

XTRACE: MAKING THE MOST OF EVERY SAMPLE IN STOCHASTIC TRACE ESTIMATION*

ETHAN N. EPPERLY[†], JOEL A. TROPP[†], AND ROBERT J. WEBBER[†]

Abstract. The implicit trace estimation problem asks for an approximation of the trace of a square matrix, accessed via matrix–vector products (matvecs). This paper designs new randomized algorithms, XTRACE and XNYSTRACE, for the trace estimation problem by exploiting both variance reduction and the exchangeability principle. For a fixed budget of matvecs, numerical experiments show that the new methods can achieve errors that are orders of magnitude smaller than existing algorithms, such as the Girard–Hutchinson estimator or the HUTCH++ estimator. A theoretical analysis confirms the benefits by offering a precise description of the performance of these algorithms as a function of the spectrum of the input matrix. The paper also develops an exchangeable estimator, XDIAG, for approximating the diagonal of a square matrix using matvecs.

Key words. Trace estimation, low-rank approximation, exchangeability, variance reduction, randomized algorithm.

AMS subject classifications. 65C05, 65F30, 68W20

1. Introduction. Over the past three decades, researchers have developed *randomized algorithms* for linear algebra problems such as trace estimation [14, 19, 26], low-rank approximation [15], and over-determined least squares [3, 30]. Many of these algorithms collect information by judicious random sampling of the problem data. As a consequence, we can design better algorithms using techniques from the theory of statistical estimation, such as variance reduction and the exchangeability principle. This paper explores how the exchangeability principle leads to faster randomized algorithms for trace estimation.

Suppose that we wish to compute a quantity $Q(\mathbf{A})$ associated with a matrix \mathbf{A} . A typical randomized algorithm might proceed as follows.

1. **Collect information** about the matrix \mathbf{A} by computing matrix–vector products $\mathbf{A}\boldsymbol{\omega}_1, \dots, \mathbf{A}\boldsymbol{\omega}_k$ with random test vectors $\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_k$.
2. **Form an estimate** of $Q(\mathbf{A})$ from the samples $\mathbf{A}\boldsymbol{\omega}_1, \dots, \mathbf{A}\boldsymbol{\omega}_k$.

The question arises: *Given the data $\mathbf{A}\boldsymbol{\omega}_1, \dots, \mathbf{A}\boldsymbol{\omega}_k$, what is an optimal estimator for $Q(\mathbf{A})$?* One property an optimal estimator must obey is the exchangeability principle:

Exchangeability principle: If the test vectors $\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_k$ are exchangeable, the minimum-variance unbiased estimator for $Q(\mathbf{A})$ is always a symmetric function of $\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_k$.

*Date: 18 January 2023.

Funding: ENE acknowledges support from the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research, Department of Energy Computational Science Graduate Fellowship under Award Number DE-SC0021110. JAT and RJW acknowledge support from the Office of Naval Research through BRC Award N00014-18-1-2363 and from the National Science Foundation through FRG Award 1952777.

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.

[†]Division of Computing and Mathematical Sciences, California Institute of Technology, Pasadena, CA 91125 USA (epperly@caltech.edu, jtropp@caltech.edu, rwebber@caltech.edu).

“Exchangeability” means that the family $(\omega_1, \dots, \omega_k)$ has the same distribution as the permuted family $(\omega_{\sigma(1)}, \dots, \omega_{\sigma(k)})$ for every permutation σ in the symmetric group S_k . In particular, an independent and identically distributed (iid) family is exchangeable.

The implication of the exchangeability principle is that our estimators should be symmetric functions of the samples, whenever possible. This idea is attributed to Halmos [16], and it plays a central role in the theory of U-statistics [21].

This paper will demonstrate that the exchangeability principle can lead to new randomized algorithms for linear algebra problems. As a case study, we will explore the problem of *implicit trace estimation*:

Implicit trace estimation problem: Given access to a square matrix A via the matrix–vector product (matvec) operation $\omega \mapsto A\omega$, estimate the trace of A .

Trace estimation plays a role in a wide range of areas, including computational statistics, statistical mechanics, and network analysis. See the survey [35] for more applications.

As we will see, it is natural to design randomized algorithms for trace estimation that use matvecs between the input matrix and random test vectors. At present, the state-of-the-art trace estimators do not satisfy the exchangeability principle. By pursuing this insight, we will develop better trace estimators. Given a fixed budget of matvecs, the new algorithms can reduce the variance of the trace estimate by several orders of magnitude. This case study highlights the importance of enforcing exchangeability in design of randomized algorithms.

1.1. Stochastic trace estimators. In this section, we outline the classic approach to randomized trace estimation based on Monte Carlo approximation. Then we introduce a more modern approach that incorporates a variance reduction strategy.

1.1.1. The Girard–Hutchinson estimator. The first randomized algorithm for trace estimation was proposed by Girard [14] and extended by Hutchinson [19].

Let $A \in \mathbb{R}^{N \times N}$ be a square input matrix. Consider an isotropic random vector $\omega \in \mathbb{R}^N$:

$$(1.1) \quad \mathbb{E}[\omega \omega^*] = I.$$

For example, we may take a random sign vector $\omega \sim \text{UNIFORM}\{\pm 1\}^N$. By isotropy,

$$\mathbb{E}[\omega^* (A\omega)] = \text{tr } A.$$

The symbol $*$ denotes the transpose. This relation suggests a Monte Carlo method.

Accordingly, the Girard–Hutchinson trace estimator takes the form

$$(1.2) \quad \hat{\text{tr}}_{\text{GH}} := \frac{1}{m} \sum_{i=1}^m \omega_i^* (A\omega_i) \quad \text{where the } \omega_i \text{ are iid copies of } \omega.$$

This estimator is exchangeable, and it is unbiased: $\mathbb{E}[\hat{\text{tr}}_{\text{GH}}] = \text{tr } A$. We can measure the quality of the estimator using the variance, $\text{Var}[\hat{\text{tr}}_{\text{GH}}]$. The variance depends on the matrix A and the distribution of ω , but it converges to zero at the Monte Carlo rate $\Theta(m^{-1})$ as we increase the number m of samples. See the survey [25, §4] for more discussion.

1.1.2. The HUTCH++ estimator. To improve on the Girard–Hutchinson estimator, several papers [12, 23, 26, 31] have advocated variance reduction techniques. The key idea is to form a low-rank approximation of the input matrix. We can compute the trace of the approximation exactly (as a control variate), so we only need to estimate the trace of the residual. This approach can attain lower variance than the Monte Carlo method.

DRAFT

Algorithm 1.1 HUTCH++ [26]**Input:** Matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ and number m of matvecs, where m is divisible by 3**Output:** Trace estimate $\hat{\text{tr}} \approx \text{tr } \mathbf{A}$

- 1: Draw iid isotropic $\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_{2m/3} \in \mathbb{R}^N$ ▷ For example, $\boldsymbol{\omega}_i \sim \text{UNIFORM}\{\pm 1\}^N$
- 2: $\mathbf{Y} \leftarrow \mathbf{A} [\boldsymbol{\omega}_{m/3+1} \ \cdots \ \boldsymbol{\omega}_{2m/3}]$ ▷ Use matvecs
- 3: $\mathbf{Q} \leftarrow \text{orth}(\mathbf{Y})$
- 4: $\mathbf{G} \leftarrow [\boldsymbol{\omega}_1 \ \cdots \ \boldsymbol{\omega}_{m/3}] - \mathbf{Q}\mathbf{Q}^* [\boldsymbol{\omega}_1 \ \cdots \ \boldsymbol{\omega}_{m/3}]$
- 5: $\hat{\text{tr}} \leftarrow \text{tr}(\mathbf{Q}^*(\mathbf{A}\mathbf{Q})) + (m/3)^{-1} \text{tr}(\mathbf{G}^*(\mathbf{A}\mathbf{G}))$ ▷ Use matvecs

The HUTCH++ estimator of Meyer, Musco, Musco, and Woodruff [26] crystallizes the variance reduction strategy. Let $\mathbf{A} \in \mathbb{R}^{N \times N}$ be a square input matrix. Given a fixed budget of m matvecs, with m divisible by 3, HUTCH++ proceeds as follows:

1. *Sample* iid isotropic vectors $\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_{2m/3} \in \mathbb{R}^N$ as in (1.1).
2. *Sketch* $\mathbf{Y} = \mathbf{A} [\boldsymbol{\omega}_{m/3+1} \ \boldsymbol{\omega}_{m/3+2} \ \cdots \ \boldsymbol{\omega}_{2m/3}]$
3. *Orthonormalize* $\mathbf{Q} = \text{orth}(\mathbf{Y})$.
4. *Output* the estimate

$$(1.3) \quad \hat{\text{tr}}_{\text{H++}} := \text{tr}(\mathbf{Q}^*(\mathbf{A}\mathbf{Q})) + \frac{1}{m/3} \sum_{i=1}^{m/3} \boldsymbol{\omega}_i^* (\mathbf{I} - \mathbf{Q}\mathbf{Q}^*) (\mathbf{A}(\mathbf{I} - \mathbf{Q}\mathbf{Q}^*) \boldsymbol{\omega}_i).$$

See [Algorithm 1.1](#) for efficient HUTCH++ pseudocode.

To illustrate how HUTCH++ takes advantage of low-rank approximation, we first observe that $\hat{\mathbf{A}} = \mathbf{Q}\mathbf{Q}^*\mathbf{A}$ is a low-rank approximation of the matrix \mathbf{A} . Indeed, the matrix $\hat{\mathbf{A}}$ coincides with the randomized SVD [15] formed from the test matrix $[\boldsymbol{\omega}_{m/3+1} \ \cdots \ \boldsymbol{\omega}_{2m/3}]$. HUTCH++ computes the trace of the low-rank approximation, which is

$$\text{tr } \hat{\mathbf{A}} = \text{tr}(\mathbf{Q}\mathbf{Q}^*\mathbf{A}) = \text{tr}(\mathbf{Q}^*\mathbf{A}\mathbf{Q}).$$

Afterward, HUTCH++ applies the Girard–Hutchinson estimator to estimate the trace of the residual

$$\text{tr}(\mathbf{A} - \hat{\mathbf{A}}) = \text{tr}((\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)\mathbf{A}) = \text{tr}((\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)\mathbf{A}(\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)).$$

Like the Girard–Hutchinson trace estimator, the HUTCH++ estimator is unbiased. In contrast to the $\Theta(m^{-1})$ variance of Girard–Hutchinson, the variance of HUTCH++ is no greater than $\mathcal{O}(m^{-2})$. In practice, the reduction in variance is conspicuous. On the other hand, the HUTCH++ estimator violates the exchangeability principle, so we recognize an opportunity to design a better algorithm.

1.2. New exchangeable trace estimators. The HUTCH++ estimator is not exchangeable because it uses some test vectors to perform low-rank approximation, while it uses other test vectors to estimate the trace of the residual. Although it might seem natural to symmetrize HUTCH++ over all splits of the test vectors, this approach is both wasteful and computationally infeasible.

Instead, we will develop a new family of exchangeable trace estimators that make the most out of every test vector. We will use *all* of the test vectors for low-rank approximation, and we will use *all* of the test vectors for estimating the trace of the residual. The key to this strategy is a leave-one-out technique that can be implemented at the same computational cost as HUTCH++. This innovation can reduce the variance by several orders of magnitude.

DRAFT

Algorithm 1.2 XTRACE: Naïve implementation**Input:** Matrix $A \in \mathbb{R}^{N \times N}$ and number m of matvecs, where m is even**Output:** Trace estimate $\widehat{\text{tr}} \approx \text{tr } A$ and error estimate $\widehat{\text{er}} \approx |\widehat{\text{tr}} - \text{tr } A|$

- 1: Draw iid isotropic $\omega_1, \dots, \omega_{m/2} \in \mathbb{R}^N$ ▷ See [subsection 2.3](#)
- 2: $Y \leftarrow A[\omega_1 \ \dots \ \omega_{m/2}]$ ▷ Use matvecs
- 3: **for** $i = 1, 2, 3, \dots, m/2$ **do**
- 4: $Q_{(i)} \leftarrow \text{orth}(Y_{-i})$ ▷ Remove i th column of Y
- 5: $\widehat{\text{tr}}_i \leftarrow \text{tr}(Q_{(i)}^*(A Q_{(i)})) + \omega_i^*(\mathbf{I} - Q_{(i)} Q_{(i)}^*)(A(\mathbf{I} - Q_{(i)} Q_{(i)}^*))\omega_i$ ▷ Use matvecs
- 6: **end for**
- 7: $\widehat{\text{tr}} \leftarrow (m/2)^{-1} \sum_{i=1}^{m/2} \widehat{\text{tr}}_i$
- 8: $\widehat{\text{er}}^2 \leftarrow ((m/2)(m/2 - 1))^{-1} \sum_{i=1}^{m/2} (\widehat{\text{tr}}_i - \widehat{\text{tr}})^2$

1.2.1. The XTRACE estimator. Our first method, called XTRACE, is an exchangeable trace estimator designed for general square matrices. It computes a family of variance-reduced trace estimators. Each estimator uses all but one test vector to form a low-rank approximation, and it uses the remaining test vector to estimate the trace of the residual. XTRACE then averages the basic estimators together to obtain an exchangeable trace estimator.

Let us give a more detailed description. Fix a square input matrix $A \in \mathbb{R}^{N \times N}$. The parameter m is the number of matvecs, where m is an even number. Draw an iid family $\omega_1, \dots, \omega_{m/2} \in \mathbb{R}^N$ of isotropic test vectors, and define the test matrix

$$\Omega = [\omega_1 \ \omega_2 \ \omega_3 \ \dots \ \omega_{m/2}].$$

Construct the orthonormal matrices

$$(1.4) \quad Q_{(i)} = \text{orth}(A\Omega_{-i}) \quad \text{for each } i = 1, \dots, m/2,$$

where Ω_{-i} is the test matrix with the i th column removed. Compute the basic trace estimators

$$(1.5) \quad \widehat{\text{tr}}_i := \text{tr}(Q_{(i)}^*(A Q_{(i)})) + \omega_i^*(\mathbf{I} - Q_{(i)} Q_{(i)}^*)(A(\mathbf{I} - Q_{(i)} Q_{(i)}^*))\omega_i \quad \text{for } i = 1, \dots, m/2.$$

The XTRACE estimator averages these basic estimators:

$$(1.6) \quad \widehat{\text{tr}}_X := \frac{1}{m/2} \sum_{i=1}^{m/2} \widehat{\text{tr}}_i.$$

The XTRACE method gives an unbiased, exchangeable estimate for the trace. [Theorem 1.1](#) provides a detailed prior bound for the variance. We can also obtain a posterior estimate for the error using the formula

$$\widehat{\text{er}}_X^2 := \frac{1}{(m/2)(m/2 - 1)} \sum_{i=1}^{m/2} (\widehat{\text{tr}}_i - \widehat{\text{tr}}_X)^2.$$

[Subsection 3.1](#) contains further discussion of the error estimate.

This procedure requires exactly m matvecs with the input matrix A . See [Algorithm 1.2](#) for direct XTRACE pseudocode. With careful attention to the linear algebra, we can develop an implementation whose computational cost is comparable with a single estimator of the form (1.5); see [subsection 2.1](#) for details.

DRAFT

Algorithm 1.3 XNYS TRACE: Naïve implementation**Input:** Psd matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ and number m of matvecs**Output:** Trace estimate $\widehat{\text{tr}} \approx \text{tr } \mathbf{A}$ and error estimate $\widehat{\text{er}} \approx |\widehat{\text{tr}} - \text{tr } \mathbf{A}|$

- 1: Draw iid isotropic $\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_m \in \mathbb{R}^N$ ▷ See [subsection 2.3](#)
- 2: $\boldsymbol{\Omega} \leftarrow [\boldsymbol{\omega}_1 \ \dots \ \boldsymbol{\omega}_m]$
- 3: $\mathbf{Y} \leftarrow \mathbf{A}\boldsymbol{\Omega}$ ▷ Use matvecs
- 4: **for** $i = 1, 2, 3, \dots, m$ **do**
- 5: $\widehat{\mathbf{A}}_i \leftarrow \mathbf{Y}_{-i}(\boldsymbol{\Omega}_{-i}^* \mathbf{Y}_{-i})^\dagger \mathbf{Y}_{-i}^*$ ▷ Remove i th column of \mathbf{Y} and $\boldsymbol{\Omega}$
- 6: $\widehat{\text{tr}}_i \leftarrow \text{tr } \widehat{\mathbf{A}}_i + \boldsymbol{\omega}_i^* ((\mathbf{A} - \widehat{\mathbf{A}}_i) \boldsymbol{\omega}_i)$ ▷ Use matvecs
- 7: **end for**
- 8: $\widehat{\text{tr}} \leftarrow m^{-1} \sum_{i=1}^m \widehat{\text{tr}}_i$
- 9: $\widehat{\text{er}}^2 \leftarrow (m(m-1))^{-1} \sum_{i=1}^m (\widehat{\text{tr}}_i - \widehat{\text{tr}})^2$

1.2.2. The XNYS TRACE estimator. Our second method, called XNYS TRACE, is an exchangeable trace estimator designed for positive-semidefinite (psd) matrices. Rather than using a randomized SVD to reduce the variance, this estimator uses a Nyström approximation [25, §14] of the psd matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$. This approximation takes the form

$$(1.7) \quad \mathbf{A}(\mathbf{X}) := \mathbf{A}\mathbf{X}(\mathbf{X}^* \mathbf{A}\mathbf{X})^\dagger (\mathbf{A}\mathbf{X})^* \quad \text{for a test matrix } \mathbf{X} \in \mathbb{R}^{N \times k}.$$

The Nyström method requires only k matvecs to compute a rank- k approximation, while the randomized SVD requires $2k$ matvecs.

Let us summarize the XNYS TRACE method. Draw iid isotropic test vectors $\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_m$, and form the test matrix $\boldsymbol{\Omega} = [\boldsymbol{\omega}_1 \ \dots \ \boldsymbol{\omega}_m]$. The basic estimators take the form

$$(1.8) \quad \widehat{\text{tr}}_i := \text{tr } \mathbf{A} \langle \boldsymbol{\Omega}_{-i} \rangle + \boldsymbol{\omega}_i^* ((\mathbf{A} - \mathbf{A} \langle \boldsymbol{\Omega}_{-i} \rangle) \boldsymbol{\omega}_i) \quad \text{for } i = 1, \dots, m.$$

As usual, $\boldsymbol{\Omega}_{-i}$ denotes the test matrix with the i th column removed. To obtain the XNYS TRACE estimator and an error estimate, we use the formulas

$$(1.9) \quad \widehat{\text{tr}}_{\text{XN}} := \frac{1}{m} \sum_{i=1}^m \widehat{\text{tr}}_i \quad \text{and} \quad \widehat{\text{er}}_{\text{XN}} := \frac{1}{m(m-1)} \sum_{i=1}^m (\widehat{\text{tr}}_i - \widehat{\text{tr}}_{\text{XN}})^2.$$

The XNYS TRACE estimator is unbiased and exchangeable. [Theorem 1.1](#) provides a bound for the variance. See [Algorithm 1.3](#) for naïve XNYS TRACE pseudocode and [subsection 2.2](#) for a more efficient approach.

The recent paper [28] describes an estimator called NYSTRÖM++ that uses a Nyström approximation to perform reduced-variance trace estimation. The NYSTRÖM++ method violates the exchangeability principle, while XNYS TRACE repairs this weakness.

1.2.3. Stochastic diagonal estimators. As an extension of XTRACE, we also propose the XDIAG algorithm for estimating the *diagonal* of an implicitly defined matrix. We will discuss this approach in [subsection 2.4](#).

1.3. Numerical experiments. To highlight the advantages of the XTRACE and XNYS TRACE estimators, we present some motivating numerical experiments. [Section 4](#) contains further numerical work.

1.3.1. Exploiting spectral decay. Our first experiment uses a synthetic input matrix to illustrate how the exchangeable estimators wring more information out of the samples.

We apply several trace estimators to a psd matrix with exponentially decreasing eigenvalues; see [subsection 4.1](#) for the details of the matrix. [Figure 1](#) reports the average error

DRAFT

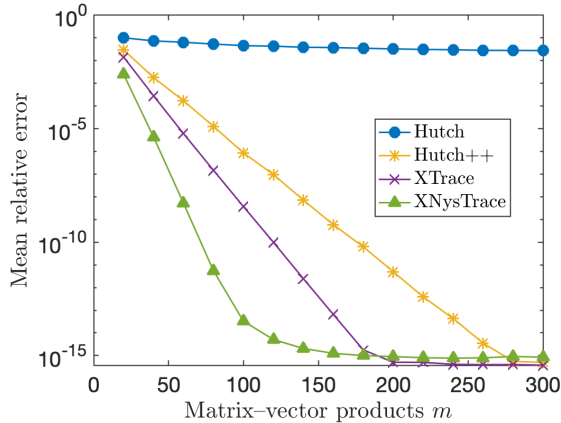


Fig. 1: **Exploiting spectral decay.** Average error of trace estimators applied to a (synthetic) psd matrix with exponentially decreasing eigenvalues. See [subsection 1.3.1](#).

over 1000 trials. The Girard–Hutchinson estimator (HUTCH) converges at the Monte Carlo rate, whereas the newer estimators all converge much faster. This improvement comes from variance reduction techniques that exploit the spectral decay. Observe that XTRACE and XNYS TRACE converge exponentially fast at $1.5\times$ and $3\times$ the rate of HUTCH++, until reaching machine precision. For a fixed budget of m matvecs, XTRACE and XNYS TRACE can reduce the error by several orders of magnitude compared to HUTCH++. Strikingly, the reduction in variance from enforcing exchangeability is almost as significant as the reduction in variance from using a low-rank approximation as a control variate.

1.3.2. Computing partition functions. Our second experiment shows how the advantages of using exchangeable estimators persist in a scientific application.

We apply several trace estimators to compute the partition function for a quantum system

$$Z(\beta) := \text{tr} \exp(-\beta \mathbf{H}),$$

where \mathbf{H} is a symmetric Hamiltonian matrix and $\beta > 0$ is an inverse temperature. Specifically, we consider the Hamiltonian matrix \mathbf{H} for the transverse-field Ising model on 18 sites, which has dimension $N = 2^{18} = 262\,144$. See [subsection 4.3](#) for details on the matrix \mathbf{H} and the partition function $Z(\beta)$. To evaluate matvecs with $\exp(-\beta \mathbf{H})$, we can use the code of Higham [18] that implements an adaptive polynomial approximation [2].

Figure 2 reports the mean estimation error over 100 trials. With just $m = 10$ matvecs, all the variance-reduced methods achieve errors that are 5 orders of magnitude smaller than the Girard–Hutchinson estimator. Furthermore, XTRACE and XNYS TRACE are faster and more precise than HUTCH++. For example, with $m = 40$ matvecs, **XTRACE is 240× more accurate**, and **XNYS TRACE is 2400× more accurate**.

1.4. Theoretical guarantees. To explain the excellent performance of the exchangeable estimators, we establish detailed theoretical guarantees. For theoretical convenience, our analysis uses standard normal test vectors. As a consequence, we can deliver explicit constants that allow us to make meaningful comparisons among the methods.

THEOREM 1.1 (Variance bounds). *Let $\mathbf{A} \in \mathbb{R}^{N \times N}$ be a square matrix. Fix the number*

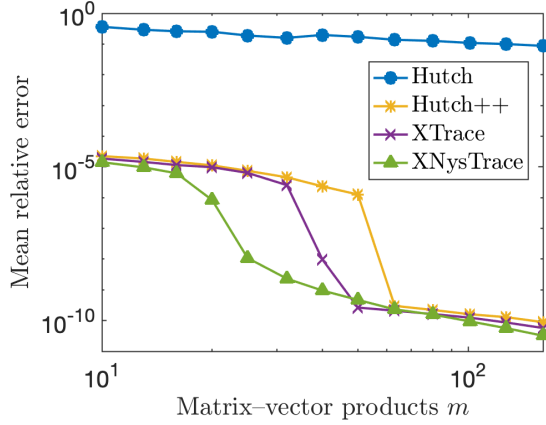


Fig. 2: **Computing a partition function.** Average error of trace estimators applied to the partition function of the transverse-field Ising model. See [subsection 1.3.2](#).

of matvecs: $m \geq 8$. The HUTCH++, XTRACE, and XNYS TRACE estimators are unbiased estimators of the trace. With standard normal test vectors $\omega \sim \text{NORMAL}(\mathbf{0}, \mathbf{I})$, these estimators satisfy the variance bounds

$$\begin{aligned} (\mathbb{E}|\widehat{\text{tr}}_{\text{Hutch++}} - \text{tr } \mathbf{A}|^2)^{1/2} &\leq \min_{r \leq m/3-2} \left(\sqrt{2} \frac{\|\mathbf{A} - [\mathbf{A}]_r\|_{\text{F}}}{\sqrt{m/3-r-1}} \right); \\ (\mathbb{E}|\widehat{\text{tr}}_{\text{XTrace}} - \text{tr } \mathbf{A}|^2)^{1/2} &\leq \sqrt{m} \min_{r \leq m/2-4} \left(2 \frac{\|\mathbf{A} - [\mathbf{A}]_r\|}{\sqrt{m/2-r-3}} + 2e \frac{\|\mathbf{A} - [\mathbf{A}]_r\|_{\text{F}}}{m/2-r-3} \right); \\ (\mathbb{E}|\widehat{\text{tr}}_{\text{XNysTrace}} - \text{tr } \mathbf{A}|^2)^{1/2} &\leq m \min_{r \leq m-6} \left(\sqrt{8} \frac{\|\mathbf{A} - [\mathbf{A}]_r\|}{m-r-5} + \sqrt{2} \frac{\|\mathbf{A} - [\mathbf{A}]_r\|_{\text{F}}}{(m-r-5)^{3/2}} + 5e^2 \frac{\|\mathbf{A} - [\mathbf{A}]_r\|_{*}}{(m-r-5)^2} \right). \end{aligned}$$

These formulas involve the Frobenius norm $\|\cdot\|_{\text{F}}$, the spectral norm $\|\cdot\|$, and the trace norm $\|\cdot\|_{*}$. The matrix $[\mathbf{A}]_r$ is a simultaneous best rank- r approximation of \mathbf{A} in these norms.

In addition, for each of these three methods, it suffices to use $m = \mathcal{O}(\eta^{-1/2})$ matvecs to achieve the variance bound

$$(1.10) \quad \text{Var}[\widehat{\text{tr}}] = \mathbb{E}|\widehat{\text{tr}} - \text{tr } \mathbf{A}|^2 \leq \eta \|\mathbf{A}\|_{*}^2 \quad \text{for } \eta \in (0, 1).$$

The proof of [Theorem 1.1](#) appears in [section 5](#).

As the number m of matvecs increases, [Theorem 1.1](#) ensures that the variance of XTRACE, XNYS TRACE, and HUTCH++ decreases at a rate of $\mathcal{O}(1/m^2)$. Therefore, these algorithms are all superior to the Girard–Hutchinson estimator, whose variance decreases at the Monte Carlo rate $\Theta(1/m)$.

[Theorem 1.1](#) also demonstrates the advantage of XTRACE and XNYS TRACE for matrices whose singular values decay rapidly. This benefit is visible from the error bounds because they allow for larger values r of the approximation rank. As an example, consider a psd matrix \mathbf{A} whose eigenvalues have exponential decay with rate $\alpha \in (0, 1)$:

$$\lambda_i(\mathbf{A}) \leq \alpha^i \quad \text{for } i = 1, 2, 3, \dots$$

DRAFT

The errors of HUTCH++, XTRACE, and XNYSTRACE decay like

$$\begin{aligned} (\mathbb{E}|\widehat{\text{tr}}_{\text{H}++} - \text{tr } \mathbf{A}|^2)^{1/2} &\leq C_1(\alpha) \alpha^{m/3}; \\ (\mathbb{E}|\widehat{\text{tr}}_{\text{X}} - \text{tr } \mathbf{A}|^2)^{1/2} &\leq \sqrt{m} C_2(\alpha) \alpha^{m/2}; \\ (\mathbb{E}|\widehat{\text{tr}}_{\text{XN}} - \text{tr } \mathbf{A}|^2)^{1/2} &\leq m C_3(\alpha) \alpha^m. \end{aligned}$$

For this class of matrix, XTRACE converges exponentially fast at $1.5\times$ the rate of HUTCH++, and XNYSTRACE converges exponentially fast at $3\times$ the rate of HUTCH++. This is precisely the behavior we observe in [Figure 1](#).

1.5. Benefits. In summary, the XTRACE and XNYSTRACE estimators have several desirable features as compared with previous approaches.

1. **Higher accuracy:** For a fixed budget of m matvecs, XTRACE and XNYSTRACE often yield errors that are *orders of magnitude smaller* than HUTCH++.
2. **Efficient algorithms:** We have designed optimized implementations of XTRACE and XNYSTRACE that only require m matvecs plus $\mathcal{O}(m^2N)$ arithmetic operations, which is the same computational cost as HUTCH++.
3. **Error estimation:** We can equip the XTRACE and XNYSTRACE estimators with reliable error estimates.

Altogether, these benefits make a compelling case that XTRACE and XNYSTRACE should be the algorithms of choice for trace estimation.

1.6. A brief history of stochastic trace estimation. Girard wrote the first paper [14] on stochastic trace estimation, in which he proposed the estimator (1.2) with test vectors drawn uniformly from a Euclidean sphere. His goal was to develop an efficient way to perform generalized cross-validation for smoothing splines. Hutchinson [19] suggested using random sign vectors instead: $\boldsymbol{\omega} \sim \text{UNIFORM}\{\pm 1\}^N$. See [25, §4] for further details.

In the last five years, researchers have developed far more efficient methods for trace estimation by incorporating variance reduction techniques. In 2017, Saibaba, Alexandrian, and Ipsen [31] proposed a biased estimator that outputs the trace of a low-rank approximation as a surrogate for $\text{tr } \mathbf{A}$. Around the same time, Gambhir, Stathopoulos, and Orginos [12] proposed a hybrid estimator, similar to HUTCH++, that outputs the trace of a low-rank approximation, $\text{tr } \widehat{\mathbf{A}}$, plus a Girard–Hutchinson estimate for $\text{tr}(\mathbf{A} - \widehat{\mathbf{A}})$. The paper [23] of Lin contains related ideas.

In 2021, Meyer et al. [26] distilled the ideas from [12, 23] to develop the HUTCH++ algorithm. They proved that HUTCH++ satisfies a worst-case variance bound of $\mathcal{O}(1/m^2)$. Meyer et al. also proposed a version of HUTCH++ that needs only a single pass over the input matrix. The follow-up paper [20] sharpens the analysis of the single-pass algorithm.

Persson, Cortinovis, and Kressner [28] have introduced several refinements to the HUTCH++ estimator. Their first improvement adaptively apportions test vectors between approximating the matrix and estimating the trace of the residual in order to meet an error tolerance. Their second contribution is NYSTRÖM++, a version of HUTCH++ for psd matrices that uses Nyström approximation.

XTRACE and XNYSTRACE build on the previous strategies of variance reduction using low-rank approximation. However, XTRACE and XNYSTRACE take a step forward by also enforcing the exchangeability principle. These algorithms push the ideas of HUTCH++ and NYSTRÖM++ to their limit by using *all* the test vectors for low-rank approximation and *all* the test vectors for residual trace estimation.

To conclude, let us mention two techniques designed for computing the trace of a standard matrix function (that is, $\text{tr } f(\mathbf{A})$). First, stochastic Lanczos quadrature [34] approximates the spectral density of \mathbf{A} , from which estimates of $\text{tr } f(\mathbf{A})$ for any function $f(\cdot)$

are immediately accessible. See [8] for a recent overview of stochastic Lanczos quadrature and related ideas. As a second approach, when the matvecs $\boldsymbol{\omega} \mapsto f(\mathbf{A})\boldsymbol{\omega}$ are computed using a Krylov subspace method [17, §13.2], the paper [7] recommends reuse of the matvecs from the Krylov subspace method for the purpose of trace estimation.

1.7. Reproducible research. Optimized MATLAB R2022b implementations of our algorithms as well as code to reproduce the experiments in this paper can be found online at <https://github.com/epperly/XTrace>.

1.8. Outline. The balance of the paper is organized as follows. Section 2 describes efficient implementations for XTRACE, XNYSTRACE, and the diagonal estimator XDIAG. Section 3 discusses error estimation and adaptive stopping, section 4 presents numerical experiments, and section 5 proves our theoretical results.

1.9. Notation. Matrices and vectors are denoted by capital and lowercase bold letters. The i th column of \mathbf{B} is expressed as \mathbf{b}_i , and the (i, j) th entry of \mathbf{B} is b_{ij} . For a matrix \mathbf{B} , we form a matrix \mathbf{B}_{-i} by deleting the i th column from \mathbf{B} . Similarly, we form \mathbf{B}_{-ij} by deleting the i th and j th columns. We work with the spectral norm $\|\cdot\|$, the Frobenius norm $\|\cdot\|_F$, and the trace norm $\|\cdot\|_*$. The symbol $\llbracket \mathbf{B} \rrbracket_r$ denotes a (simultaneous) best rank- r approximation of \mathbf{B} with respect to any unitarily invariant norm.

2. Efficient implementation of exchangeable trace estimators. This section works through some issues that arise in the implementation of the XTRACE and XNYSTRACE estimators. Subsections 2.1 and 2.2 show how to implement the new estimators efficiently using insights from numerical linear algebra. Subsection 2.3 discusses a method of renormalizing the test vectors that improves the accuracy of XTRACE and XNYSTRACE. Subsection 2.4 develops the XDIAG estimator.

2.1. Computing XTRACE. In this section, we develop an efficient implementation of the XTRACE estimator from subsection 1.2.1. Recall that $\mathbf{A} \in \mathbb{R}^{N \times N}$ is a general square matrix, and introduce the test matrix $\boldsymbol{\Omega} = [\boldsymbol{\omega}_1 \ \dots \ \boldsymbol{\omega}_{m/2}] \in \mathbb{R}^{N \times (m/2)}$.

First, we form the matrix product $\mathbf{A}\boldsymbol{\Omega}$ and compute the orthogonal decomposition $\mathbf{A}\boldsymbol{\Omega} = \mathbf{Q}\mathbf{R}$. Following [10, App. A], we make the critical observation that the basis matrix $\mathbf{Q}_{(i)} = \text{orth}(\mathbf{A}\boldsymbol{\Omega}_{-i})$ is related to the full basis matrix \mathbf{Q} by a rank-one update:

$$(2.1) \quad \mathbf{Q}_{(i)}\mathbf{Q}_{(i)}^* = \mathbf{Q}(\mathbf{I} - \mathbf{s}_i\mathbf{s}_i^*)\mathbf{Q}^* \quad \text{where} \quad \mathbf{R}_{-i}^*\mathbf{s}_i = \mathbf{0} \quad \text{and} \quad \|\mathbf{s}_i\|_{\ell_2} = 1.$$

Thus, the rank-one update requires a unit vector $\mathbf{s}_i \in \mathbb{R}^N$ in the null space of \mathbf{R}_{-i}^* .

Let us exhibit an efficient algorithm that simultaneously computes all the vectors \mathbf{s}_i for $1 \leq i \leq m/2$. We argue that the matrix $\mathbf{S} = [\mathbf{s}_1 \ \dots \ \mathbf{s}_{m/2}]$ can be represented as

$$(2.2) \quad \mathbf{S} = (\mathbf{R}^*)^{-1}\mathbf{D},$$

where the diagonal matrix \mathbf{D} enforces the normalization of the columns of \mathbf{S} . Indeed, since $\mathbf{R}^*\mathbf{S} = \mathbf{D}$ is diagonal, the i th column of $\mathbf{R}_{-i}^*\mathbf{S}$ is the zero vector. We reach the desired conclusion $\mathbf{R}_{-i}^*\mathbf{s}_i = \mathbf{0}$.

In summary, given the full basis \mathbf{Q} , we can use (2.2) to compute all the vectors \mathbf{s}_i needed to construct the orthogonal projectors $\mathbf{Q}_{(i)}\mathbf{Q}_{(i)}^*$ for $1 \leq i \leq m/2$ appearing in (2.1). This calculation requires just $\mathcal{O}(m^3)$ operations, which is dominated by the cost of solving $m/2$ triangular linear systems. It follows that the XTRACE estimator can be computed in just $\mathcal{O}(m^2N)$ operations, which is the same asymptotic cost as HUTCH++. For full details, see the efficient MATLAB implementation in the supplementary materials Program SM2.1.

DRAFT

2.2. Computing XNYS TRACE. We can design a similar method to compute the XNYS TRACE estimator (1.9) efficiently. Let $\mathbf{A} \in \mathbb{R}^{N \times N}$ be a psd matrix and define the test matrix $\mathbf{\Omega} = [\boldsymbol{\omega}_1 \ \dots \ \boldsymbol{\omega}_m]$.

As before, we compute the orthogonal decomposition $\mathbf{A}\mathbf{\Omega} = \mathbf{Q}\mathbf{R}$. We can express the Nyström approximation (1.7) in the form

$$\mathbf{A}\langle \mathbf{\Omega} \rangle = \mathbf{Q}\mathbf{R}\mathbf{H}^{-1}\mathbf{R}^*\mathbf{Q}^*,$$

where $\mathbf{H} = \mathbf{\Omega}^*\mathbf{A}\mathbf{\Omega}$. After deleting the i th column from $\mathbf{\Omega}$, the resulting Nyström approximation satisfies

$$(2.3) \quad \mathbf{A}\langle \mathbf{\Omega}_{-i} \rangle = \mathbf{Q}\mathbf{R}_{-i}\mathbf{H}_{(i)}^{-1}\mathbf{R}_{-i}^*\mathbf{Q}^*,$$

where $\mathbf{H}_{(i)}$ is \mathbf{H} upon deletion of its i th row and column. To compute $\mathbf{A}\langle \mathbf{\Omega}_{-i} \rangle$ efficiently, we recognize that (2.3) can be expressed as a rank-one update:

$$(2.4) \quad \mathbf{A}\langle \mathbf{\Omega}_{-i} \rangle = \mathbf{Q}\mathbf{R} \left(\mathbf{H}^{-1} - \frac{\mathbf{H}^{-1}\mathbf{e}_i\mathbf{e}_i^*\mathbf{H}^{-1}}{\mathbf{e}_i^*\mathbf{H}^{-1}\mathbf{e}_i} \right) \mathbf{R}^*\mathbf{Q}^*.$$

Taking advantage of the rank-one update formula (2.4), we can form the XNYS TRACE estimator using just m matvecs and $\mathcal{O}(m^2N)$ post-processing operations. An efficient MATLAB implementation appears in the supplementary materials [Program SM2.2](#), which incorporates several additional tricks taken from [22, 33] to improve its numerical stability.

2.3. Normalization of test vectors. For the best performance of XTRACE and XNYS TRACE, we recommend the following normalization of the test vectors. First, we draw the test vectors from a spherically symmetric distribution, such as $\boldsymbol{\omega} \sim \text{NORMAL}(\mathbf{0}, \mathbf{I})$. We use these test vectors to form the matrix $\mathbf{Q}_{(i)}$ or Nyström approximation $\mathbf{A}\langle \mathbf{\Omega}_{-i} \rangle$. Second, when computing the basic trace estimates, we **normalize** the test vectors. In XTRACE, we compute

$$\boldsymbol{\mu}_i := (\mathbf{I} - \mathbf{Q}_{(i)}\mathbf{Q}_{(i)}^*)\boldsymbol{\omega}_i \quad \text{and} \quad \mathbf{v}_i := \sqrt{N - \text{rank}(\mathbf{Q}_{(i)})} \cdot \boldsymbol{\mu}_i / \|\boldsymbol{\mu}_i\|_{\ell_2}.$$

To obtain the i th trace estimate, we form

$$\widehat{\text{tr}}_i := \text{tr}(\mathbf{Q}_{(i)}^*(\mathbf{A}\mathbf{Q}_{(i)})) + \mathbf{v}_i^*(\mathbf{A}\mathbf{v}_i).$$

For XNYS TRACE, we set

$$\mathbf{P}_{(i)} := \text{orth}\mathbf{\Omega}_{-i} \quad \text{and} \quad \boldsymbol{\mu}_i := (\mathbf{I} - \mathbf{P}_{(i)}\mathbf{P}_{(i)}^*)\boldsymbol{\omega}_i \quad \text{and} \quad \mathbf{v}_i := \sqrt{N - \text{rank}(\mathbf{P}_{(i)})} \cdot \boldsymbol{\mu}_i / \|\boldsymbol{\mu}_i\|_{\ell_2}.$$

Then define the basic trace estimates

$$\widehat{\text{tr}}_i := \text{tr}\mathbf{A}\langle \mathbf{\Omega}_{-i} \rangle + \mathbf{v}_i^*(\mathbf{A}\mathbf{v}_i).$$

The normalization removes a source of variance related to the random lengths of the vectors $\boldsymbol{\mu}_i$, improving the accuracy compared to unnormalized Gaussian test vectors or uniform random vectors on the sphere. We compare this normalization approach against alternative distributions for test vectors in [subsection 4.2](#).

2.4. Diagonal estimation. In the spirit of Girard and Hutchinson, the paper [5] of Bekas, Kokiopoulou, and Saad (BKS) develops an estimator for the diagonal of an implicit matrix:

$$(2.5) \quad \widehat{\text{diag}}_{\text{BKS}} = \frac{\sum_{i=1}^m \boldsymbol{\omega}_i \odot (\mathbf{A}\boldsymbol{\omega}_i)}{\sum_{i=1}^m \boldsymbol{\omega}_i \odot \boldsymbol{\omega}_i}.$$

DRAFT

Here, \odot denotes the entrywise product and the division is performed entrywise. The recent paper [4] proposes a biased estimator for the diagonal, called DIAG++, that is inspired by BKS and HUTCH++.

We have observed that the exchangeability principle leads to an unbiased diagonal estimator with lower variance. Our diagonal estimator, XDIAG, takes the form

$$\widehat{\text{diag}}_X = \frac{1}{m/2} \sum_{i=1}^{m/2} \left[\text{diag}(\mathbf{Q}_{(i)}(\mathbf{Q}_{(i)}^* \mathbf{A})) + \frac{\boldsymbol{\omega}_i \odot (\mathbf{I} - \mathbf{Q}_{(i)} \mathbf{Q}_{(i)}^*)(\mathbf{A} \boldsymbol{\omega}_i)}{\boldsymbol{\omega}_i \odot \boldsymbol{\omega}_i} \right],$$

where $\mathbf{Q}_{(i)}$ is defined in (1.4). In contrast to XTRACE, the XDIAG estimator requires matvecs with \mathbf{A}^* in addition to matvecs with \mathbf{A} . The same ideas from subsection 2.1 allow us to implement XDIAG in $\mathcal{O}(m^2 N)$ operations; an implementation is provided in the supplementary materials (Program SM2.3).

3. Error estimation and adaptive stopping. Our exchangeable estimators depend on averaging over a family of basic estimators, and we can reuse the basic estimators to compute a reliable posterior approximation for the error (subsection 3.1). The error estimate allows us to develop adaptive methods for selecting the number m of matvecs to achieve a specified error tolerance (subsection 3.2). These refinements are very important for practical implementations.

3.1. Error estimation. The XTRACE and XNYS TRACE estimators are both formed as averages of individual trace estimates $\widehat{\text{tr}}_1, \dots, \widehat{\text{tr}}_\ell$ where

$$\begin{aligned} \widehat{\text{tr}}_i &= \text{tr}(\mathbf{Q}_{(i)}^* \mathbf{A} \mathbf{Q}_{(i)}) + \boldsymbol{\omega}_i^* (\mathbf{I} - \mathbf{Q}_{(i)} \mathbf{Q}_{(i)}^*) \mathbf{A} (\mathbf{I} - \mathbf{Q}_{(i)} \mathbf{Q}_{(i)}^*) \boldsymbol{\omega}_i, & \ell &:= \frac{m}{2} & (\text{XTRACE}); \\ \widehat{\text{tr}}_i &= \text{tr}(\mathbf{A} \langle \boldsymbol{\Omega}_{-i} \rangle) + \boldsymbol{\omega}_i^* (\mathbf{A} - \mathbf{A} \langle \boldsymbol{\Omega}_{-i} \rangle) \boldsymbol{\omega}_i, & \ell &:= m & (\text{XNYS TRACE}). \end{aligned}$$

The scaled variance of the individual trace estimates $\widehat{\text{tr}}_i$ provides a useful estimate for the squared error in the trace estimate:

$$(3.1) \quad \widehat{\text{err}}^2 := \frac{1}{\ell(\ell-1)} \sum_{i=1}^{\ell} |\widehat{\text{tr}}_i - \widehat{\text{tr}}|^2 \approx |\text{tr}(\mathbf{A}) - \widehat{\text{tr}}|^2 \quad \text{where} \quad \widehat{\text{tr}} = \frac{1}{\ell} \sum_{i=1}^{\ell} \widehat{\text{tr}}_i.$$

The next result contains an analysis of this posterior error estimator.

PROPOSITION 3.1 (Error estimate). *The error estimate (3.1) satisfies*

$$\mathbb{E} \widehat{\text{err}}^2 = \frac{1 - \text{Cor}(\widehat{\text{tr}}_1, \widehat{\text{tr}}_2)}{1 + (\ell - 1) \text{Cor}(\widehat{\text{tr}}_1, \widehat{\text{tr}}_2)} \cdot \mathbb{E} |\text{tr}(\mathbf{A}) - \widehat{\text{tr}}|^2.$$

We have written $\text{Cor}(\cdot, \cdot)$ for the correlation of two random variables.

The proof and some additional discussion of the correlation $\text{Cor}(\widehat{\text{tr}}_1, \widehat{\text{tr}}_2)$ appears in subsection 5.6. In practice, we find that the individual trace estimators have a small positive correlation, so we typically *underestimate* the true error by a small amount. Thus, the posterior error estimate is a valuable tool. For an illustration, see Figure 5a in subsection 4.3.

3.2. Adaptive stopping. In practice, we often wish to choose the number m of matvecs adaptively to estimate $\text{tr} \mathbf{A}$ up to a prescribed accuracy level:

$$|\text{tr}(\mathbf{A}) - \widehat{\text{tr}}| \leq \varepsilon \cdot |\text{tr} \mathbf{A}| \quad \text{for } \varepsilon \in (0, 1).$$

One simple way to achieve this tolerance is through a *doubling strategy*:

DRAFT

1. To initialize, collect m_0 matvecs $\mathbf{A}\boldsymbol{\omega}_1, \dots, \mathbf{A}\boldsymbol{\omega}_{m_0}$, and set $j \leftarrow 0$.
2. Use $\mathbf{A}\boldsymbol{\omega}_1, \dots, \mathbf{A}\boldsymbol{\omega}_{m_j}$ to form a trace estimate $\widehat{\text{tr}}^{(j)}$ and an error estimate $\widehat{\text{err}}^{(j)}$.
3. If $\widehat{\text{err}}^{(j)} \leq \varepsilon \cdot |\widehat{\text{tr}}^{(j)}|$, then stop.
4. Collect m_j additional matvecs $\mathbf{A}\boldsymbol{\omega}_{m_j+1}, \dots, \mathbf{A}\boldsymbol{\omega}_{2m_j}$. Set $j \leftarrow j+1$ and $m_j \leftarrow 2m_{j-1}$. Go to step 2.

The doubling strategy requires at most twice the optimal number of matvecs to meet the tolerance and maintains the $\mathcal{O}(m^2 N)$ computational cost of XTRACE and XNYSTRACE. We implement this approach in our experiments to produce [Figure 5b](#).

4. Numerical experiments. This section presents a numerical evaluation of XTRACE, XNYSTRACE, and XDIAG. [Subsection 4.1](#) compares different trace estimators on synthetic matrices, [subsection 4.2](#) evaluates XNYSTRACE with different distributions for the test vectors, [subsection 4.3](#) applies XTRACE and XNYSTRACE to computations in quantum statistical physics, and [subsection 4.4](#) applies XDIAG to computations in network science. Code to reproduce the experiments can be found at <https://github.com/epperly/XTrace>.

4.1. Comparison of trace estimators. The first experiment is designed to compare the accuracy of six trace estimators that each use m matvecs:

- HUTCH: The Girard–Hutchinson estimator [\(1.2\)](#).
- LRA: The Saibaba et al. [\[31\]](#) estimator (without additional subspace iteration): $\text{tr} \widehat{\mathbf{A}}$ where $\widehat{\mathbf{A}} = \mathbf{Q}\mathbf{Q}^* \mathbf{A}$ and $\mathbf{Q} = \text{orth}(\mathbf{A}\boldsymbol{\Omega})$.
- The HUTCH++ estimator [\(1.3\)](#).
- The NYSTRÖM++ estimator [\[28\]](#). We use the implementation provided by the authors of [\[28\]](#), modified to the test vector $\boldsymbol{\omega} \sim \text{UNIFORM}\{\pm 1\}^N$.
- The XTRACE estimator [\(1.6\)](#).
- The XNYSTRACE estimator [\(1.9\)](#).

To create a fair comparison, we apply all six estimators using a test matrix whose entries are uniformly random signs: $\boldsymbol{\Omega} \sim \text{UNIFORM}\{\pm 1\}^{N \times \ell}$, as was used in HUTCH++. The additional benefit for XTRACE and XNYSTRACE of using normalized, spherically symmetric test vectors is explored in [subsection 4.2](#). The supplementary materials ([section SM1](#)) contain additional comparisons with the adaptive HUTCH++ algorithm of [\[28\]](#); this comparison requires a more complicated experimental setup.

We apply each of these estimators to randomly generated matrices of the form

$$\mathbf{A}(\boldsymbol{\lambda}) = \mathbf{U} \text{diag}(\boldsymbol{\lambda}) \mathbf{U}^*$$

where \mathbf{U} is a Haar random orthogonal matrix. We use four choices for the eigenvalues $\boldsymbol{\lambda}$:

- flat: $\boldsymbol{\lambda} = (3 - 2(i - 1)/(N - 1) : i = 1, 2, \dots, N)$.
- poly: $\boldsymbol{\lambda} = (i^{-2} : i = 1, 2, \dots, N)$.
- exp: $\boldsymbol{\lambda} = (0.7^i : i = 0, 1, \dots, N - 1)$.
- step: $\boldsymbol{\lambda} = (\underbrace{1, \dots, 1}_{50 \text{ times}}, \underbrace{10^{-3}, \dots, 10^{-3}}_{N - 50 \text{ times}})$.

We fix the matrix dimension $N = 1000$, and we report the relative error $|\text{tr}(\mathbf{A}) - \widehat{\text{tr}}| / \text{tr}(\mathbf{A})$ averaged over 1000 trials.

Discussion. The variance-reduced trace estimators dramatically outperform HUTCH, except on the flat instance ([Figure 3a](#)). The implication is that HUTCH only makes sense when estimating the trace of a matrix with a nearly flat spectrum. For the flat instance, the performance of LRA is especially poor because LRA is a biased estimator that substantially underestimates the trace.

Across all the instances, XTRACE produces smaller errors than HUTCH++, sometimes by orders of magnitude. For the exp instance ([Figure 3c](#)), the error of XTRACE decays expo-

DRAFT

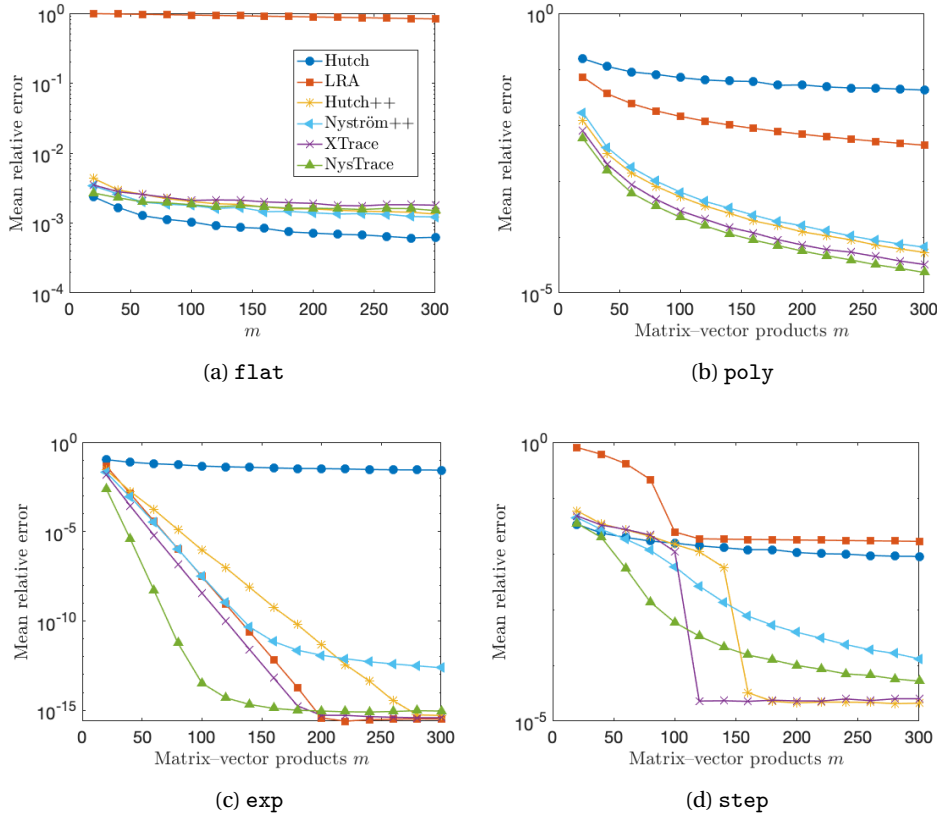


Fig. 3: **Synthetic instances.** Average relative error of trace estimators for matrices with particular spectral profiles using random sign test vectors. See [subsection 4.1](#) and [Figure 4](#).

nentially fast at a rate $1.5\times$ faster than HUTCH++. The superiority of XTRACE is also visible for the `step` instance ([Figure 3d](#)) where XTRACE achieves accuracy 10^{-4} with just $m \approx 120$ matvecs, as compared to $m \approx 160$ for HUTCH++.

XNYSTRACE is frequently the most accurate of the trace estimation methods. For the `exp` instance ([Figure 3c](#)), XNYSTRACE converges at a rate $2\times$ faster than XTRACE and NYSTRÖM++ and $3\times$ faster than HUTCH++. However, XNYSTRACE (and NYSTRÖM++) can exhibit poor performance for matrices that possess a long tail of slowly decreasing eigenvalues (see the `step` instance in [Figure 3d](#)). To understand this phenomenon, observe that the error bounds for XNYSTRACE depend on the trace norm, which is sensitive to slow eigenvalue decay ([Theorem 1.1](#)). We can improve the performance by using the normalization approach of [subsection 2.3](#), as we detail in the next section.

4.2. Choice of test vectors. In [subsection 2.3](#), we recommended an implementation of XTRACE and XNYSTRACE that uses rotationally invariant test vectors for low-rank approximation and normalizes the distinguished test vector used for trace estimation.

This section shows how this method can improve over estimators that lack the normalization step. We compare against test vectors from the random sign distribution $\omega \sim \text{UNIFORM}\{\pm 1\}^N$, the standard normal distribution $\omega \sim \text{NORMAL}(\mathbf{0}, \mathbf{I})$, and the uniform distribution on the sphere $\omega \sim \text{UNIFORM}\{\mathbf{x} \in \mathbb{R}^N : \|\mathbf{x}\|_{\ell_2} = \sqrt{N}\}$. The differences among the

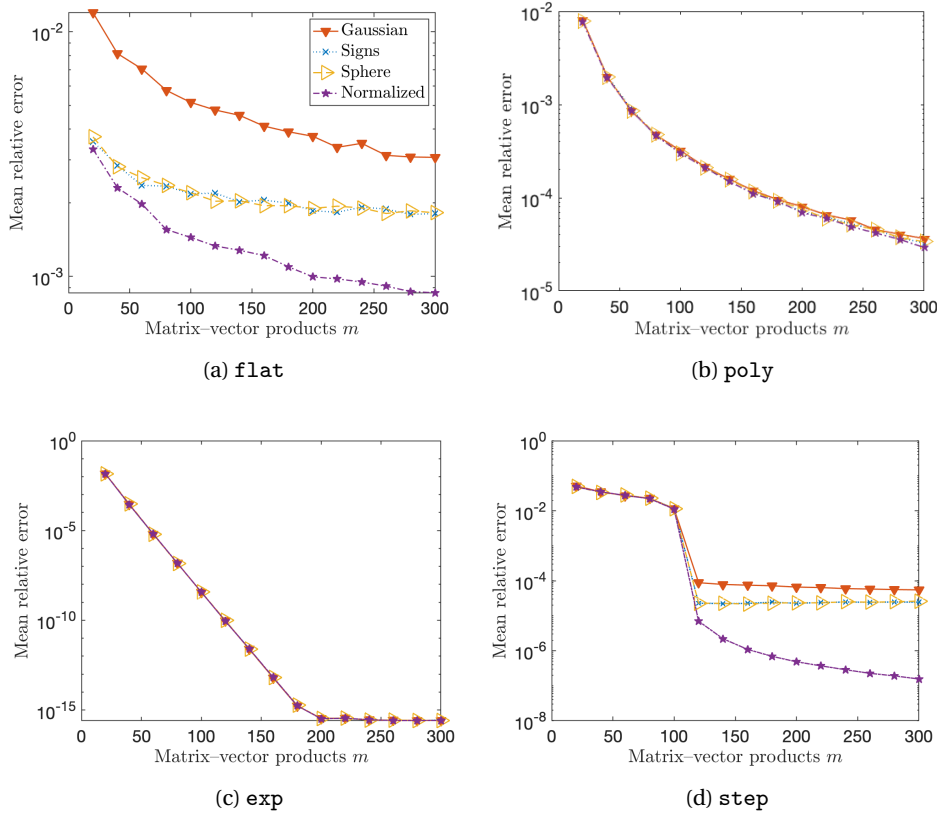


Fig. 4: **Normalization of test vectors.** The average relative error of the XTRACE estimator using normalized test vectors (subsection 2.3) as compared with alternative (unnormalized) test vector distributions. See subsection 4.2.

distributions are only visible for matrices whose spectrum has flat segments, as in the `flat` and `step` examples. For these examples, the normalization strategy is conspicuously the best, followed by the uniform sign and uniform sphere distributions, with the standard normal distribution lagging well behind.

4.3. Application: Quantum statistical mechanics. Our next experiment shows the benefits of using XTRACE and XNYSTRACE for an application in quantum physics. To compute a phase diagram, we must evaluate a large number of trace estimators. Our exchangeable estimators reduce the number of matvecs needed to achieve a desired tolerance, and we can use the posterior error estimator to adaptively determine the minimum number m of matvecs.

The average energy of a quantum system with a symmetric Hamiltonian matrix $\mathbf{H} \in \mathbb{R}^{N \times N}$ at inverse temperature $\beta > 0$ is

$$E(\beta) = \frac{1}{Z(\beta)} \text{tr}[\mathbf{H} \exp(-\beta \mathbf{H})] \quad \text{where } Z(\beta) = \text{tr} \exp(-\beta \mathbf{H}).$$

The quantity $Z(\beta)$ is the partition function, introduced in subsection 1.3. We observe that $\text{tr}[\mathbf{H} \exp(-\beta \mathbf{H})]$ and $\text{tr} \exp(-\beta \mathbf{H})$ are ideal candidates for estimation using XTRACE and

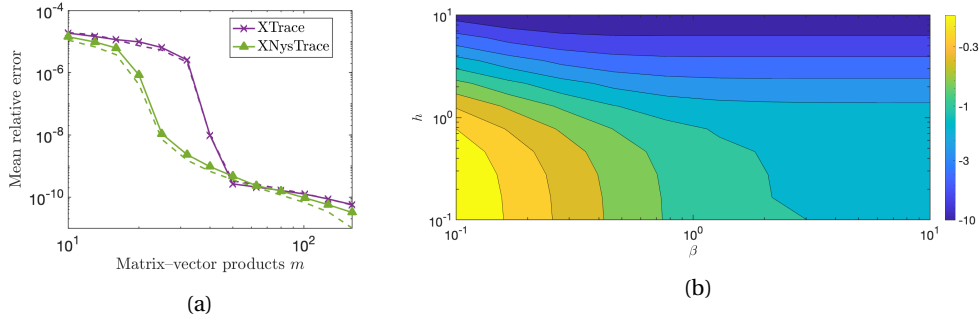


Fig. 5: **Quantum statistical mechanics.** *Left:* Mean relative error (solid lines) and mean posterior error estimates (dashed lines) when computing the partition function $Z(\beta = 0.6, h = 10)$. *Right:* Average energy $E(\beta, h)/n$ per site for $\beta, h \in [10^{-1}, 10^1]$ as computed by XNYS TRACE. See [subsection 4.3](#).

XNYS TRACE, since the matrix exponential leads to rapidly decaying eigenvalues.

We apply XTRACE and XNYS TRACE to compute the partition function and energy for the transverse field Ising model (TFIM) for a periodic 1D chain [29], which is specified by the Hamiltonian matrix

$$(4.1) \quad \mathbf{H} = - \sum_{i=1}^n \mathbf{Z}_i \mathbf{Z}_{i+1} - h \sum_{i=1}^n \mathbf{X}_i \in \mathbb{R}^{2^n \times 2^n}.$$

Here, \mathbf{X}_i and \mathbf{Z}_i denote Pauli operators acting on the i th site; that is,

$$\mathbf{X}_i = \mathbf{I}_{2 \times 2}^{\otimes (n-i-1)} \otimes \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \otimes \mathbf{I}_{2 \times 2}^{\otimes (n-i)}, \quad \mathbf{Z}_i = \mathbf{I}_{2 \times 2}^{\otimes (n-i-1)} \otimes \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} \otimes \mathbf{I}_{2 \times 2}^{\otimes (n-i)},$$

and $\mathbf{Z}_{n+1} = \mathbf{Z}_1$ by periodicity. The eigenvalues of \mathbf{H} are known exactly [24, eqs. (16)–(17)], which allows us to precisely evaluate the error of stochastic estimates of $Z(\beta) = Z(\beta, h)$ and $E(\beta) = E(\beta, h)$. Before applying stochastic trace estimation, we shift the Hamiltonian matrix by a constant $b = (1 + h)n$ so that $\mathbf{H} + b\mathbf{I}$ is positive semidefinite.

Figure 5a shows the errors of XTRACE and XNYS TRACE when computing the partition function $Z(\beta, h)$ of the TFIM with $n = 18$, $h = 10$, and $\beta = 0.6$; this is the same setting as in Figure 2. The thick lines indicate the average errors over 10 trials, while the dashed lines indicate the average error estimates introduced in subsection 3.1. We observe that the error estimates closely track the true errors, differing by a factor of at most 3.2.

Figure 5b shows the average energy $E(\beta, h)/n$ per site for parameters $\beta, h \in [10^{-1}, 10^1]$, up to a relative error of 10^{-3} . We compute the energy by using XNYS TRACE, together with the doubling strategy from subsection 3.2. To ensure the robustness of the doubling strategy, we use a slightly stricter tolerance $\varepsilon = 10^{-4}$ than our desired accuracy of 10^{-3} .

4.4. Application: Networks. One of the basic problems in network science is to measure the centrality of each node in a graph. We focus on two centrality measures, which can be defined in terms of the adjacency matrix \mathbf{M} :

1. The *number of triangles* [1] incident on node i is given by $\Delta_i(\mathbf{M}) = \frac{1}{2}(\mathbf{M}^3)_{ii}$.
2. The *subgraph centrality* [11] of node i is defined as $\text{SC}_i := (\exp(\mathbf{M}))_{ii}$.

Both centrality measures are the diagonal entries of functions of the adjacency matrix.

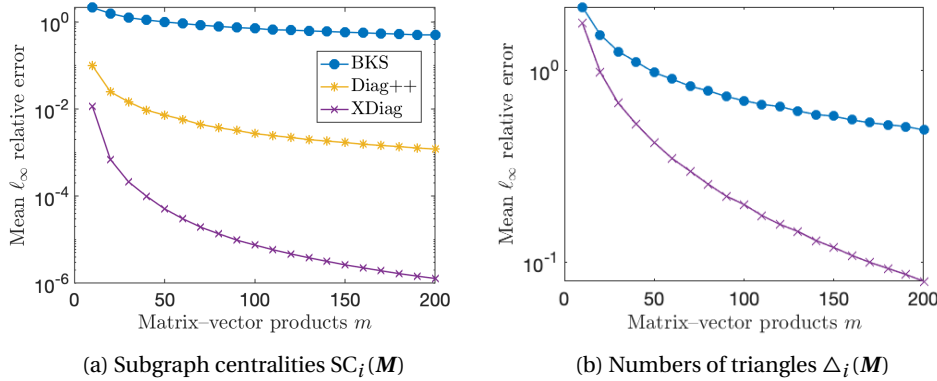


Fig. 6: **Networks.** Mean error of BKS diagonal estimator and XDIAG for the subgraph centralities (*left*) and triangle numbers (*right*) for the yeast graph. See [subsection 4.4](#).

Using the BKS, DIAG++, and XDIAG diagonal estimators, we estimate these centrality measures for the protein–protein interaction network for budding yeast [6], available in the SuiteSparse collection [9]. (Following [4], DIAG++ is omitted for the triangle problem because $\mathbf{M}^3/2$ is not psd.) We evaluate the quality of our estimates using the relative ℓ_∞ error

$$\text{error}(\widehat{\text{diag}}) := \frac{\max_{1 \leq i \leq N} |a_{ii} - \widehat{\text{diag}}_i|}{\max_{1 \leq i \leq N} |a_{ii}|},$$

averaged over 1000 trials. [Figure 6](#) shows the results. For the subgraph centrality problem after $m = 200$ matvecs, XDIAG is more accurate than BKS by five orders of magnitude and more accurate than DIAG++ by three orders of magnitude.

5. Theoretical analysis. In this final section, we prove [Theorem 1.1](#), which provides refined error bounds for three trace estimators, and we prove [Proposition 3.1](#), which describes the behavior of the posterior error estimator.

5.1. XTRACE variance bound. To begin, let us establish an initial variance bound for the XTRACE estimator. This result shows that the variance depends on the error in a low-rank approximation of the input matrix. Later, we will bound these errors using standard results for the randomized SVD.

PROPOSITION 5.1 (XTRACE error). *Fix $\mathbf{A} \in \mathbb{R}^{N \times N}$, and consider the XTRACE estimator $\widehat{\text{tr}}_X$ defined in (1.6) with a test matrix $\mathbf{\Omega} = [\boldsymbol{\omega}_1 \ \dots \ \boldsymbol{\omega}_{m/2}]$ consisting of $m/2$ standard normal test vectors. The estimator is unbiased: $\mathbb{E} \widehat{\text{tr}}_X = \text{tr} \mathbf{A}$. Moreover, the variance satisfies*

$$\text{Var}[\widehat{\text{tr}}_X] = \mathbb{E} |\widehat{\text{tr}}_X - \text{tr} \mathbf{A}|^2 \leq \frac{2}{m/2} \mathbb{E} \|(\mathbf{I} - \mathbf{Q}_{(1)} \mathbf{Q}_{(1)}^*) \mathbf{A}\|_F^2 + 4 \mathbb{E} \|(\mathbf{I} - \mathbf{Q}_{(12)} \mathbf{Q}_{(12)}^*) \mathbf{A}\|^2,$$

where $\mathbf{Q}_{(i)} = \text{orth}(\mathbf{A}\boldsymbol{\Omega}_{-i})$ and $\mathbf{Q}_{(ij)} = \text{orth}(\mathbf{A}\boldsymbol{\Omega}_{-ij})$.

Proof. For all indices $1 \leq i, j \leq m/2$, we abbreviate the orthogonal projectors $\mathbf{\Pi}_i := \mathbf{Q}_{(i)} \mathbf{Q}_{(i)}^*$ and $\mathbf{\Pi}_{ij} := \mathbf{Q}_{(ij)} \mathbf{Q}_{(ij)}^*$. Note that $\boldsymbol{\omega}_i$ is independent from $\mathbf{\Pi}_i$, while $(\boldsymbol{\omega}_i, \boldsymbol{\omega}_j)$ is independent from $\mathbf{\Pi}_{ij}$.

For each $1 \leq i \leq m/2$, we can check that the basic estimator $\hat{\tau}_i$ defined in (1.5) is unbiased. To do so, condition on $\mathbf{\Omega}_{-i}$ and average over $\boldsymbol{\omega}_i$ to arrive at the identity

$$(5.1) \quad \begin{aligned} \mathbb{E}[\hat{\tau}_i \mid \mathbf{\Omega}_{-i}] &= \mathbb{E}[\text{tr}(\mathbf{\Pi}_i \mathbf{A} \mathbf{\Pi}_i) + \boldsymbol{\omega}_i^* (\mathbf{I} - \mathbf{\Pi}_i) \mathbf{A} (\mathbf{I} - \mathbf{\Pi}_i) \boldsymbol{\omega}_i \mid \mathbf{\Omega}_{-i}] \\ &= \mathbb{E}[\text{tr}(\mathbf{\Pi}_i \mathbf{A} \mathbf{\Pi}_i) + \text{tr}((\mathbf{I} - \mathbf{\Pi}_i) \mathbf{A} (\mathbf{I} - \mathbf{\Pi}_i)) \mid \mathbf{\Omega}_{-i}] = \text{tr} \mathbf{A}. \end{aligned}$$

The second equality uses the fact that each test vector $\boldsymbol{\omega}_i$ is isotropic and independent from $\mathbf{\Pi}_i$, which is a function of $\mathbf{\Omega}_{-i}$. The third equality follows when we cycle the traces and invoke the fact that the projector $\mathbf{I} - \mathbf{\Pi}_i$ is idempotent. We confirm that $\hat{\tau}_i$ is unbiased by applying the tower law to take the total expectation of (5.1). The full estimator $\hat{\tau}_X$ is unbiased because it is an average of the unbiased estimators $\hat{\tau}_i$.

Next, to bound the variance, we use the exchangeability of $\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_{m/2}$ to compute

$$\text{Var}[\hat{\tau}_X] = \text{Var}\left[\frac{1}{m/2} \sum_{i=1}^{m/2} \hat{\tau}_i\right] = \underbrace{\left(\frac{1}{m/2}\right) \text{Var}[\hat{\tau}_1]}_A + \underbrace{\left(1 - \frac{1}{m/2}\right) \text{Cov}[\hat{\tau}_1, \hat{\tau}_2]}_B.$$

Hence, $\text{Var}[\hat{\tau}_X]$ is the weighted average of a variance term A and a covariance term B.

To evaluate the variance term A, condition on $\mathbf{\Omega}_{-1}$ and average over $\boldsymbol{\omega}_1$. Thus,

$$\begin{aligned} A &= \text{Var}[\hat{\tau}_1] = \mathbb{E}[\text{Var}[\hat{\tau}_1 \mid \mathbf{\Omega}_{-1}]] + \text{Var}[\mathbb{E}[\hat{\tau}_1 \mid \mathbf{\Omega}_{-1}]] \\ &= \mathbb{E}[\text{Var}[\text{tr}(\mathbf{\Pi}_1^* \mathbf{A} \mathbf{\Pi}_1) + \boldsymbol{\omega}_1^* (\mathbf{I} - \mathbf{\Pi}_1) \mathbf{A} (\mathbf{I} - \mathbf{\Pi}_1) \boldsymbol{\omega}_1 \mid \mathbf{\Omega}_{-1}]] \\ &= \mathbb{E}[\text{Var}[\boldsymbol{\omega}_1^* (\mathbf{I} - \mathbf{\Pi}_1) \mathbf{A} (\mathbf{I} - \mathbf{\Pi}_1) \boldsymbol{\omega}_1 \mid \mathbf{\Omega}_{-1}]] \leq 2 \mathbb{E}\|(\mathbf{I} - \mathbf{\Pi}_1) \mathbf{A} (\mathbf{I} - \mathbf{\Pi}_1)\|_{\mathbb{F}}^2. \end{aligned}$$

The first relation is the chain rule for the variance. To pass to the second line, we invoke the fact (5.1) that the conditional expectation is constant. To pass to the third line, we drop the trace, which is conditionally constant. The last relation follows from a direct calculation using the facts that $\boldsymbol{\omega}_1$ is standard normal and independent from $\mathbf{\Pi}_1$.

To bound the covariance term B, it is helpful to isolate the part of the covariance that only depends on $\boldsymbol{\omega}_1$ and $\boldsymbol{\omega}_2$. We rely on the following observation. For any (random) matrix $\mathbf{X} \in \mathbb{R}^{N \times N}$ that is independent from $\boldsymbol{\omega}_1$ and $\boldsymbol{\omega}_2$, we may calculate that

$$(5.2) \quad \mathbb{E}[(\hat{\tau}_1 - \text{tr} \mathbf{A} + \text{tr} \mathbf{X} - \boldsymbol{\omega}_1^* \mathbf{X} \boldsymbol{\omega}_1)(\hat{\tau}_2 - \text{tr} \mathbf{A} + \text{tr} \mathbf{X} - \boldsymbol{\omega}_2^* \mathbf{X} \boldsymbol{\omega}_2)]$$

$$(5.3) \quad = \mathbb{E}[(\hat{\tau}_1 - \text{tr} \mathbf{A} + \text{tr} \mathbf{X} - \boldsymbol{\omega}_1^* \mathbf{X} \boldsymbol{\omega}_1)(\hat{\tau}_2 - \text{tr} \mathbf{A})]$$

$$(5.4) \quad = \mathbb{E}[(\hat{\tau}_1 - \text{tr} \mathbf{A})(\hat{\tau}_2 - \text{tr} \mathbf{A})] = \text{Cov}[\hat{\tau}_1, \hat{\tau}_2].$$

To pass to (5.3), we condition on $\mathbf{\Omega}_{-1}$ and average over $\boldsymbol{\omega}_1$, exploiting the fact (5.1) that $\hat{\tau}_1$ is an unbiased estimator of $\text{tr} \mathbf{A}$, conditional on $\mathbf{\Omega}_{-1}$. To pass to (5.4), condition on $\mathbf{\Omega}_{-2}$ and average over $\boldsymbol{\omega}_2$.

To continue, select the particular random matrix $\mathbf{X} = (\mathbf{I} - \mathbf{\Pi}_{12}) \mathbf{A} (\mathbf{I} - \mathbf{\Pi}_{12})$. Applying the Cauchy-Schwarz inequality, we find that

$$\begin{aligned} B &= \text{Cov}[\hat{\tau}_1, \hat{\tau}_2] = \mathbb{E}[(\hat{\tau}_1 - \text{tr} \mathbf{A} + \text{tr} \mathbf{X} - \boldsymbol{\omega}_1^* \mathbf{X} \boldsymbol{\omega}_1)(\hat{\tau}_2 - \text{tr} \mathbf{A} + \text{tr} \mathbf{X} - \boldsymbol{\omega}_2^* \mathbf{X} \boldsymbol{\omega}_2)] \\ &\leq \mathbb{E}|\hat{\tau}_1 - \text{tr} \mathbf{A} + \text{tr} \mathbf{X} - \boldsymbol{\omega}_1^* \mathbf{X} \boldsymbol{\omega}_1|^2 \\ &= \mathbb{E} \text{Var}[\boldsymbol{\omega}_1^* [(\mathbf{I} - \mathbf{\Pi}_1) \mathbf{A} (\mathbf{I} - \mathbf{\Pi}_1) - \mathbf{X}] \boldsymbol{\omega}_1 \mid \mathbf{\Omega}_{-1}] \\ &\leq 2 \mathbb{E}\|(\mathbf{I} - \mathbf{\Pi}_1) \mathbf{A} (\mathbf{I} - \mathbf{\Pi}_1) - \mathbf{X}\|_{\mathbb{F}}^2. \end{aligned}$$

Since $\text{range}(\mathbf{Q}_{(12)}) \subseteq \text{range}(\mathbf{Q}_{(1)})$, we have the relations $\mathbf{\Pi}_{12} = \mathbf{\Pi}_1 \mathbf{\Pi}_{12}$ and $\mathbf{I} - \mathbf{\Pi}_1 = (\mathbf{I} - \mathbf{\Pi}_1)(\mathbf{I} -$

DRAFT

$\mathbf{\Pi}_{12}$). Since $\mathbf{X} = (\mathbf{I} - \mathbf{\Pi}_{12})\mathbf{A}(\mathbf{I} - \mathbf{\Pi}_{12})$, it follows that

$$\begin{aligned} \|(\mathbf{I} - \mathbf{\Pi}_1)\mathbf{A}(\mathbf{I} - \mathbf{\Pi}_1) - \mathbf{X}\|_{\text{F}}^2 &= \|(\mathbf{I} - \mathbf{\Pi}_1)\mathbf{X}(\mathbf{I} - \mathbf{\Pi}_1) - \mathbf{X}\|_{\text{F}}^2 = \|(\mathbf{I} - \mathbf{\Pi}_1)\mathbf{X}\mathbf{\Pi}_1 + \mathbf{\Pi}_1\mathbf{X}\|_{\text{F}}^2 \\ &= \|(\mathbf{I} - \mathbf{\Pi}_1)\mathbf{X}\mathbf{\Pi}_1\|_{\text{F}}^2 + \|\mathbf{\Pi}_1\mathbf{X}\|_{\text{F}}^2 \leq \|\mathbf{X}\mathbf{\Pi}_1\|_{\text{F}}^2 + \|\mathbf{\Pi}_1\mathbf{X}\|_{\text{F}}^2 \\ &= \|\mathbf{X}(\mathbf{\Pi}_1 - \mathbf{\Pi}_{12})\|_{\text{F}}^2 + \|(\mathbf{\Pi}_1 - \mathbf{\Pi}_{12})\mathbf{X}\|_{\text{F}}^2 \\ &= \|\mathbf{X}(\mathbf{\Pi}_1 - \mathbf{\Pi}_{12})\|^2 + \|(\mathbf{\Pi}_1 - \mathbf{\Pi}_{12})\mathbf{X}\|^2 \leq 2\|\mathbf{X}\|^2. \end{aligned}$$

We invoke the Pythagorean theorem to pass to the second line. To reach the third line, exploit the representations $\mathbf{X} = \mathbf{X}(\mathbf{I} - \mathbf{\Pi}_{12})$ and $\mathbf{X} = (\mathbf{I} - \mathbf{\Pi}_{12})\mathbf{X}$. To pass to the fourth line, note that $\mathbf{\Pi}_1 - \mathbf{\Pi}_{12}$ is a rank-one orthogonal projector; the Frobenius norm and spectral norm coincide for rank-one matrices. Combining the last two displays, we deduce that

$$\mathbf{B} = \text{Cov}[\hat{\text{tr}}_1, \hat{\text{tr}}_2] \leq 4\mathbb{E}\|(\mathbf{I} - \mathbf{\Pi}_{12})\mathbf{A}(\mathbf{I} - \mathbf{\Pi}_{12})\|^2.$$

Combining the estimates for \mathbf{A} and \mathbf{B} , we achieve the stated bound for the variance. \square

5.2. HUTCH++ variance bound. By a similar argument, we can obtain an initial variance bound for the HUTCH++ estimator. This result is more elementary because it does not require us to account for interactions between the simple estimators.

PROPOSITION 5.2 (HUTCH++ error). *Fix $\mathbf{A} \in \mathbb{R}^{N \times N}$, and consider the HUTCH++ estimator $\hat{\text{tr}}_{\text{H}++}$ defined in (1.3) with $2m/3$ standard normal test vectors. The estimator is unbiased: $\mathbb{E}\hat{\text{tr}}_{\text{H}++} = \text{tr } \mathbf{A}$. Moreover, the variance satisfies*

$$\mathbb{E}|\hat{\text{tr}}_{\text{H}++} - \text{tr } \mathbf{A}|^2 \leq \frac{2}{m/3}\mathbb{E}\|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)\mathbf{A}\|_{\text{F}}^2,$$

where $\mathbf{Q} = \text{orth}(\mathbf{A}\mathbf{\Omega})$ and $\mathbf{\Omega} = [\boldsymbol{\omega}_{m/3+1} \ \dots \ \boldsymbol{\omega}_{2m/3}]$.

Proof. The idea is to condition on the low-rank approximation $\mathbf{Q}\mathbf{Q}^*\mathbf{A}$ and invoke the chain rule for the variance, as in the proof of Proposition 5.1. See the argument in [27, Thm. 10], which was supplied by the second author of this paper. \square

5.3. XNYS TRACE variance bound. Last, we establish an initial variance bound for the XNYS TRACE estimator. This result shows how the variance depends on the error in a randomized Nyström approximation. Later, we will use recent results for the Nyström approximation to obtain a complete variance bound.

PROPOSITION 5.3 (XNYS TRACE error). *Let $\mathbf{A} \in \mathbb{R}^{N \times N}$ be psd. Consider the XNYS TRACE estimator $\hat{\text{tr}}_{\text{NX}}$ with m standard normal test vectors as defined in (1.9). The estimator is unbiased: $\mathbb{E}\hat{\text{tr}}_{\text{NX}} = \text{tr } \mathbf{A}$. Moreover, the variance satisfies the bound*

$$\mathbb{E}|\hat{\text{tr}}_{\text{NX}} - \text{tr } \mathbf{A}|^2 \leq \frac{2}{m}\mathbb{E}\|\mathbf{A} - \mathbf{A}\langle \mathbf{\Omega}_{-1} \rangle\|_{\text{F}}^2 + 2\mathbb{E}\|\mathbf{A} - \mathbf{A}\langle \mathbf{\Omega}_{-12} \rangle\|_{\text{F}}^2.$$

Proof. The proof resembles the proof of Proposition 5.1 but is slightly simpler. The unbiasedness of XNYS TRACE follows from a short computation similar to (5.1).

To control the variance, we calculate that

$$\text{Var}[\hat{\text{tr}}_{\text{NX}}] = \text{Var}\left[\frac{1}{m}\sum_{i=1}^m \hat{\text{tr}}_i\right] = \underbrace{\left(\frac{1}{m}\right)\text{Var}[\hat{\text{tr}}_1]}_{\mathbf{A}} + \underbrace{\left(1 - \frac{1}{m}\right)\text{Cov}[\hat{\text{tr}}_1, \hat{\text{tr}}_2]}_{\mathbf{B}}.$$

The variance term is exactly

$$\mathbf{A} = \text{Var}[\hat{\text{tr}}_1] = 2\mathbb{E}\|\mathbf{A} - \mathbf{A}\langle \mathbf{\Omega}_{-1} \rangle\|_{\text{F}}^2.$$

DRAFT

To bound the covariance term B , we set $\mathbf{X} = \mathbf{A} - \mathbf{A}\langle\boldsymbol{\Omega}_{-12}\rangle$ in (5.2). Applying the Cauchy-Schwarz inequality, we find that

$$\begin{aligned} B &= \text{Cov}[\widehat{\text{tr}}_1, \widehat{\text{tr}}_2] = \mathbb{E}[(\widehat{\text{tr}}_1 - \text{tr } \mathbf{A} + \text{tr } \mathbf{X} - \boldsymbol{\omega}_1^* \mathbf{X} \boldsymbol{\omega}_1)(\widehat{\text{tr}}_2 - \text{tr } \mathbf{A} + \text{tr } \mathbf{X} - \boldsymbol{\omega}_2^* \mathbf{X} \boldsymbol{\omega}_2)] \\ &\leq \mathbb{E}|\widehat{\text{tr}}_1 - \text{tr } \mathbf{A} + \text{tr } \mathbf{X} - \boldsymbol{\omega}_1^* \mathbf{X} \boldsymbol{\omega}_1|^2 \\ &= \mathbb{E} \text{Var}[\boldsymbol{\omega}_1^* [\mathbf{A}\langle\boldsymbol{\Omega}_{-12}\rangle - \mathbf{A}\langle\boldsymbol{\Omega}_{-1}\rangle] \boldsymbol{\omega}_1 \mid \boldsymbol{\Omega}_{-1}] \\ &= 2 \mathbb{E} \|\mathbf{A}\langle\boldsymbol{\Omega}_{-1}\rangle - \mathbf{A}\langle\boldsymbol{\Omega}_{-12}\rangle\|_{\text{F}}^2. \end{aligned}$$

The psd matrix $\mathbf{A}\langle\boldsymbol{\Omega}_{-1}\rangle - \mathbf{A}\langle\boldsymbol{\Omega}_{-12}\rangle$ has rank one, and it is bounded above by $\mathbf{A} - \mathbf{A}\langle\boldsymbol{\Omega}_{-12}\rangle$ in the psd order. Therefore,

$$B = \text{Cov}[\widehat{\text{tr}}_1, \widehat{\text{tr}}_2] \leq 2 \mathbb{E} \|\mathbf{A} - \mathbf{A}\langle\boldsymbol{\Omega}_{-12}\rangle\|^2.$$

Combine the displays to complete the proof. \square

5.4. Error bounds for low-rank approximations. To prove the main result, [Theorem 1.1](#), we need two auxiliary lemmas. First, we present error bounds for randomized SVD and randomized Nyström approximation, drawn from the recent paper [32].

LEMMA 5.4 (Randomized SVD and randomized Nyström error). *Fix a matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$, and draw a standard normal matrix $\boldsymbol{\Omega} \in \mathbb{R}^{N \times k}$. For any $r \leq k - 2$, the randomized SVD error is bounded by*

$$\begin{aligned} \mathbb{E} \|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)\mathbf{A}\|^2 &\leq \frac{k+r-1}{k-r-1} \left(\|\mathbf{A} - [\mathbf{A}]_r\|^2 + \frac{e^2}{k-r} \|\mathbf{A} - [\mathbf{A}]_r\|_{\text{F}}^2 \right), \\ \mathbb{E} \|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)\mathbf{A}\|_{\text{F}}^2 &\leq \frac{k-1}{k-r-1} \|\mathbf{A} - [\mathbf{A}]_r\|_{\text{F}}^2, \end{aligned}$$

where $\mathbf{Q} = \text{orth}(\mathbf{A}\boldsymbol{\Omega})$.

Assume that $\mathbf{A} \in \mathbb{R}^{N \times N}$ is a psd matrix. For any $r \leq k - 4$, the randomized Nyström error is bounded by

$$\begin{aligned} (\mathbb{E} \|\mathbf{A} - \mathbf{A}\langle\boldsymbol{\Omega}\rangle\|^2)^{1/2} &\leq \frac{k+r-1}{k-r-3} \left(\|\mathbf{A} - [\mathbf{A}]_r\| + \frac{\sqrt{3}e^2}{k-r} \|\mathbf{A} - [\mathbf{A}]_r\|_* \right), \\ (\mathbb{E} \|\mathbf{A} - \mathbf{A}\langle\boldsymbol{\Omega}\rangle\|_{\text{F}}^2)^{1/2} &\leq \frac{k-2}{k-r-3} \left(\|\mathbf{A} - [\mathbf{A}]_r\|_{\text{F}} + \frac{1}{\sqrt{k-r}} \|\mathbf{A} - [\mathbf{A}]_r\|_* \right). \end{aligned}$$

Second, we report a standard fact about the decay rate of the singular values, which is also exploited in [27, Lem. 13] and [13, Lem. 7]. We omit the easy proof.

FACT 5.5. *For any matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ and any $r \geq 1$,*

$$\|\mathbf{A} - [\mathbf{A}]_r\| \leq \frac{\|\mathbf{A}\|_*}{r+1}, \quad \|\mathbf{A} - [\mathbf{A}]_r\|_{\text{F}} \leq \frac{\|\mathbf{A}\|_*}{2\sqrt{r}}.$$

5.5. The complete variance bound. To establish the main result, [Theorem 1.1](#), we begin with the initial variance bounds and introduce the results from [Lemma 5.4](#) and [Fact 5.5](#).

Proof of [Theorem 1.1](#). We recognize that all the terms in the error formulas in [Propositions 5.1](#) to [5.3](#) reflect the squared approximation error in a randomized SVD or a randomized Nyström approximation. Therefore, we can apply the error bounds in [Lemma 5.4](#) to obtain more explicit error representations. For HUTCH++, when $r \leq m/3 - 2$,

$$\mathbb{E} |\widehat{\text{tr}}_{\text{H}++} - \text{tr } \mathbf{A}|^2 \leq \frac{2}{m/3} \mathbb{E} \|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)\mathbf{A}\|_{\text{F}}^2 \leq \frac{2}{m/3 - r - 1} \|\mathbf{A} - [\mathbf{A}]_r\|_{\text{F}}^2.$$

DRAFT

For XTRACE, when $r \leq m/2 - 4$,

$$\begin{aligned} \mathbb{E} |\widehat{\text{tr}}_X - \text{tr } \mathbf{A}|^2 &\leq \frac{2}{m/2} \mathbb{E} \|(\mathbf{I} - \mathbf{Q}_{(1)} \mathbf{Q}_{(1)}^*) \mathbf{A}\|_{\text{F}}^2 + 4 \mathbb{E} \|(\mathbf{I} - \mathbf{Q}_{(12)} \mathbf{Q}_{(12)}^*) \mathbf{A}\|_{\text{F}}^2 \\ &\leq \frac{4m}{m/2 - r - 3} \|\mathbf{A} - \llbracket \mathbf{A} \rrbracket_r\|^2 + \frac{4e^2 m}{(m/2 - r - 3)^2} \|\mathbf{A} - \llbracket \mathbf{A} \rrbracket_r\|_{\text{F}}^2. \end{aligned}$$

Last, for XNYSTRACE, when $r \leq m - 6$,

$$\begin{aligned} (\mathbb{E} |\widehat{\text{tr}}_{\text{XN}} - \text{tr } \mathbf{A}|^2)^{1/2} &\leq \frac{\sqrt{2}}{\sqrt{m}} (\mathbb{E} \|\mathbf{A} - \mathbf{A} \langle \boldsymbol{\Omega}_{-1} \rangle\|_{\text{F}}^2)^{1/2} + \sqrt{2} (\mathbb{E} \|\mathbf{A} - \mathbf{A} \langle \boldsymbol{\Omega}_{-12} \rangle\|_{\text{F}}^2)^{1/2} \\ &\leq \frac{\sqrt{8}m}{m - r - 5} \|\mathbf{A} - \llbracket \mathbf{A} \rrbracket_r\| + \frac{\sqrt{2}m}{(m - r - 5)^{3/2}} \|\mathbf{A} - \llbracket \mathbf{A} \rrbracket_r\|_{\text{F}} + \frac{5e^2 m}{(m - r - 5)^2} \|\mathbf{A} - \llbracket \mathbf{A} \rrbracket_r\|_*. \end{aligned}$$

Thus, we confirm the detailed error bounds in [Theorem 1.1](#).

All that remains is to show that each trace estimator $\widehat{\text{tr}}$ satisfies

$$(5.5) \quad (\mathbb{E} |\widehat{\text{tr}} - \text{tr } \mathbf{A}|^2)^{1/2} \leq \frac{C}{m} \|\mathbf{A}\|_*,$$

for an absolute constant C . To that end, apply [Fact 5.5](#) to bound $\|\mathbf{A} - \llbracket \mathbf{A} \rrbracket_r\|$ and $\|\mathbf{A} - \llbracket \mathbf{A} \rrbracket_r\|_{\text{F}}$ in terms of $\|\mathbf{A} - \llbracket \mathbf{A} \rrbracket_r\|_*$. For HUTCH++, we set $r = \lfloor m/6 \rfloor - 1$. For XTRACE, we set $r = \lfloor m/4 \rfloor - 2$. For XNYSTRACE, we set $r = \lfloor m/2 \rfloor - 3$. Simplifying yields (5.5) for each estimator. \square

5.6. Proof of [Proposition 3.1](#). Last, we must argue that the posterior error estimator $\widehat{\text{er}}$ defined in (3.1) reflects the actual error. We instate the notation from [Proposition 3.1](#).

For both XTRACE and XNYSTRACE, each individual trace estimate $\widehat{\text{tr}}_i$ is unbiased. As a consequence, the variance takes the form

$$\mathbb{E} |\widehat{\text{tr}} - \text{tr } \mathbf{A}|^2 = \text{Var}(\widehat{\text{tr}}) = \frac{1}{\ell^2} \sum_{i,j=1}^{\ell} \text{Cov}(\widehat{\text{tr}}_i, \widehat{\text{tr}}_j).$$

A short calculation yields

$$\mathbb{E} \widehat{\text{er}}^2 = \frac{1}{\ell^2(\ell - 1)} \sum_{i,j=1}^{\ell} [\text{Var}(\widehat{\text{tr}}_j) - \text{Cov}(\widehat{\text{tr}}_i, \widehat{\text{tr}}_j)].$$

Since the samples $\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_\ell$ are exchangeable, the variance is the same for each j and the covariance is the same for all $i \neq j$. Therefore,

$$\begin{aligned} \mathbb{E} |\widehat{\text{tr}} - \text{tr}(\mathbf{A})|^2 &= \frac{1}{\ell} \text{Var}(\widehat{\text{tr}}_1) + \frac{\ell - 1}{\ell} \text{Cov}(\widehat{\text{tr}}_1, \widehat{\text{tr}}_2), \\ \mathbb{E} \widehat{\text{er}}^2 &= \frac{1}{\ell} \text{Var}(\widehat{\text{tr}}_1) - \frac{1}{\ell} \text{Cov}(\widehat{\text{tr}}_1, \widehat{\text{tr}}_2). \end{aligned}$$

The result follows when we take the ratio of these two quantities and simplify. \square

As a final comment, we observe that the calculations in [subsection 5.1](#) show for *symmetric matrices* that the XTRACE correlations are bounded by

$$\text{Cor}(\widehat{\text{tr}}_1, \widehat{\text{tr}}_2) \leq 2 \frac{\mathbb{E} \|(\mathbf{I} - \mathbf{Q}_{(12)} \mathbf{Q}_{(12)}^*) \mathbf{A} (\mathbf{I} - \mathbf{Q}_{(12)} \mathbf{Q}_{(12)}^*)\|_{\text{F}}^2}{\mathbb{E} \|(\mathbf{I} - \mathbf{Q}_{(1)} \mathbf{Q}_{(1)}^*) \mathbf{A} (\mathbf{I} - \mathbf{Q}_{(1)} \mathbf{Q}_{(1)}^*)\|_{\text{F}}^2},$$

where $\mathbf{Q}_{(1)}$ and $\mathbf{Q}_{(12)}$ are defined in [Proposition 5.1](#). These correlations are small for matrices with slow rates of singular value decay, i.e., when $\|\mathbf{A} - \llbracket \mathbf{A} \rrbracket_{m/2-1}\|_{\text{F}} \gg \|\mathbf{A} - \llbracket \mathbf{A} \rrbracket_{m/2-2}\|_{\text{F}}$.

DRAFT

In practice, we observe the correlations to be small even for matrices with singular values which decay more quickly. As an example, for the matrix with exponentially decaying eigenvalues in Figure 1, the XTRACE correlations (measured over 10^4 independent runs of the algorithm) are no higher than 0.06 and the average error estimate is correct up to a factor of 1.2.

Acknowledgments. We thank Eitan Levin for helpful discussions regarding the fast implementation of XTRACE.

REFERENCES

- [1] M. AL HASAN AND V. S. DAVE, *Triangle counting in large networks: A review*, WIREs Data Mining and Knowledge Discovery, 8 (2018), p. e1226, <https://doi.org/10.1002/widm.1226>. 15
- [2] A. H. AL-MOHY AND N. J. HIGHAM, *Computing the action of the matrix exponential, with an application to exponential integrators*, SIAM Journal on Scientific Computing, 33 (2011), pp. 488–511, <https://doi.org/10.1137/100788860>. 6
- [3] H. AVRON, P. MAYMOUNKOV, AND S. TOLEDO, *Blendenpik: Supercharging LAPACK's least-squares solver*, SIAM Journal on Scientific Computing, 32 (2010), pp. 1217–1236, <https://doi.org/10.1137/090767911>. 1
- [4] R. A. BASTON AND Y. NAKATSUKASA, *Stochastic diagonal estimation: Probabilistic bounds and an improved algorithm*, Jan. 2022, <https://arxiv.org/abs/2201.10684>. 11, 16
- [5] C. BEKAS, E. KOKIOPOULOU, AND Y. SAAD, *An estimator for the diagonal of a matrix*, Applied Numerical Mathematics, 57 (2007), pp. 1214–1229, <https://doi.org/10.1016/j.apnum.2007.01.003>. 10
- [6] D. BU, Y. ZHAO, L. CAI, H. XUE, X. ZHU, H. LU, J. ZHANG, S. SUN, L. LING, N. ZHANG, G. LI, AND R. CHEN, *Topological structure analysis of the protein–protein interaction network in budding yeast*, Nucleic Acids Research, 31 (2003), pp. 2443–2450, <https://doi.org/10.1093/nar/gkg340>. 16
- [7] T. CHEN AND E. HALLMAN, *Krylov-aware stochastic trace estimation*, May 2022, <https://arxiv.org/abs/2205.01736>. 9
- [8] T. CHEN, T. TROGDON, AND S. UBARU, *Randomized matrix-free quadrature for spectrum and spectral sum approximation*, Apr. 2022, <https://arxiv.org/abs/2204.01941>. 9
- [9] T. DAVIS AND Y. HU, *The University of Florida sparse matrix collection*, ACM Transactions on Mathematical Software, 38 (2011), pp. 1–25, <https://doi.org/10.1145/2049662.2049663>. 16
- [10] E. N. EPPERLY AND J. A. TROPP, *Jackknife variability estimation for randomized matrix computations*, July 2022, <https://arxiv.org/abs/2207.06342v2>. 9
- [11] E. ESTRADA, *The many facets of the Estrada indices of graphs and networks*, SeMA Journal, 79 (2022), pp. 57–125, <https://doi.org/10.1007/s40324-021-00275-w>. 15
- [12] A. S. GAMBHIR, A. STATHOPOULOS, AND K. ORGINOS, *Deflation as a method of variance reduction for estimating the trace of a matrix inverse*, SIAM Journal on Scientific Computing, 39 (2017), pp. A532–A558, <https://doi.org/10.1137/16M1066361>. 2, 8
- [13] A. C. GILBERT, M. J. STRAUSS, J. A. TROPP, AND R. VERSHYNIN, *One sketch for all: Fast algorithms for compressed sensing*, in Proceedings of the Thirty-Ninth Annual ACM Symposium on Theory of Computing, June 2007, pp. 237–246, <https://doi.org/10.1145/1250790.1250824>. 19
- [14] A. GIRARD, *A fast “Monte-Carlo cross-validation” procedure for large least squares problems with noisy data*, Numerische Mathematik, 56 (1989), pp. 1–23, <https://doi.org/10.1007/BF01395775>. 1, 2, 8
- [15] N. HALKO, P.-G. MARTINSSON, AND J. A. TROPP, *Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions*, SIAM Review, 53 (2011), pp. 217–288, <https://doi.org/10.1137/090771806>. 1, 3
- [16] P. R. HALMOS, *The theory of unbiased estimation*, The Annals of Mathematical Statistics, 17 (1946), pp. 34–43, <https://doi.org/10.1214/aoms/1177731020>. 2
- [17] N. J. HIGHAM, *Functions of Matrices: Theory and Computation*, SIAM, Philadelphia, 2008, <https://doi.org/10.1137/1.9780898717778>. 9
- [18] N. J. HIGHAM, *Matrix exponential times a vector*, Nov. 2010. Available at <https://www.mathworks.com/matlabcentral/fileexchange/29576-matrix-exponential-times-a-vector> (accessed 10/25/2022). 6
- [19] M. F. HUTCHINSON, *A stochastic estimator of the trace of the influence matrix for Laplacian smoothing splines*, Communications in Statistics - Simulation and Computation, 18 (1989), pp. 1059–1076, <https://doi.org/10.1080/03610918908812806>. 1, 2, 8
- [20] S. JIANG, H. PHAM, D. P. WOODRUFF, AND Q. ZHANG, *Optimal sketching for trace estimation*, in 35th Conference on Neural Information Processing Systems., 2021, p. 13. 8
- [21] V. S. KOROLJUK AND Y. V. BOROVSKICH, *Theory of U-Statistics*, Springer Netherlands, Dordrecht, 1994, <https://doi.org/10.1007/978-94-017-3515-5>. 2

DRAFT

- [22] 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 Transactions On Mathematical Software, 43 (2017), <https://doi.org/10.1145/3004053>. 10
- [23] L. LIN, *Randomized estimation of spectral densities of large matrices made accurate*, Numerische Mathematik, 136 (2017), pp. 183–213, <https://doi.org/10.1007/s00211-016-0837-7>. 2, 8
- [24] C. LITENS, *Transverse Field Ising Model with Different Boundary Conditions*, Bachelor’s Thesis, Stockholm University, Feb. 2019. Available at http://staff.fysik.su.se/~ardonne/files/theses/bachelor-thesis_christopher-litens.pdf. 15
- [25] P.-G. MARTINSSON AND J. A. TROPP, *Randomized numerical linear algebra: Foundations and algorithms*, Acta Numerica, 29 (2020), pp. 403–572, <https://doi.org/10.1017/S0962492920000021>. 2, 5, 8
- [26] R. A. MEYER, C. MUSCO, C. MUSCO, AND D. P. WOODRUFF, *Hutch++: Optimal stochastic trace estimation*, in Symposium on Simplicity in Algorithms, SIAM, Jan. 2021. 1, 2, 3, 8
- [27] R. A. MEYER, C. MUSCO, C. MUSCO, AND D. P. WOODRUFF, *Hutch++: Optimal Stochastic Trace Estimation*, June 2021, <https://arxiv.org/abs/2002.11457>. 18, 19
- [28] D. PERSSON, A. CORTINOVIS, AND D. KRESSNER, *Improved variants of the Hutch++ algorithm for trace estimation*, SIAM Journal on Matrix Analysis and Applications, (2022), pp. 1162–1185, <https://doi.org/10.1137/21M1447623>. 5, 8, 12, 23
- [29] P. PFEUTY, *The one-dimensional Ising model with a transverse field*, Annals of Physics, 57 (1970), pp. 79–90, [https://doi.org/10.1016/0003-4916\(70\)90270-8](https://doi.org/10.1016/0003-4916(70)90270-8). 15
- [30] V. ROKHLIN AND M. TYGERT, *A fast randomized algorithm for overdetermined linear least-squares regression*, Proceedings of the National Academy of Sciences, 105 (2008), pp. 13212–13217, <https://doi.org/10.1073/pnas.0804869105>. 1
- [31] A. K. SAIBABA, A. ALEXANDERIAN, AND I. C. F. IPSEN, *Randomized matrix-free trace and log-determinant estimators*, Numerische Mathematik, 137 (2017), pp. 353–395, <https://doi.org/10.1007/s00211-017-0880-z>. 2, 8, 12
- [32] J. A. TROPP AND R. J. WEBBER, *Randomized algorithms for low-rank matrix approximation: Design, analysis, and applications*, Manuscript in preparation, (2023). 19
- [33] 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. 10
- [34] S. UBARU, J. CHEN, AND Y. SAAD, *Fast estimation of $\text{tr}(f(A))$ via stochastic Lanczos quadrature*, SIAM Journal on Matrix Analysis and Applications, 38 (2017), pp. 1075–1099, <https://doi.org/10.1137/16M1104974>. 8
- [35] S. UBARU AND Y. SAAD, *Applications of trace estimation techniques*, in High Performance Computing in Science and Engineering, 2017, https://doi.org/10.1007/978-3-319-97136-0_2. 2

SUPPLEMENTARY MATERIAL

SM1. Comparison with Adaptive HUTCH++. The adaptive HUTCH++ algorithm of [28] flexibly apportions test vectors between low-rank approximation and residual trace estimation, based on the estimated singular values of the input matrix. This algorithm interpolates between the Girard–Hutchinson estimator, the HUTCH++ estimator, and a purely low-rank approximation-based trace estimate. This algorithm also estimates the total number of matvecs m to meet an (absolute) error tolerance:

$$|\widehat{\text{tr}} - \text{tr}(\mathbf{A})| \leq \varepsilon_{\text{abs}} \quad \text{except with probability } 1 - 2\delta,$$

where δ is a parameter chosen by the user.

We compare XTRACE with adaptive HUTCH++ in [Figure SM1](#). To produce these plots, we first run adaptive HUTCH++ with error tolerances $\varepsilon_{\text{abs}} = \varepsilon \cdot \text{tr} \mathbf{A}$ and failure probability parameter $\delta := 1/10$. For each value of the error tolerance ε , we plot the mean error (solid line) and desired accuracy ε against the mean number m of matvecs required by the algorithm. We then run XTRACE with a variable number of matvecs m and report the mean error over 1000 trials. Overall, we find that XTRACE performs similarly adaptive HUTCH++ or better than HUTCH++ by up to an order of magnitude.

SM2. MATLAB implementations. We present MATLAB R2022b implementations of XTRACE, XNYS TRACE, and XDIAG in [Programs SM2.1](#) to [SM2.3](#).

DRAFT

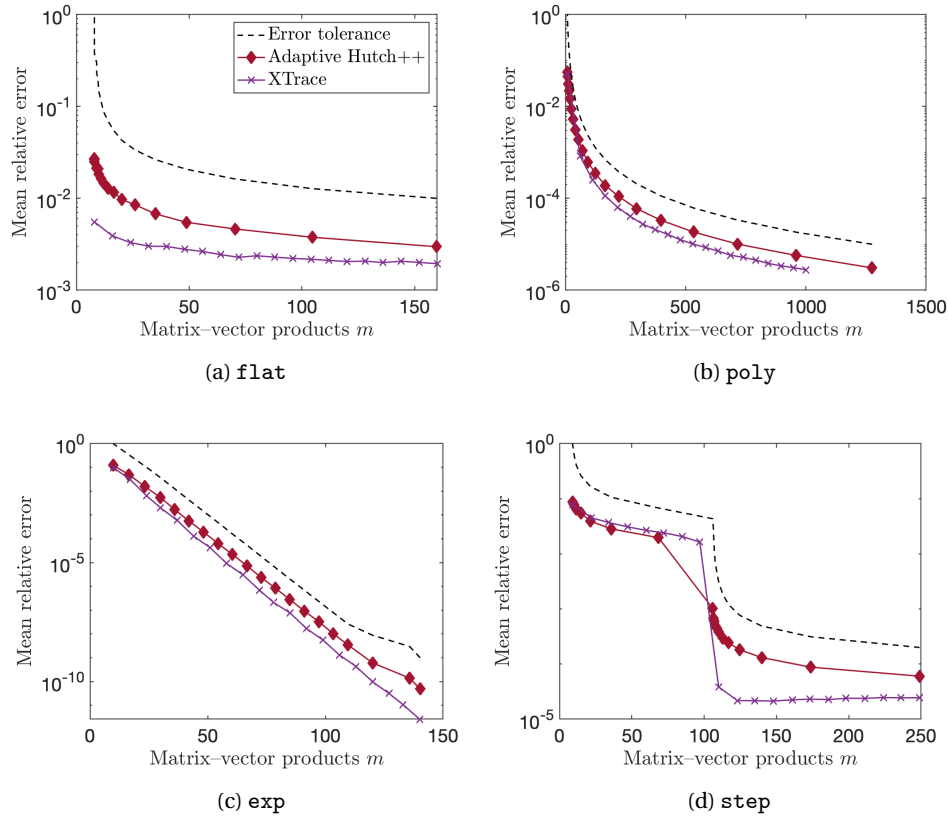


Fig. SM1: **Comparison with adaptive HUTCH++**. Average relative error using XTRACE and adaptive HUTCH++ (solid lines) and user-chosen error tolerance for Adaptive HUTCH++ (dashed lines). See [section SM1](#).

```

function [t,err] = xtrace(A, m)
N = size(A,1); m = floor(m/2);

%% Helper functions
cnormc = @(M) M ./ vecnorm(M,2,1);
diag_prod = @(A,B) sum(conj(A).*B,1).'; % computes diag(A'*B)

%% Randomized SVD
Om = sqrt(N)*cnormc(randn(N, m));
Y = A*Om;
[Q,R] = qr(Y,0);

%% Normalization
W = Q'*Om;
S = cnormc(inv(R)');
scale = (N - m + 1) ./ (N - (vecnorm(W)')^2 ...
    + abs(diag_prod(S,W) .* vecnorm(S)')^2);

%% Quantities needed for trace estimation
Z = A*Q; H = Q'*Z; HW = H*W; T=Z'*Om;
dSW = diag_prod(S, W); dSHS = diag_prod(S, H*S);
dTW = diag_prod(T, W); dWHW = diag_prod(W, HW);
dSRmHW = diag_prod(S, R-HW); dTmHRS = diag_prod(T-H'*W,S);

%% Trace estimate
ests = trace(H)*ones(m,1) - dSHS + (dWHW - dTW + dTmHRS .* dSW...
    + abs(dSW).^2 .* dSHS + conj(dSW) .* dSRmHW) .* scale;
t = mean(ests);
err = std(ests)/sqrt(m);
end

```

Program SM2.1: MATLAB 2022b implementation for XTRACE

DRAFT

```

function [t,err] = xnystrace(A, m)
N = size(A,1);

%% Helper functions
cnormc = @(M) M ./ vecnorm(M,2,1);
diag_prod = @(A,B) sum(conj(A).*B,1).'; % computes diag(A'*B)

%% Nystrom
Om = sqrt(N) * cnormc(randn(N,m));
Y = A*Om;
nu = eps*norm(Y,'fro')/sqrt(N);
Y = Y + nu*Om; % Shift for numerical stability
[Q,R] = qr(Y,0);
H = Om'*Y; C = chol((H+H')/2);
B = R/C; % Nystrom approx is Q*B*B'*Q'

%% Normalization
[QQ,RR] = qr(Om,0);
WW = QQ'*Om;
SS = cnormc(inv(RR)');
scale = (N - m + 1) ./ (N - vecnorm(WW).^2 ...
    + abs(diag_prod(SS,WW)') .* vecnorm(SS).^2);

%% Trace estimate
W = Q'*Om; S = (B/C') .* (diag(inv(H))').^(-0.5);
dSW = diag_prod(S, W).';
ests = norm(B,'fro')^2 - vecnorm(S).^2 + abs(dSW).^2 .* scale...
    - nu*N;
t = mean(ests);
err = std(ests)/sqrt(m);
end

```

Program SM2.2: MATLAB 2022b implementation for XNYS TRACE

DRAFT

```

function d = xdiag(A, m)
N = size(A,1); m = floor(m/2);

%% Helper functions
cnormc = @(M) M ./ vecnorm(M,2,1);
diag_prod = @(A,B) sum(conj(A).*B,1).'; % computes diag(A'*B)

%% Randomized SVD
Om = -3 + 2*randi(2,n,m); % Random signs
Y = A*Om;
[Q,R] = qr(Y,0);

% Quantities needed for trace estimation
Z = A'*Q; T=Z'*Om;
S = cnormc(inv(R)');
dQZ = diag_prod(Q',Z'); dQSSZ = diag_prod((Q*S)',(Z*S)');
dOmQT = diag_prod(Om',(Q*T).'); dOmY = diag_prod(Om',Y. ');
dOmQSST = diag_prod(Om',(Q*S*diag(diag_prod(S,T))). ');

% Trace estimate
d = dQZ + (-dQSSZ+dOmY-dOmQT+dOmQSST)/m;
end

```

Program SM2.3: MATLAB 2022b implementation for XDIAG

DRAFT