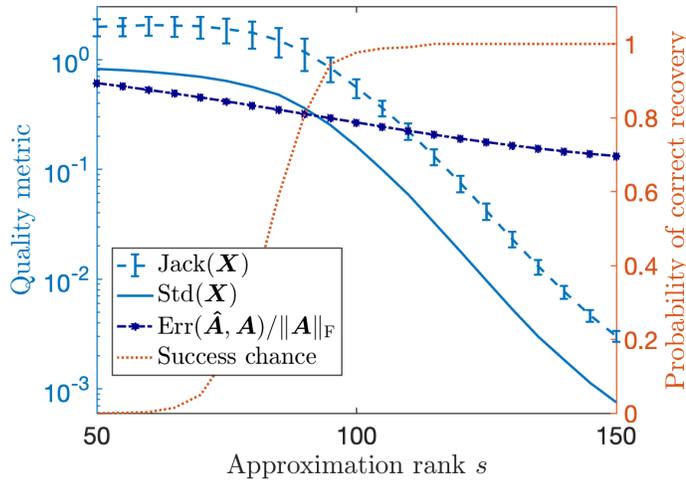


FIG. 2. *Matrix jackknife for spectral clustering.* Top: Example of correct (left) and incorrect (right) clusterings. Bottom: Standard deviation  $\text{Std}(\mathbf{X})$  and its jackknife estimate  $\text{Jack}(\mathbf{X})$  (left axis) and success probability of spectral clustering (right axis) versus Nyström approximation rank  $50 \leq s \leq 150$ .

Figure 2 shows the results. The estimate  $\text{Jack}(\mathbf{X})$  overestimates the standard deviation by a modest amount (a factor of five at most). The jackknife is not quantitatively sharp, but it is a reliable indicator of whether the variance is high or low.

The jackknife estimate  $\text{Jack}(\mathbf{X})$  can be used to determine whether the Nyström-based approximations to the eigenvectors are accurate enough for the clustering task. At rank  $s = 50$ , clustering is never performed correctly, and the jackknife estimate  $\text{Jack}(\mathbf{X}) \approx 2$  is high. As the approximation rank  $s$  is increased, clustering begins to succeed with higher and higher probability and the jackknife estimate  $\text{Jack}(\mathbf{X})$  decreases, indicating reduced variability.

The success of the jackknife estimate  $\text{Jack}(\mathbf{X})$  should be contrasted with the failure of the Nyström error  $\text{Err}(\hat{\mathbf{A}}, \mathbf{A})$  as a useful diagnostic for spectral clustering correctness. The Nyström error  $\text{Err}(\hat{\mathbf{A}}, \mathbf{A})$  decreases at a steady rate as  $s$  is increased—unlike the jackknife estimate, the plot of  $\text{Err}(\hat{\mathbf{A}}, \mathbf{A})$  does not show an inflection point indicating the transition between clustering failure and success.

On the basis of these experiments, we propose two possible uses for the jackknife estimate in a spectral clustering workflow:

- **User warning.** If the jackknife variance estimate is high, provide a warning to the user. This allows the user to determine and fix the problem for themselves by changing the Nyström parameters  $s, q$  or the spectral clustering parameters  $n_{\text{cen}}, n_{\text{dim}}$ .

- **Adaptive stopping.** Choose parameters  $s$  or  $q$  for the Nyström approximation adaptively at runtime by increasing these parameters until  $\text{Jack}(\mathbf{X})$  falls below a tolerance (e.g., 0.1).

These uses demonstrate the potential for jackknife variance estimation to be helpful in incorporating randomized low-rank approximation into general-purpose software.

**1.2.5. Benefits of matrix jackknife variance estimation.** The spectral clustering example demonstrates a number of virtues for matrix jackknife variance estimation:

- **Flexibility.** Matrix jackknife variance estimation can be applied to a very general target function  $\mathbf{X}(\boldsymbol{\Omega})$  depending on a random test matrix  $\boldsymbol{\Omega}$ . This allows the jackknife to be applied to a wide array of randomized low-rank approximation algorithms and allows the user to design the variance estimation target for their application.
- **Efficiency.** By using optimized algorithms (section 4 and subsection SM1.2), computation of the jackknife variance estimate can be very fast. For instance, for the clustering problem, the  $\mathcal{O}(s^2d)$  cost of the jackknife estimate is dwarfed by the  $\mathcal{O}(sd^2)$  cost of the clustering procedure. For  $s = 150$ , computing the jackknife variance estimate amounts to less than 3% of the total runtime.

**1.3. Outline.** Having introduced our two diagnostics, we present each in more generality; section 2 discusses leave-one-out error estimation and section 3 discusses the matrix jackknife. Section 4 discusses efficient computations for both of these diagnostics applied to the randomized SVD and Nyström approximation. Section 5 contains numerical experiments, and section 6 extends the matrix jackknife to higher Schatten norms ( $p > 2$ ).

**1.4. Notation.** We work over the field  $\mathbb{K} = \mathbb{R}$  or  $\mathbb{K} = \mathbb{C}$ . Notations  $*$ ,  $\dagger$ , and  $\|\cdot\|_F$  denote the conjugate transpose, Moore–Penrose pseudoinverse, and Frobenius norm. The expectation of a random variable  $X$  is denoted  $\mathbb{E}X$ , and its variance is defined as  $\text{Var}(X) := \mathbb{E}|X - \mathbb{E}X|^2$ . We adopt the convention that nonlinear operators bind before the expectation; for example,  $\mathbb{E}X^2 := \mathbb{E}(X^2)$ . The variance of a random matrix is given by (1.7).

**2. Leave-one-out error estimation for low-rank approximation.** In this section, we present the leave-one-out error estimation technique introduced in subsection 1.1 for more general randomized matrix approximations.

**2.1. The estimator.** Let  $\mathbf{A} \in \mathbb{K}^{d_1 \times d_2}$  be a matrix we seek to approximate by a randomized approximation  $\mathbf{X}$ . We are interested in a general class of algorithms which collect information about the matrix  $\mathbf{A}$  by matrix-vector products

$$\mathbf{A}\boldsymbol{\omega}_1, \dots, \mathbf{A}\boldsymbol{\omega}_s$$

with random test vectors  $\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_s$ . Many algorithms are defined for an arbitrary number of test vectors  $s$ , allowing us to construct error estimates by leaving out a test vector, resulting in an approximation  $\mathbf{X}_{s-1}$  defined using only  $s-1$  vectors. This motivates the following abstract setup:

- Let  $\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_s$  be independent and identically distributed (iid) random vectors in  $\mathbb{K}^{d_2}$  that are *isotropic*:  $\mathbb{E}[\boldsymbol{\omega}_j \boldsymbol{\omega}_j^*] = \mathbf{I}$ .
- Let  $\mathbf{X}$  denote one of two matrix estimators defined for  $s$  or  $s-1$  inputs:

$$\mathbf{X} : (\mathbb{K}^{d_2})^s \rightarrow \mathbb{K}^{d_1 \times d_2} \quad \text{or} \quad \mathbf{X} : (\mathbb{K}^{d_2})^{s-1} \rightarrow \mathbb{K}^{d_1 \times d_2}.$$

- Define estimates  $\mathbf{X}_s := \mathbf{X}(\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_s)$  and  $\mathbf{X}_{s-1} := \mathbf{X}(\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_{s-1})$ .

Examples of estimators which fit this description include randomized Nyström approximation, the randomized SVD [7], and randomized block Krylov iteration [14, 23].

We seek to approximate the mean-square error

$$\text{MSE}(\mathbf{X}_{s-1}, \mathbf{A}) = \mathbb{E} \|\mathbf{A} - \mathbf{X}_{s-1}\|_{\mathbb{F}}^2$$

of the  $(s - 1)$ -sample approximation  $\mathbf{X}_{s-1}$  as a proxy for the mean-square error  $\text{MSE}(\mathbf{X}_s)$  of the  $s$ -sample approximation  $\mathbf{X}_s$ . Define the *leave-one-out mean-square error estimate*

$$(2.1) \quad \widehat{\text{Err}}^2(\mathbf{X}_{s-1}, \mathbf{A}) := \frac{1}{s} \sum_{j=1}^s \left\| (\mathbf{A} - \mathbf{X}^{(j)}) \boldsymbol{\omega}_j \right\|^2,$$

where the replicates  $\mathbf{X}^{(j)}$  are

$$\mathbf{X}^{(j)} = \mathbf{X}(\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_{j-1}, \boldsymbol{\omega}_{j+1}, \dots, \boldsymbol{\omega}_s) \quad \text{for } j = 1, \dots, s.$$

This estimator is an unbiased estimator for  $\text{MSE}(\mathbf{X}_{s-1}, \mathbf{A})$ .

**THEOREM 2.1** (leave-one-out error estimator). *With the prevailing notation,*

$$\text{MSE}(\mathbf{X}_{s-1}, \mathbf{A}) = \mathbb{E} \widehat{\text{Err}}^2(\mathbf{X}_{s-1}, \mathbf{A}).$$

*Proof.* For each  $j$ ,  $\mathbf{X}^{(j)}$  and  $\boldsymbol{\omega}_j$  are independent. Consequently, letting  $\mathbb{E}_j$  denote an expectation over the randomness in  $\boldsymbol{\omega}_j$  alone, we compute

$$\begin{aligned} \mathbb{E}_j \left\| (\mathbf{A} - \mathbf{X}^{(j)}) \boldsymbol{\omega}_j \right\|^2 &= \mathbb{E}_j [\boldsymbol{\omega}_j^* (\mathbf{A} - \mathbf{X}^{(j)})^* (\mathbf{A} - \mathbf{X}^{(j)}) \boldsymbol{\omega}_j] \\ &= \mathbb{E}_j \text{tr} [(\mathbf{A} - \mathbf{X}^{(j)})^* (\mathbf{A} - \mathbf{X}^{(j)}) \boldsymbol{\omega}_j \boldsymbol{\omega}_j^*] \\ &= \text{tr} ((\mathbf{A} - \mathbf{X}^{(j)})^* (\mathbf{A} - \mathbf{X}^{(j)}) \mathbb{E} [\boldsymbol{\omega}_j \boldsymbol{\omega}_j^*]) \\ &= \text{tr} ((\mathbf{A} - \mathbf{X}^{(j)})^* (\mathbf{A} - \mathbf{X}^{(j)})) = \|\mathbf{A} - \mathbf{X}^{(j)}\|_{\mathbb{F}}^2. \end{aligned}$$

The first line is an identity for the Frobenius norm, the second line is the cyclic property of the trace, the third line is the independence of  $\mathbf{X}^{(j)}$  and  $\boldsymbol{\omega}_j$ , and the fourth line is the isotropic property of  $\boldsymbol{\omega}_j$ . Thus, by the tower property of conditional expectation, we conclude

$$\begin{aligned} \mathbb{E} \widehat{\text{Err}}^2(\mathbf{X}_{s-1}, \mathbf{A}) &= \frac{1}{s} \sum_{j=1}^s \mathbb{E} [\mathbb{E}_j \left\| (\mathbf{A} - \mathbf{X}^{(j)}) \boldsymbol{\omega}_j \right\|^2] \\ &= \frac{1}{s} \sum_{j=1}^s \mathbb{E} \left\| (\mathbf{A} - \mathbf{X}^{(j)}) \boldsymbol{\omega}_j \right\|_{\mathbb{F}}^2 = \text{MSE}(\mathbf{X}_{s-1}, \mathbf{A}). \end{aligned}$$

This confirms the theorem. □

Numerical evidence for the quality of this error estimate is provided in Figures 1 and 3. With efficient algorithms (section 4), the leave-one-out error estimator is rapid to compute for the randomized SVD and Nyström approximation.

**2.2. Alternatives.** Two alternatives to the leave-one-out error estimator are worth mentioning. First, in many situations, it may be possible and computationally

cheap to simply compute the error  $\|\mathbf{A} - \mathbf{X}\|_F$  directly. For instance, if  $\mathbf{X}$  is the approximation produced by the randomized SVD [7], then

$$\|\mathbf{A} - \mathbf{X}\|_F^2 = \|\mathbf{A}\|_F^2 - \|\mathbf{X}\|_F^2,$$

which facilitates fast computation of the error if  $\mathbf{A}$  is a dense or sparse matrix stored in memory. The leave-one-out error estimator should only be used if direct computation of the error is not possible or is too expensive. This is the case, for example, in the black-box setting where one has access to  $\mathbf{A}$  only through the matrix-vector product  $\boldsymbol{\omega} \mapsto \mathbf{A}\boldsymbol{\omega}$  and adjoint-vector product  $\boldsymbol{\omega} \mapsto \mathbf{A}^*\boldsymbol{\omega}$  operations.

A second alternative is the Girard–Hutchinson norm estimator (1.4), discussed in subsection 1.1. The leave-one-out estimator improves on the Girard–Hutchinson estimator as the leave-one-out estimate does not require any additional matrix-vector products and automatically improves in quality as the number of vectors  $s$  is increased.

**3. Matrix jackknife variance estimation.** This section outlines our proposal for *matrix jackknife variance estimation* for more general randomized matrix algorithms. Subsection 3.1 reviews jackknife variance estimation for scalar quantities. We introduce and analyze the matrix jackknife variance estimator in subsection 3.2. Subsections 3.3 to 3.5 discuss potential applications of the matrix jackknife and complementary topics.

**3.1. Tukey’s jackknife variance estimator and the Efron–Stein–Steele inequality.** To motivate our matrix jackknife proposal, we begin by presenting the jackknife variance estimator [26] for scalar estimators due to Tukey in subsection 3.1.1. In subsection 3.1.2, we discuss the Efron–Stein–Steele inequality used in its analysis.

**3.1.1. Tukey’s jackknife variance estimator.** Consider the problem of estimating the variance of a statistical estimator computed from  $s$  random samples. We assume it makes sense to evaluate the estimator with fewer than  $s$  samples, as is the case for many classical estimators like the sample mean and variance. This motivates the following setup:

- Let  $\omega_1, \dots, \omega_s$  be iid random elements taking values in a measurable space  $\Omega$ .
- Let  $f$  denote either one of two estimators, defined for  $s$  or  $s - 1$  arguments:

$$f : \Omega^s \rightarrow \mathbb{K} \quad \text{or} \quad f : \Omega^{s-1} \rightarrow \mathbb{K}.$$

- Assume that  $f$  is invariant to a reordering of its inputs:

$$f(\omega_1, \dots, \omega_s) = f(\omega_{\pi(1)}, \dots, \omega_{\pi(s)}) \quad \text{for any permutation } \pi.$$

- Define estimates  $E_{s-1} := f(\omega_1, \dots, \omega_{s-1})$  and  $E_s := f(\omega_1, \dots, \omega_s)$ .

We think of  $E_s$  as a statistic computed from a collection of samples  $\omega_1, \dots, \omega_s$ . We can also evaluate the statistic with only  $s - 1$  samples, resulting in  $E_{s-1}$ .

Tukey’s jackknife variance estimator provides an estimate for  $\text{Var}(E_{s-1})$ , which serves as a proxy for the variance of the  $s$ -sample estimator  $E_s$ . Define jackknife replicates and mean

$$E^{(j)} := f(\omega_1, \dots, \omega_{j-1}, \omega_{j+1}, \dots, \omega_s) \quad \text{for each } j = 1, 2, \dots, s; \quad E^{(\cdot)} := \frac{1}{s} \sum_{j=1}^s E^{(j)}.$$

The quantities  $E^{(1)}, \dots, E^{(s)}$  represent the statistic recomputed with each of the samples  $\omega_1, \dots, \omega_s$  left out in turn.

Tukey’s estimator for  $\text{Var}(E_{s-1})$  is given by

$$(3.1) \quad \widehat{\text{Var}}(E_{s-1}) := \sum_{j=1}^s \left| E^{(j)} - E^{(\cdot)} \right|^2.$$

Observe that Tukey’s estimator (3.1) is the sample variance of the jackknife replicates  $E^{(j)}$  up to a normalizing constant. The form of Tukey’s estimator suggests that the distribution of the jackknife replicates somehow approximates the distribution of the estimator. This intuition can be formalized using the Efron–Stein–Steele inequality.

**3.1.2. Efron–Stein–Steele inequality.** To analyze Tukey’s estimator, we rely on an inequality of Efron and Stein [4], which was improved by Steele [21].

**FACT 3.1** (Efron–Stein–Steele inequality). *Let  $\omega_1, \dots, \omega_s \in \Omega$  be independent elements in a measurable space  $\Omega$ , and let  $f : \Omega^s \rightarrow \mathbb{K}$  be measurable. Let  $(\omega'_j : j = 1, \dots, s)$  be an independent copy of  $(\omega_j : j = 1, \dots, s)$ . Then*

$$(3.2) \quad \text{Var}(f(\omega_1, \dots, \omega_s)) \leq \frac{1}{2} \sum_{i=1}^s \mathbb{E} \left| f(\omega_1, \dots, \omega_s) - f(\omega_1, \dots, \omega_{j-1}, \omega'_j, \omega_{j+1}, \dots, \omega_s) \right|^2.$$

The complex-valued version of the inequality presented here follows from the more standard version for real values by treating the real and imaginary parts separately.

In the setting of Tukey’s estimator (3.1), the samples  $\omega_1, \dots, \omega_{s-1}$  are identically distributed and the function  $f$  depends symmetrically on its arguments. Therefore, the last sample  $\omega_s$  can be used to fill the role of each  $\omega'_j$  in the right-hand side of (3.2). As a consequence, the Efron–Stein–Steele inequality shows that

$$(3.3) \quad \begin{aligned} \text{Var}(E_{s-1}) &\leq \frac{1}{2} \sum_{i=1}^{s-1} \mathbb{E} \left( E^{(j)} - E^{(s)} \right)^2 = \frac{1}{2s} \sum_{i,j=1}^s \mathbb{E} \left( E^{(i)} - E^{(j)} \right)^2 \\ &= \sum_{j=1}^s \mathbb{E} \left( E^{(j)} - E^{(\cdot)} \right)^2 = \mathbb{E} \widehat{\text{Var}}(E_{s-1}). \end{aligned}$$

To move from the first line to the second, we expand the square, use the definition of the mean  $E^{(\cdot)} = s^{-1} \sum_{j=1}^s E^{(j)}$ , and regroup terms. This computation shows that Tukey’s variance estimator (3.1) *overestimates* the true variance on average.

**3.2. The matrix jackknife estimator of variance.** Suppose we are interested in the variance of the output  $\mathbf{X} \in \mathbb{K}^{d_1 \times d_2}$  to a randomized matrix algorithm. Similar to the scalar setting, we assume that  $\mathbf{X}$  is a function of independent samples  $\omega_1, \dots, \omega_s$  and that it makes sense to evaluate  $\mathbf{X}$  with fewer than  $s$  samples. The formal setup is as follows:

- Let  $\omega_1, \dots, \omega_s$  be iid random elements in a measurable space  $\Omega$ .
- Let  $\mathbf{X}$  denote one of two matrix estimators defined for  $s$  or  $s - 1$  inputs:

$$\mathbf{X} : \Omega^s \rightarrow \mathbb{K}^{d_1 \times d_2} \quad \text{or} \quad \mathbf{X} : \Omega^{s-1} \rightarrow \mathbb{K}^{d_1 \times d_2}.$$

- Assume  $\mathbf{X}$  is invariant to reordering of its inputs:

$$\mathbf{X}(\omega_1, \dots, \omega_s) = \mathbf{X}(\omega_{\pi(1)}, \dots, \omega_{\pi(s)}) \quad \text{for any permutation } \pi.$$

- Define estimates  $\mathbf{X}_s := \mathbf{X}(\omega_1, \dots, \omega_s)$  and  $\mathbf{X}_{s-1} := \mathbf{X}(\omega_1, \dots, \omega_{s-1})$ .

For a randomized low-rank approximation algorithm, the samples  $\omega_1, \dots, \omega_s$  might represent the columns of a test matrix  $\mathbf{\Omega}$ , as they did in subsection 1.2.

We are interested in estimating the variance of  $\mathbf{X}_{s-1}$  as a proxy for the variance of  $\mathbf{X}_s$ . We expect that adding additional samples will refine the approximation and thus reduce its variance. Define jackknife replicates  $\mathbf{X}^{(j)}$  and their average  $\mathbf{X}^{(\cdot)}$ ,

$$\mathbf{X}^{(j)} = \mathbf{X}(\omega_1, \dots, \omega_{j-1}, \omega_{j+1}, \dots, \omega_s) \quad \text{for each } j = 1, \dots, s; \quad \mathbf{X}^{(\cdot)} := \frac{1}{s} \sum_{j=1}^s \mathbf{X}^{(j)}.$$

We propose the matrix jackknife estimate

$$(3.4) \quad \text{Jack}^2(\mathbf{X}_{s-1}) := \sum_{j=1}^s \left\| \mathbf{X}^{(j)} - \mathbf{X}^{(\cdot)} \right\|_F^2$$

for the variance  $\text{Var}(\mathbf{X}_{s-1})$ . The estimator  $\text{Jack}(\mathbf{X}_{s-1})$  can be efficiently computed for several randomized low-rank approximations, as we shall demonstrate in section 4. Similar to the classic jackknife variance estimator, we can use the Efron–Stein–Steele inequality to show that this variance estimate is an overestimate on average.

**THEOREM 3.2 (matrix jackknife).** *With the prevailing notation,*

$$(3.5) \quad \text{Var}(\mathbf{X}_{s-1}) \leq \mathbb{E} \text{Jack}^2(\mathbf{X}_{s-1}).$$

*Proof.* Fix a pair of indices  $1 \leq m \leq d_1$  and  $1 \leq n \leq d_2$ . Applying (3.3) to the  $(m, n)$ -matrix entry  $(\mathbf{X}_{s-1})_{mn}$ , we observe

$$\mathbb{E} |(\mathbf{X}_{s-1})_{mn} - \mathbb{E}(\mathbf{X}_{s-1})_{mn}|^2 \leq \mathbb{E} \sum_{j=1}^s \left| \mathbf{X}_{mn}^{(j)} - \mathbf{X}_{mn}^{(\cdot)} \right|^2.$$

Summing this equation over all  $1 \leq m \leq d_1$  and  $1 \leq n \leq d_2$  yields the stated result.  $\square$

When the jackknife variance estimate is small, Theorem 3.2 shows the variance of the approximation is also small. Empirical evidence (section 5) suggests that  $\text{Var}(\mathbf{X}_{s-1})$  and  $\mathbb{E} \text{Jack}^2(\mathbf{X}_{s-1})$  tend to be within an order of magnitude for the algorithms we considered.

It is natural to ask whether we can develop and theoretically jackknife estimates of the bias of randomized matrix algorithms to complement our variance estimate, perhaps using the natural analogue of Quenouille’s scalar jackknife bias estimate [17]. This is an interesting question for future work. As we have demonstrated in subsection 1.2 and will further demonstrate in section 5, our variance estimate already provides useful and actionable information for randomized matrix algorithms.

**3.3. Uses for matrix jackknife variance estimator.** Variance is a useful diagnostic for randomized matrix approximations. When the variance is large, it suggests that one of two issues has arisen:

1. More samples are needed to refine the approximation.
2. The underlying approximation problem is badly conditioned.

In either case, the jackknife variance estimate can provide evidence that the computed output should not be trusted.

We anticipate the primary use case for matrix jackknife variance estimation will be for computations using eigenvectors or singular vectors computed by randomized low-rank approximation algorithms such as the randomized SVD and Nyström approximation. Spectral computations with randomized algorithms currently lack effective posterior estimates, making jackknife variance estimation one of the few available

tools to assess the quality of the outputs of such computations at runtime. In the context of spectral computations, matrix jackknife variance estimation can be used to adaptively determine the approximation rank  $s$  needed to achieve outputs of sufficiently high quality. This was demonstrated in subsection 1.2, where we used jackknife variance estimation to determine how large to pick  $s$  in a spectral clustering context. As we will later demonstrate in section 5, we can also use the jackknife variance estimation to detect ill-disposed eigenvectors and to get *coordinatewise* variance estimates for singular vector computations.

**3.4. Matrix jackknife versus scalar jackknife.** Sometimes, we are only interested in scalar outputs of a randomized matrix computation, such as eigenvalues or singular values, entries of eigenvectors or singular vectors, or the trace. In these cases, it might be more efficient to directly apply Tukey’s variance estimator (3.1) to assess the variance of these scalars. The matrix jackknife may still be a useful tool because it gives *simultaneous* variance estimates over many scalar quantities. As examples, the matrix jackknife estimates the maximum variance over all linear functionals:

$$\mathbb{E} \max_{\|C\|_F \leq 1} |\text{tr}(C\mathbf{X}_{s-1}) - \text{tr}(C\mathbb{E}\mathbf{X}_{s-1})|^2 = \mathbb{E} \|\mathbf{X}_{s-1} - \mathbb{E}\mathbf{X}_{s-1}\|_F^2 \leq \mathbb{E} \text{Jack}^2(\mathbf{X}_{s-1}).$$

The matrix jackknife gives the following variance estimate for the singular values:

$$\mathbb{E} \sum_{j=1}^{\min(d_1, d_2)} |\sigma_j(\mathbf{X}_{s-1}) - \sigma_j(\mathbb{E}\mathbf{X}_{s-1})|^2 = \mathbb{E} \|\mathbf{X}_{s-1} - \mathbb{E}\mathbf{X}_{s-1}\|_F^2 \leq \mathbb{E} \text{Jack}^2(\mathbf{X}_{s-1}).$$

Thus, the matrix jackknife is appealing even when one is interested in multiple scalar-valued functions of the matrix approximation  $\mathbf{X}_{s-1}$ . In addition, efficient algorithms for matrix jackknife variance estimation, as detailed in section 4, are useful for scalar jackknife variance estimation of functionals of a randomized matrix approximation.

**3.5. Related work: Bootstrap for randomized matrix algorithms.** Bootstrap resampling [3, section 5], a close relative of the jackknife, has seen several applications to matrix computations. An early use case was to provide confidence intervals for eigenvalues and eigenvectors of sample covariance matrices [22, section 7.2].

A more recent line of work, led by Lopes and collaborators, applies bootstrap resampling to randomized matrix algorithms [9, 10, 11, 30]. The work closest to ours [9] uses the bootstrap to provide asymptotically sharp estimates of the error quantiles with regards to general error metrics for each singular value and singular vector computed by a *Monte Carlo-type* sketched SVD algorithm. In this special case, the bootstrap provides more fine-grained information than the matrix jackknife. Unfortunately, this sketched SVD is a poor computational method because its error decays at the Monte Carlo rate.

The main benefit of our matrix jackknife approach is that it is effective for very general matrix algorithms, such as the randomized SVD and Nyström approximation. These algorithms are used widely in practice because they achieve spectral accuracy, producing errors comparable with the best low-rank approximation [7]. As we demonstrate in section SM2, a straightforward application of the bootstrap variance of Efron [3, section 5.1] to the randomized SVD can produce standard deviation estimates which are incorrect by over four orders of magnitude.

**4. Case studies in low-rank approximation.** In this section, we develop *efficient* computational procedures to compute the jackknife variance estimate and leave-one-out error estimate for two randomized low-rank approximations, randomized Nyström approximation (subsection 4.1) and the randomized SVD (subsection 4.2).

**4.1. Nyström approximation.** Given a test matrix  $\Omega \in \mathbb{K}^{d \times s}$ , consider again the Nyström approximation (1.1) with  $q$  steps of subspace iteration (1.2) applied to a psd matrix  $\mathbf{A} \in \mathbb{K}^{d \times d}$ :

$$\mathbf{X} = \mathbf{X}(\Omega) := \mathbf{A} \langle \Phi \rangle = \mathbf{A} \Phi (\Phi^* \mathbf{A} \Phi)^\dagger (\mathbf{A} \Phi)^*, \quad \text{where } \Phi = \mathbf{A}^q \Omega.$$

The Nyström approximation  $\mathbf{X}$  is the best psd approximation to  $\mathbf{A}$  spanned by  $\mathbf{A} \Phi$  with a psd residual. We focus on the case where  $\Omega$  is populated with iid, isotropic columns, such as when  $\Omega$  is a standard Gaussian matrix.

We work with the Nyström approximation in eigenvalue decomposition form:

$$\mathbf{X} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^*,$$

where  $\mathbf{V} \in \mathbb{K}^{d \times s}$  has orthonormal columns and  $\mathbf{\Lambda} \in \mathbb{R}_+^{s \times s}$  is diagonal. To facilitate efficient computations of our diagnostics, we can compute the Nyström approximation in eigendecomposition form as follows:

1. Draw a test matrix  $\Omega \in \mathbb{K}^{n \times s}$  with iid isotropic columns.
2. Apply subspace iteration  $\Phi = \mathbf{A}^q \Omega$ .
3. Compute the product  $\mathbf{Y} = \mathbf{A} \Omega$ .
4. Orthonormalize  $\mathbf{Q} = \text{orth}(\mathbf{Y})$  using economy QR factorization  $\mathbf{Y} = \mathbf{Q} \mathbf{R}$ .
5. Compute  $\mathbf{H} = \Omega^* \mathbf{Y}$  and Cholesky factorize  $\mathbf{H} = \mathbf{C}^* \mathbf{C}$ .
6. Obtain a singular value decomposition  $\mathbf{R} \mathbf{C}^{-1} = \mathbf{U} \mathbf{\Sigma} \mathbf{Z}^*$ .
7. Set  $\mathbf{\Lambda} := \mathbf{\Sigma}^2$  and  $\mathbf{V} := \mathbf{Q} \mathbf{U}$ .

Algorithm 4.1 provides an implementation with  $q = 0$  with tricks to improve its numerical stability adapted from [8, 24]. For  $q > 0$ , it may be necessary to introduce additional orthogonalization steps for reasons of numerical stability [20, Algorithm 5.2].

Treating the Nyström approximation  $\mathbf{X}$  as a symmetric function of the iid, isotropic columns  $\omega_1, \dots, \omega_s$  of the test matrix  $\Omega$ ,

$$\mathbf{X} = \mathbf{X}(\omega_1, \dots, \omega_s),$$

we can apply both the leave-one-out error estimator and jackknife variance estimation to  $\mathbf{X}$ . Define replicates

$$\mathbf{X}^{(j)} = \mathbf{X}(\omega_1, \dots, \omega_{j-1}, \omega_{j+1}, \dots, \omega_s).$$

To compute the replicates efficiently, we use the update formula [5, equation (2.4)]

$$(4.1a) \quad \mathbf{X}^{(j)} = \mathbf{V} (\mathbf{\Lambda} - \mathbf{t}_j \mathbf{t}_j^*) \mathbf{V}^*,$$

where  $\mathbf{t}_1, \dots, \mathbf{t}_s$  are the columns of the matrix

$$(4.1b) \quad \mathbf{T} = \mathbf{U}^* \mathbf{R} \mathbf{H}^{-1} \cdot \text{diag} \left( (\mathbf{H}^{-1})_{ii}^{-1/2} : i = 1, 2, \dots, s \right).$$

A derivation of this formula is provided in subsection A.1.

The update formula facilitates efficient algorithms for the leave-one-out error estimator and jackknife variance estimates for the Nyström approximation and derived quantities like spectral projectors and truncation of  $\mathbf{X}$  to rank  $r < s$ . To not belabor the point by presenting all possible variations, Algorithm 4.1 presents an implementation of Nyström approximation without subspace iteration (i.e.,  $q = 0$ ) with the leave-one-out error estimate  $\widehat{\text{Err}}(\mathbf{X}, \mathbf{A})$ . The computation of the error estimate is a simple addition to the algorithm, requiring just a single line and taking only  $\mathcal{O}(s^3)$  operations, independent of the size  $d$  of the input matrix. Further variants are discussed in subsection SM1.1 and a MATLAB implementation is provided in Program SM1.

---

**Algorithm 4.1.** Nyström approximation ( $q = 0$ ) with leave-one-out error estimate.

---

**Input:**  $\mathbf{A} \in \mathbb{K}^{d \times d}$  to be approximated and approximation rank  $s$

**Output:** Factors  $\mathbf{V} \in \mathbb{K}^{d \times s}$  and  $\widehat{\mathbf{\Lambda}} \in \mathbb{R}^{s \times s}$  defining a rank- $s$  approximation  $\mathbf{X} = \mathbf{V}\widehat{\mathbf{\Lambda}}\mathbf{V}^*$  and leave-one-out error estimate  $\widehat{\text{Err}}$

- 1:  $\mathbf{\Omega} \leftarrow \text{randn}(d, s)$
  - 2:  $\mathbf{Y} \leftarrow \mathbf{A}\mathbf{\Omega}$
  - 3:  $\nu \leftarrow \epsilon_{\text{mach}} \|\mathbf{Y}\|$  and  $\mathbf{Y} \leftarrow \mathbf{Y} + \nu\mathbf{\Omega}$  ▷ Shift for numerical stability
  - 4:  $(\mathbf{Q}, \mathbf{R}) \leftarrow \text{qr}(\mathbf{Y}, \text{'econ'})$  ▷ Economy QR factorization
  - 5:  $\mathbf{H} \leftarrow \mathbf{\Omega}^* \mathbf{Y}$
  - 6:  $\mathbf{C} \leftarrow \text{chol}((\mathbf{H} + \mathbf{H}^*)/2)$  ▷ Upper triangular Cholesky decomposition  $\mathbf{H} = \mathbf{C}^* \mathbf{C}$
  - 7:  $(\mathbf{U}, \mathbf{\Sigma}, \sim) \leftarrow \text{svd}(\mathbf{R}\mathbf{C}^{-1})$  ▷ Triangular solve
  - 8:  $\mathbf{\Lambda} \leftarrow \max(\mathbf{\Sigma}^2 - \nu\mathbf{I}, 0)$  ▷ Entrywise maximum, shift back for numerical stability
  - 9:  $\mathbf{V} \leftarrow \mathbf{Q}\mathbf{U}$
  - 10:  $\widehat{\text{Err}} \leftarrow \left\| (\mathbf{R}\mathbf{C}^{-1})\mathbf{C}^{-*} \cdot \text{diag}\{(\mathbf{H}_{ii}^{-1})^{-1} : i = 1, \dots, s\} \right\|_{\text{F}} / \sqrt{s}$
- 

**4.2. Randomized SVD.** The randomized SVD computes a rank- $s$  approximation to  $\mathbf{A} \in \mathbb{K}^{d_1 \times d_2}$  formed as an economy SVD  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ , where  $\mathbf{U} \in \mathbb{K}^{d_1 \times s}$  and  $\mathbf{V} \in \mathbb{K}^{d_2 \times s}$  have orthonormal columns and  $\mathbf{\Sigma} \in \mathbb{R}_+^{s \times s}$  is diagonal. With  $q \geq 0$  steps of subspace iteration, the algorithm proceeds as follows:

1. Draw a *test matrix*  $\mathbf{\Omega} \in \mathbb{C}^{d_2 \times s}$  with iid isotropic columns.
2. Compute the product  $\mathbf{Y} = (\mathbf{A}\mathbf{A}^*)^q \mathbf{A}\mathbf{\Omega}$ .
3. Orthonormalize  $\mathbf{Q} = \text{orth}(\mathbf{Y})$  using economy QR factorization  $\mathbf{Y} = \mathbf{Q}\mathbf{R}$ .
4. Form the matrix  $\mathbf{C} = \mathbf{Q}^* \mathbf{A}$ .
5. Compute an economy SVD  $\mathbf{C} = \mathbf{W}\mathbf{\Sigma}\mathbf{V}^*$ .
6. Set  $\mathbf{U} = \mathbf{Q}\mathbf{W}$ .

The output  $\mathbf{X}$  is a symmetric function of the iid, isotropic columns  $\omega_1, \dots, \omega_s$  of  $\mathbf{\Omega}$ ,

$$\mathbf{X} = \mathbf{X}(\omega_1, \dots, \omega_s),$$

making it a candidate for leave-one-out error estimation and jackknife variance estimation.

To compute the replicates efficiently, we will use the following update formula for the  $\mathbf{Q}$  matrix in the randomized SVD [5, equation (2.1)]:

$$(4.2) \quad \mathbf{Q}^{(j)} \left( \mathbf{Q}^{(j)} \right)^* = \mathbf{Q} (\mathbf{I} - \mathbf{t}_j \mathbf{t}_j^*) \mathbf{Q}^*,$$

where  $\mathbf{Q}^{(j)}$  denotes the  $\mathbf{Q}$  matrix produced by the randomized SVD algorithm executed without the  $j$ th column of  $\mathbf{\Omega}$  and the vectors  $\mathbf{t}_1, \dots, \mathbf{t}_s$  are the normalized columns of  $(\mathbf{R}^*)^{-1}$ . See Appendix A.2 for a derivation. With this formula, the replicates are easily computed,

$$(4.3) \quad \mathbf{X}^{(j)} = \mathbf{Q}^{(j)} \left( \mathbf{Q}^{(j)} \right)^* \mathbf{A} = \mathbf{Q} (\mathbf{I} - \mathbf{t}_j \mathbf{t}_j^*) \mathbf{Q}^* \mathbf{A} = \mathbf{U} (\mathbf{I} - \mathbf{W}^* \mathbf{t}_j \mathbf{t}_j^* \mathbf{W}) \mathbf{\Sigma} \mathbf{V}^*.$$

The update formula enables efficient algorithms for the leave-one-out error estimator and jackknife variance estimator for the approximation  $\mathbf{X}$  and derived quantities like projectors onto singular subspaces and truncation of  $\mathbf{X}$  to rank  $r < s$ . As one example, Algorithm 4.2 gives an implementation the randomized SVD with no subspace iteration ( $q = 0$ ) with the leave-one-out error estimator  $\widehat{\text{Err}}(\mathbf{X}, \mathbf{A})$ . The leave-one-out error estimator requires just two lines and runs in  $\mathcal{O}(s^3)$  operations. Further variants

---

**Algorithm 4.2.** Randomized SVD ( $q = 0$ ) with jackknife variance estimate.

---

**Input:**  $\mathbf{A} \in \mathbb{K}^{d_1 \times d_2}$  to be approximated and approximation rank  $s$

**Output:** Factors  $\mathbf{U} \in \mathbb{K}^{d_1 \times s}$ ,  $\mathbf{\Sigma} \in \mathbb{R}^{s \times s}$ , and  $\mathbf{V} \in \mathbb{K}^{d_2 \times s}$  defining a rank- $s$  approximation  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ , leave-one-out error estimate  $\widehat{\text{Err}} = \widehat{\text{Err}}(\mathbf{X}, \mathbf{A})$

1:  $\mathbf{\Omega} \leftarrow \text{randn}(d_2, s)$

2:  $\mathbf{Y} \leftarrow \mathbf{A}\mathbf{\Omega}$

3:  $(\mathbf{Q}, \mathbf{R}) \leftarrow \text{qr}(\mathbf{Y}, \text{'econ'})$

▷ Economy QR factorization

4:  $\mathbf{C} \leftarrow \mathbf{Q}^* \mathbf{A}$

5:  $(\mathbf{W}, \mathbf{\Sigma}, \mathbf{V}) \leftarrow \text{svd}(\mathbf{C}, \text{'econ'})$

6:  $\mathbf{U} \leftarrow \mathbf{Q}\mathbf{W}$

7:  $\mathbf{G} \leftarrow (\mathbf{R}^*)^{-1}$

8:  $\widehat{\text{Err}} \leftarrow \left( s^{-1} \sum_{j=1}^s \|\mathbf{G}(:, j)\|^2 \right)^{1/2}$

---

are discussed in subsection SM1.3 and a MATLAB implementation is provided in Program SM5.

**5. Numerical experiments.** In this section, we showcase numerical examples that demonstrate the effectiveness of the matrix jackknife variance estimate and leave-one-out error estimate for the Nyström approximation and randomized SVD. All numerical experiments work over the real numbers,  $\mathbb{K} = \mathbb{R}$ .

**5.1. Experimental setup.** To evaluate our diagnostics for matrices with different spectral characteristics, we consider synthetic test matrices from [24, section 5]:

$$\text{(NoisyLR)} \quad \mathbf{A} = \text{diag}(\underbrace{1, \dots, 1}_{R \text{ times}}, 0, \dots, 0) + \xi d^{-1} \mathbf{G}\mathbf{G}^* \in \mathbb{R}^{d \times d}.$$

$$\text{(ExpDecay)} \quad \mathbf{A} = \text{diag}(\underbrace{1, \dots, 1}_{R \text{ times}}, 10^{-q}, 10^{-2q}, \dots, 10^{-(d-R)q}) \in \mathbb{R}^{d \times d}.$$

Here,  $\xi, q \in \mathbb{R}$  are parameters, and  $\mathbf{G} \in \mathbb{R}^{d \times d}$  is a standard Gaussian matrix. Using diagonal test matrices is justified by the observation that the randomized SVD, Nyström approximation, and our diagnostics are orthogonally invariant when  $\mathbf{\Omega}$  is a standard Gaussian matrix, which we use. We also consider matrices from applications:

- **Velocity.** We consider a matrix  $\mathbf{A} \in \mathbb{R}^{25096 \times 1000}$  whose columns are snapshots of the velocity and pressure from simulations of a fluid flow past a cylinder. We thank Beverley McKeon and Sean Symon for this data.
- **Spectral clustering.** We consider the matrix  $\mathbf{A}$  from the spectral clustering example in subsection 1.2.4.

We apply the jackknife variance estimate and leave-one-out error estimate to the randomized Nyström approximation and randomized SVD with a standard Gaussian test matrix  $\mathbf{\Omega}$  for a range of values for the approximation rank  $s$ . For each value of  $s$ , we estimate the mean error  $\text{Err}$  (1.5), standard deviation  $\text{Std}$  (1.8), mean jackknife estimate  $\text{Jack}$ , or mean leave-one-out error estimator  $\widehat{\text{Err}}$  using 1000 independent trials. Error bars on all figures show one standard deviation.

**5.2. Leave-one-out error estimator for randomized SVD.** First, we apply leave-one-out error estimation to estimate the error for the randomized SVD,

$$(5.1) \quad \mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*.$$

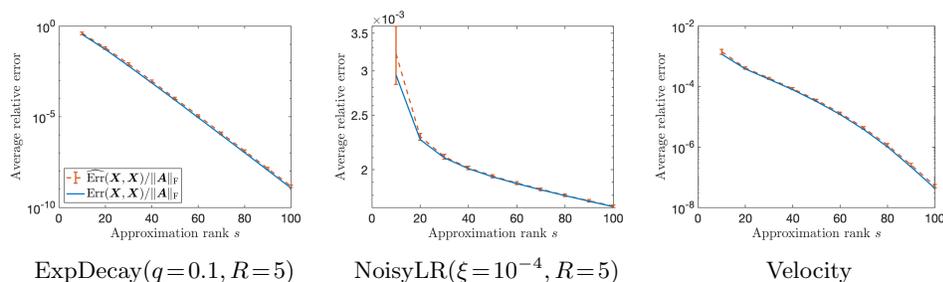


FIG. 3. *Leave-one-out error estimator for randomized SVD.* Error and error estimate for randomized SVD ( $q=0$ ) approximation for three different matrices.

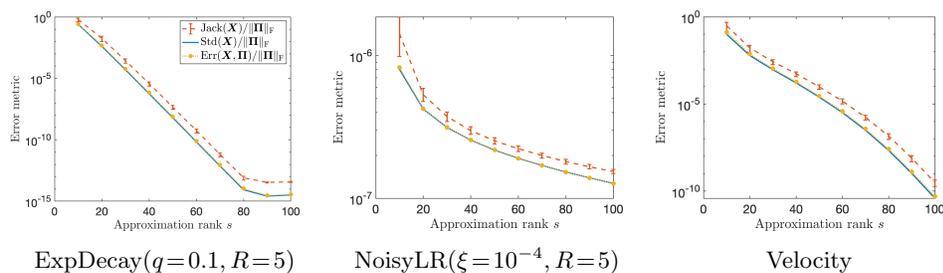


FIG. 4. *Matrix jackknife for projectors onto singular subspaces.* Error, standard deviation, and jackknife estimate for randomized SVD ( $q=0$ ) approximation  $\mathbf{X}$  (5.2) to the projector  $\mathbf{\Pi}$  onto the span of the five dominant right singular vectors for three different matrices.

We set  $q=0$ . Figure 3 shows the results for three examples in the previous section. In all cases, error estimate  $\bar{\text{Err}}(\mathbf{X}, \mathbf{A})$  tracks the true error  $\text{Err}(\mathbf{X}, \mathbf{A})$  closely. Additional examples and analogous plots for randomized Nyström approximation are provided in section SM3.

**5.3. Matrix jackknife for projectors onto singular subspaces.** Consider the task of computing the projector  $\mathbf{\Pi}$  onto the dominant five-dimensional right singular subspace to a matrix  $\mathbf{A}$ . The randomized SVD  $\mathbf{A} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$  yields the approximation

$$(5.2) \quad \mathbf{X} = \mathbf{V}(:, 1:5)\mathbf{V}(:, 1:5)^*.$$

Figure 4 shows the mean error  $\text{Err}(\mathbf{X}, \mathbf{\Pi})$ , standard deviation  $\text{Std}(\mathbf{X})$ , and jackknife estimate  $\text{Jack}(\mathbf{X})$  for the same three test matrices as previously. We again set  $q=0$ . Consistent with Theorem 3.2, the jackknife estimate  $\text{Jack}(\mathbf{X})$  is an *overestimate* of  $\text{Std}(\mathbf{X})$  by a factor of  $2\times$  to  $8\times$ . While the jackknife is not quantitatively sharp, it provides an order of magnitude estimate of the standard deviation and is a useful diagnostic for the quality of the computed output. Additional examples and plots for randomized Nyström approximations of projectors onto invariant subspaces are provided in section SM3.

**5.4. Application: Diagnosing ill-conditioning for spectral clustering.** Jackknife variance estimates can be used to identify situations when a computational task is ill-posed or ill-conditioned. In such cases, refining the approximation (i.e., increasing  $s$  or  $q$ ) may be of little help to improve the quality of the computation.

As an example, consider the spectral clustering application from subsection 1.2. In Nyström-accelerated spectral clustering, we use the dominant  $n_{\text{dim}}$  eigenvectors of a Nyström approximation  $\mathbf{V}\mathbf{A}\mathbf{V}^*$  as coordinates for k-means clustering. For spectral clustering to be reliable, we should pick the parameter  $n_{\text{dim}}$  such that these coordinates are well-conditioned; that is, they should not be highly sensitive to small changes in the normalized kernel matrix  $\mathbf{A}$ . For the example in subsection 1.2.4, the five largest eigenvalues of  $\mathbf{A}$  are

$$\begin{aligned}\lambda_1 &= 0.9999999999999999, \\ \lambda_2 &= 0.99999998639842, \\ \lambda_3 &= 0.999999940523446, \\ \lambda_4 &= 0.999999931126177, \\ \lambda_5 &= 0.997867975285136.\end{aligned}$$

The first four eigenvalues agree up to eight digits of accuracy, with the fifth eigenvalue separated by  $\lambda_4 - \lambda_5 \approx 2 \times 10^{-3}$ . Based on these values, the natural parameter setting would be  $n_{\text{dim}} = 4$ , as the first four eigenvalues are nearly indistinguishable but are well-separated from the fifth.

When we use Nyström-accelerated spectral clustering, we do not have access to the true eigenvalues of the matrix  $\mathbf{A}$  and thus can have difficulties selecting the parameter  $n_{\text{dim}}$  appropriately. Fortunately, the jackknife variance estimate can help warn the user of a poor choice for  $n_{\text{dim}}$ . Let

$$(5.3) \quad \mathbf{X} = \mathbf{V}(:, 1:n_{\text{dim}})\mathbf{V}(:, 1:n_{\text{dim}})^*$$

be the orthoprojector onto the dominant  $n_{\text{dim}}$ -dimensional invariant subspace of the Nyström approximation  $\mathbf{V}\mathbf{A}\mathbf{V}^*$ . Figure 5 shows the standard deviation  $\text{Std}(\mathbf{X})$  and its jackknife estimate  $\text{Jack}(\mathbf{X})$  for both  $n_{\text{dim}} = 3$  and  $n_{\text{dim}} = 4$ . As in subsection 1.2.4,

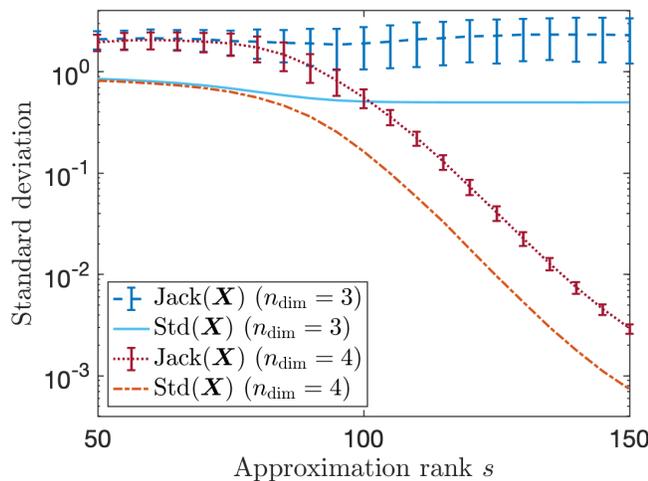


FIG. 5. *Detecting ill-conditioning.* Standard deviation and jackknife estimate for the Nyström spectral projector (5.3) associated with the dominant  $n_{\text{dim}}$ -dimensional invariant subspace for  $n_{\text{dim}} \in \{3, 4\}$  for the spectral clustering matrix.

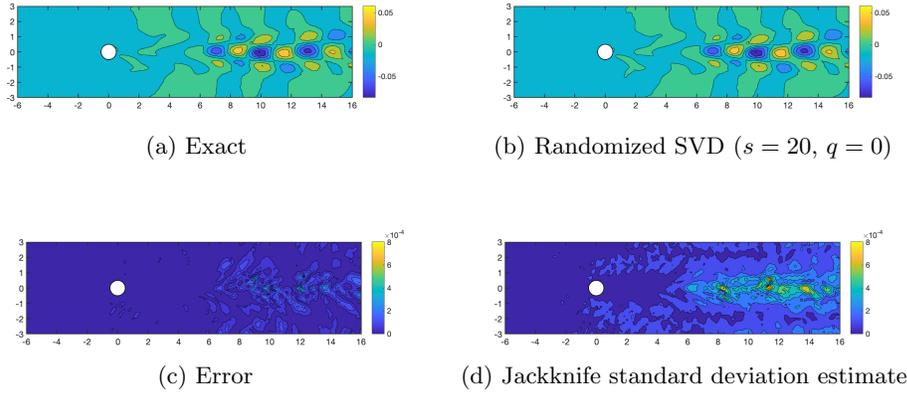


FIG. 6. *Jackknife for singular vector.* These panels assess the entrywise errors in the streamwise velocity from the fifth left singular vector of the velocity test matrix. Panel (a) shows the exact answer, and panel (b) shows the estimate produced by the randomized SVD. Panels (c) and (d) display the error and the Tukey jackknife standard deviation estimate. The jackknife estimate presents a descriptive portrait of where the error is localized.

use  $q = 3$  and  $50 \leq s \leq 150$ . For the good parameter setting  $n_{\text{dim}} = 4$ , the variance decreases sharply as  $s$  is increased. For the bad choice  $n_{\text{dim}} = 3$ , the variance remains persistently high, even as the approximation is refined. This provides evidence to the user that  $\mathbf{X}$  is poorly conditioned and allows the user to fix this by changing the parameter  $n_{\text{dim}}$ .

**5.5. Application: POD modes.** Jackknife variance estimation can be used to give more fine-grained information about a randomized matrix computation. In Figure 6, we compute a (scalar) jackknife variance estimate for the absolute value of each entry of the fifth left singular vector of the velocity matrix computed using the randomized SVD with  $s = 20$  and  $q = 0$ . (The absolute value is introduced to avoid sign ambiguities.)

Left singular vectors of a matrix of simulation data are known as POD modes and are useful for data visualization and model reduction. Since variance is a lower bound on the mean-square error, the coordinatewise jackknife variance estimates can be used as a diagnostic to help identify regions of a POD mode that have high error. This is demonstrated in Figure 6; the jackknife estimate is not quantitatively sharp but paints a descriptive portrait of where the errors are localized.

**6. Extension: Variance estimates for higher Schatten norms.** The variance estimate  $\text{Jack}(\mathbf{X})$  serves as an estimate for the *Frobenius-norm* variance

$$\text{Jack}^2(\mathbf{X}) \approx \text{Var}(\mathbf{X}) = \mathbb{E} \|\mathbf{X} - \mathbb{E} \mathbf{X}\|_F^2.$$

Often, it is more desirable to have error or variance estimates for Schatten norms  $\|\cdot\|_p$  with  $p > 2$ , defined as the  $\ell_p$  norm of the singular values:

$$\|\mathbf{B}\|_p^p := \sum_{j=1}^{\min(d_1, d_2)} \sigma_j^p(\mathbf{B}).$$

One can also construct jackknife estimates for the variance in higher Schatten norms, although the estimates take more intricate forms. For this section, fix an even number

$p \geq 2$  and assume the same setup as subsection 3.2 with the additional stipulation that the samples  $\omega_1, \dots, \omega_s$  take values in a Polish space  $\Omega$ .

The jackknife variance estimate is defined as follows. Consider matrix-valued jackknife *variance proxies*:

$$\widehat{\mathbf{Var}}_1(\mathbf{X}_{s-1}) := \frac{1}{2} \sum_{j=1}^{s-1} \left| \mathbf{X}^{(s)} - \mathbf{X}^{(j)} \right|^2 \quad \text{and} \quad \widehat{\mathbf{Var}}_2(\mathbf{X}_{s-1}) := \frac{1}{2} \sum_{j=1}^{s-1} \left| \left( \mathbf{X}^{(s)} - \mathbf{X}^{(j)} \right)^* \right|^2.$$

Here,  $|\mathbf{B}| = (\mathbf{B}^* \mathbf{B})^{1/2}$  denotes the matrix modulus. Define the Schatten  $p$ -norm variance estimate

$$\text{Jack}_p(\mathbf{X}_{s-1}) := 2^{-\frac{1}{p}} \sqrt{2(p-1)} \left( \left\| \widehat{\mathbf{Var}}_1 \right\|_{p/2}^{p/2} + \left\| \widehat{\mathbf{Var}}_2 \right\|_{p/2}^{p/2} \right)^{1/p}.$$

This quantity seeks to approximate

$$\text{Jack}_p^p(\mathbf{X}_{s-1}) \approx \mathbb{E} \left\| \mathbf{X}_{s-1} - \mathbb{E} \mathbf{X}_{s-1} \right\|_p^p.$$

A matrix generalization of the Efron–Stein–Steele inequality [15, Theorem 4.2] shows this jackknife variance estimate overestimates the Schatten  $p$ -norm variance in the sense

$$\mathbb{E} \left\| \mathbf{X}_{s-1} - \mathbb{E} \mathbf{X}_{s-1} \right\|_p^p \leq \mathbb{E} \text{Jack}_p^p(\mathbf{X}_{s-1}).$$

The techniques we introduce in section 4 can be extended in a natural way to compute  $\text{Jack}_p(\mathbf{X}_{s-1})$  efficiently for the randomized SVD and Nyström approximation.

**Appendix A. Derivation of update formulas.** In this section, we provide proofs of the update formulas (4.1) and (4.2).

**A.1. Proof of (4.1).** Assume without loss of generality that  $j = s$ , instate the notation of subsection 4.1, and assume  $\mathbf{H}$  is invertible. The key ingredient is the following consequence of the Banachiewicz inversion formula [16, equation (0.7.2)]:

$$\begin{bmatrix} (\mathbf{H}^{(s)})^{-1} & \mathbf{0} \\ \mathbf{0}^* & 0 \end{bmatrix} = \mathbf{H}^{-1} - \frac{\mathbf{H}^{-1} \mathbf{e}_s \mathbf{e}_s^* \mathbf{H}^{-1}}{\mathbf{e}_s^* \mathbf{H}^{-1} \mathbf{e}_s}.$$

Using this formula and letting  $\mathbf{R}_{-s}$  denote  $\mathbf{R}$  without its  $s$ th column, we compute

$$\mathbf{X}^{(s)} = \mathbf{Q} \mathbf{R}_{-s} (\mathbf{H}^{(s)})^{-1} \mathbf{R}_{-s}^* \mathbf{Q}^* = \mathbf{X} - \frac{\mathbf{Q} \mathbf{R} \mathbf{H}^{-1} \mathbf{e}_s \mathbf{e}_s^* \mathbf{H}^{-1} \mathbf{R}^* \mathbf{Q}^*}{\mathbf{e}_s^* \mathbf{H}^{-1} \mathbf{e}_s}.$$

Since  $\mathbf{Q} = \mathbf{V} \mathbf{U}^*$  and  $\mathbf{X} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^*$ , we thus have

$$\mathbf{X}^{(s)} = \mathbf{V} \left( \mathbf{\Lambda} - \frac{\mathbf{U}^* \mathbf{R} \mathbf{H}^{-1} \mathbf{e}_s \mathbf{e}_s^* \mathbf{H}^{-1} \mathbf{R}^* \mathbf{Q}^* \mathbf{U}}{\mathbf{e}_s^* \mathbf{H}^{-1} \mathbf{e}_s} \right) \mathbf{V}^* = \mathbf{V} (\mathbf{\Lambda} - \mathbf{t}_s \mathbf{t}_s^*) \mathbf{V}^*,$$

where  $\mathbf{t}_s$  is defined in (4.1b). The formula is established.

**A.2. Proof of (4.2).** Fix  $j \in \{1, \dots, s\}$ , instate the notation of subsection 4.2, and assume  $\mathbf{R}$  is invertible. First observe that

$$\mathbf{Y}^{(j)} = \mathbf{Q} \mathbf{R}_{-j},$$

where  $\mathbf{R}_{-j}$  is  $\mathbf{R}$  without its  $j$ th column. To compute  $\mathbf{Q}^{(j)}$ , we need an economy QR factorization of  $\mathbf{Y}^{(j)} = \mathbf{Q}\mathbf{R}_{-j}$ . To this end, compute a (full) QR decomposition of  $\mathbf{R}_{-j}$ :

$$(A.1) \quad \mathbf{R}_{-j} = [\mathbf{Q}' \quad \mathbf{t}_j] \begin{bmatrix} \mathbf{R}^{(j)} \\ \mathbf{0}^* \end{bmatrix},$$

where  $\mathbf{Q}' \in \mathbb{K}^{s \times (s-1)}$ ,  $\mathbf{R}^{(j)} \in \mathbb{K}^{(s-1) \times (s-1)}$ , and  $\mathbf{t}_j \in \mathbb{K}^s$ . Then

$$\mathbf{Y}^{(j)} = \mathbf{Q}^{(j)}\mathbf{R}^{(j)} \quad \text{for} \quad \mathbf{Q}^{(j)} = \mathbf{Q}\mathbf{Q}'$$

is an economy QR decomposition of  $\mathbf{Y}^{(j)}$ . Since  $[\mathbf{Q}' \quad \mathbf{t}_j]$  is orthogonal, we have

$$\mathbf{I} = [\mathbf{Q}' \quad \mathbf{t}_j] [\mathbf{Q}' \quad \mathbf{t}_j]^* = \mathbf{Q}'(\mathbf{Q}')^* + \mathbf{t}_j\mathbf{t}_j^* \implies \mathbf{Q}'(\mathbf{Q}')^* = \mathbf{I} - \mathbf{t}_j\mathbf{t}_j^*.$$

Thus,

$$\mathbf{Q}^{(j)}(\mathbf{Q}^{(j)})^* = \mathbf{Q}\mathbf{Q}'(\mathbf{Q}')^*\mathbf{Q}^* = \mathbf{Q}(\mathbf{I} - \mathbf{t}_j\mathbf{t}_j^*)\mathbf{Q}^*.$$

Finally, observe that (A.1) implies that  $\mathbf{t}_j$  is orthogonal to the columns of  $\mathbf{R}_{-j}$  so

$$\mathbf{t}_j^*\mathbf{R} \quad \text{is a nonzero multiple of } \mathbf{e}_j^*.$$

Therefore,  $\mathbf{t}_j$  is proportional to the  $j$ th column of  $(\mathbf{R}^*)^{-1}$ . Since  $\mathbf{t}_j$  is a column of an orthogonal matrix, it is equal to a normalized version of the  $j$ th column of  $(\mathbf{R}^*)^{-1}$ .

**Acknowledgment.** We thank Robert Webber for his advice and feedback.

**Disclaimer.** This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government nor any agency thereof, nor any of their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof.

REFERENCES

- [1] D. ARTHUR AND S. VASSILVITSKII, *k-means++: The advantages of careful seeding*, in Proceedings of the 18th Annual ACM-SIAM Symposium on Discrete Algorithms, SIAM, Philadelphia, 2007, pp. 1027–1035.
- [2] P. DRINEAS AND M. W. MAHONEY, *RandNLA: Randomized numerical linear algebra*, Commun. ACM, 59 (2016), pp. 80–90.
- [3] B. EFRON, *The Jackknife, the Bootstrap and Other Resampling Plans*, SIAM, Philadelphia, 1982.
- [4] B. EFRON AND C. STEIN, *The jackknife estimate of variance*, Ann. Statist., 9 (1981), pp. 586–596.
- [5] E. N. EPPERLY, J. A. TROPP, AND R. J. WEBBER, *XTRACE: Making the most out of every sample in stochastic trace estimation*, SIAM J. Matrix Anal. Appl., 45 (2024), pp. 1–23, <https://doi.org/10.1137/23M1548323>.
- [6] A. GITTENS AND M. MAHONEY, *Revisiting the Nyström method for improved large-scale machine learning*, in Proceedings of the 30th International Conference on Machine Learning, PMLR, 2013, pp. 567–575.

- [7] N. HALKO, P.-G. MARTINSSON, AND J. A. TROPP, *Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions*, SIAM Rev., 53 (2011), pp. 217–288.
- [8] H. LI, G. C. LINDERMAN, A. SZLAM, K. P. STANTON, Y. KLUGER, AND M. TYGERT, *Algorithm 971: An implementation of a randomized algorithm for principal component analysis*, ACM Trans. Math. Software, 43 (2017), 28.
- [9] M. LOPES, N. B. ERICHSON, AND M. MAHONEY, *Error estimation for sketched SVD via the bootstrap*, in Proceedings of the 37th International Conference on Machine Learning, Proc. Mach. Learn. Res. 119, PMLR, 2020, pp. 6382–6392.
- [10] M. E. LOPES, N. B. ERICHSON, AND M. W. MAHONEY, *Bootstrapping the operator norm in high dimensions: Error estimation for covariance matrices and sketching*, Bernoulli, 29 (2023), pp. 428–450.
- [11] M. E. LOPES, S. WANG, AND M. W. MAHONEY, *A bootstrap method for error estimation in randomized matrix multiplication*, J. Mach. Learn. Res., 20 (2019), 39.
- [12] P.-G. MARTINSSON AND J. A. TROPP, *Randomized numerical linear algebra: Foundations and algorithms*, Acta Numer., 29 (2020), pp. 403–572.
- [13] R. MURRAY, J. DEMMEL, M. W. MAHONEY, N. B. ERICHSON, M. MELNICHENKO, O. A. MALIK, L. GRIGORI, P. LUSZCZEK, M. DEREZINSKI, M. E. LOPES, T. LIANG, H. LUO, AND J. DONGARRA, *Randomized Numerical Linear Algebra: A Perspective on the Field with an Eye to Software*, <http://arxiv.org/abs/2302.11474>, 2023.
- [14] C. MUSCO AND C. MUSCO, *Randomized block Krylov methods for stronger and faster approximate singular value decomposition*, in Proceedings of the 28th International Conference on Neural Information Processing Systems, Vol. 1, MIT Press, Cambridge, MA, 2015, pp. 1396–1404.
- [15] D. PAULIN, L. MACKEY, AND J. A. TROPP, *Efron–Stein inequalities for random matrices*, Ann. Probab., 44 (2016), pp. 3431–3473.
- [16] S. PUNTANEN AND G. P. H. STYAN, *Historical introduction: Issai Schur and the early development of the Schur complement*, in The Schur Complement and Its Applications, Numer. Methods Algorithms 4, F. Zhang, ed., Springer, New York, 2005, pp. 1–16.
- [17] M. H. QUENOUILLE, *Approximate tests of correlation in time-series*, J. R. Stat. Soc. Ser. B. Methodol., 11 (1949), pp. 68–84.
- [18] R. RAMAKRISHNAN, P. O. DRAL, M. RUPP, AND O. A. VON LILIEFELD, *Quantum chemistry structures and properties of 134 kilo molecules*, Scientific Data, 1 (2014), 140022.
- [19] L. RUDDIGKEIT, R. VAN DEURSEN, L. C. BLUM, AND J.-L. REYMOND, *Enumeration of 166 billion organic small molecules in the chemical universe database GDB-17*, J. Chem. Inf. Model., 52 (2012), pp. 2864–2875.
- [20] Y. SAAD, *Numerical Methods for Large Eigenvalue Problems*, Manchester University Press, Manchester, UK, 1992.
- [21] J. M. STEELE, *An Efron–Stein inequality for nonsymmetric statistics*, Ann. Statist., 14 (1986), pp. 753–758.
- [22] R. J. TIBSHIRANI AND B. EFRON, *An Introduction to the Bootstrap*, Monogr. Statist. Appl. Probab. 57, CRC Press, Boca Raton, FL, 1993.
- [23] J. A. TROPP AND R. J. WEBBER, *Randomized Algorithms for Low-Rank Matrix Approximation: Design, Analysis, and Applications*, 2023, <https://doi.org/10.48550/arXiv.2306.12418>.
- [24] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *Fixed-rank approximation of a positive-semidefinite matrix from streaming data*, in Advances in Neural Information Processing Systems, Vol. 30, 2017, pp. 1225–1234.
- [25] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *Streaming low-rank matrix approximation with an application to scientific simulation*, SIAM J. Sci. Comput., 41 (2019), pp. A2430–A2463.
- [26] J. TUKEY, *Bias and confidence in not quite large samples*, Ann. Math. Statist., 29 (1958), p. 614.
- [27] U. VON LUXBURG, *A tutorial on spectral clustering*, Stat. Comput., 17 (2007), pp. 395–416.
- [28] C. K. I. WILLIAMS AND M. SEEGER, *Using the Nyström method to speed up kernel machines*, in Proceedings of the 13th International Conference on Neural Information Processing Systems, MIT Press, Cambridge, MA, 2000, pp. 661–667.
- [29] D. P. WOODRUFF, *Sketching as a tool for numerical linear algebra*, Found. Trends Theoret. Comput. Sci., 10 (2014), pp. 1–157.
- [30] J. YAO, N. B. ERICHSON, AND M. E. LOPES, *Error estimation for random Fourier features*, in Proceedings of the 26th International Conference on Artificial Intelligence and Statistics, Proc. Mach. Learn. Res. 206, PMLR, 2023, pp. 2348–2364.