

HOW ACCURATELY SHOULD I COMPUTE IMPLICIT MATRIX-VECTOR PRODUCTS WHEN APPLYING THE HUTCHINSON TRACE ESTIMATOR?

JIE CHEN*

Abstract. The Hutchinson estimator defines an estimate of the trace of a matrix M , based on a bilinear form with independent vectors y of zero-mean unit-variance uncorrelated entries. This technique is particularly useful when M is only implicitly given but the matrix-vector product My can be efficiently computed without M being explicitly formed. Well-known examples in practice are $M = A^{-1}$, and more generally, $M = f(A)$. We study in this paper the conditions under which the numerical error incurred in computing My is comparable with the statistical uncertainty caused by the randomness of y . With these conditions, we also derive the sufficient number of random vectors that guarantees a relative error bound given any desired probability. For the purpose of obtaining easily computable conditions, we focus on the use of random vectors consisting of normal variables, a precursor technique attributed to Girard by Hutchinson. As demonstrated in many practical scenarios, normal variables are as effective as symmetric Bernoulli variables (a more common definition under the name of Hutchinson), but are advantageous in that they enjoy a simultaneous estimation of the estimator variance.

Key words. Matrix trace, Hutchinson estimator, matrix inverse, matrix function

AMS subject classifications. 65C05, 65F10, 65F60

1. Introduction. The trace of a large, implicit matrix M finds many applications in scientific computing. In estimation theory, if A is the Fisher information matrix of an unbiased estimator, then the trace of $M = A^{-1}$ gives a lower bound of the total variance of the estimator (see the Cramér–Rao inequality; e.g., [8]). Similarly, the log-determinant of a covariance matrix A , which is equivalent to the trace of $M = \log A$, appears naturally in the maximization of Gaussian log-likelihoods [1, 29, 17]; across disciplines, this term serves as a barrier in interior point methods for solving semidefinite programs when A is the semidefinite constraint [30, 7]. In electronic structures, the trace of the Fermi–Dirac function

$$f_{\text{FD}}(A) = \left[I + \exp\left(\frac{A - \mu I}{kT}\right) \right]^{-1} \quad (1.1)$$

gives the average number of electrons in a quantum system at chemical potential μ and temperature T , where A is the discretized Hamiltonian and k is the Boltzmann constant [27, 6]. Additionally, applications in lattice quantum chromodynamics [3, 28, 32], density of states [5, 23, 22], and uncertainty quantification [4, 20] comprise a limited, yet informative, list that illustrates the importance of trace computation.

When the $n \times n$ matrix M is implicitly defined through a given matrix A , it may not be computationally economic, or even viable, to first form M before extracting the trace. If, on the other hand, matrix-vector products with M are relatively inexpensive to compute, then n such products suffice the recovery of the trace: $\text{tr}(M) = \sum_{i=1}^n e_i^T M e_i$, where e_i is the i th column of the identity matrix. If, however, n is so large that even n matrix-vector products are too expensive to form, most of the existing work approximates the trace based on the following stochastic approach.

THEOREM 1.1 (Hutchinson [19]). *Let $M \in \mathbb{R}^{n \times n}$ be symmetric and $Y \in \mathbb{R}^n$ be a multivariate random variable of zero mean and unit covariance. Then, for a sample*

*IBM T. J. Watson Research Center, Yorktown Heights, NY 10598. Email: chenjie@us.ibm.com

y of Y ,

$$\mathbb{E}[y^T M y] = \text{tr}(M) \quad \text{and} \quad \text{Var}(y^T M y) = 2 \text{tr}(M^2) + \sum_{i=1}^n (\mathbb{E}[Y_i^4] - 3) M_{ii}^2.$$

Remark. The theorem straightforwardly extends to more general cases of M . For example, if M is unsymmetric, one may symmetrize the matrix to obtain the same trace: $\text{tr}(M) = \text{tr}(M + M^T)/2$. As another example, if M is Hermitian but not real, then $\text{tr}(M) = \text{tr}(\Re(M))$.

Practical uses of Theorem 1.1 form a sample average with N iid samples y_i , $i = 1, \dots, N$, such that the variance is reduced by a factor of N :

$$\mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N y_i^T M y_i \right] = \mathbb{E}[y^T M y] \quad \text{with} \quad \text{Var} \left(\frac{1}{N} \sum_{i=1}^N y_i^T M y_i \right) = \frac{1}{N} \text{Var}(y^T M y). \quad (1.2)$$

This technique, while extremely useful, poses an often neglected issue on the numerical accuracy of the evaluation of the $M y_i$'s. Consider, for example, $M y_i = A^{-1} y_i$. To one extreme, if the stochastic approximation (1.2) is highly accurate, meaning that the variance is sufficiently small, and if A is so ill conditioned that even a stable direct method for solving $A x_i = y_i$ leaves a comparably large backward error, then the overall departure of the estimate from the truth $\text{tr}(M)$ is dominated by numerical bias. Whereas such a scenario rarely occurs in practice, the opposite scenario is certainly not uncommon: the estimate yields a moderate variance, in the sense that it agrees with the truth on a small number of digits. Then, it is of little use to solve linear systems $A x_i = y_i$ highly accurately (if possible). Instead, it suffices for one to use an iterative solver (possibly enhanced by using block iterations for improving convergence [24, 26, 13]) that terminates at a moderate residual.

1.1. Stochastic uncertainty and numerical error. Hence, the subject of this paper is to study the balance of stochastic uncertainty of the Hutchinson estimator and the numerical error incurred in the evaluation of M -vector products. We derive practical conditions that make the two sources of errors comparable. To this end, we first need to establish the concrete meaning of “comparable” on a statistical basis.

Let $h(y)$ be an unbiased estimator of some quantity μ , with variance σ as in

$$\mathbb{E}_y[h(y)] = \mu \quad \text{and} \quad \text{Var}_y(h(y)) = \sigma^2.$$

Moreover, let y_i , $i = 1, \dots, N$ be N independent samples from the same distribution and define

$$h_0 = \frac{1}{N} \sum_{i=1}^N h(y_i).$$

Clearly, h_0 as an estimator is also unbiased. The central limit theorem states that $\sqrt{N}(h_0 - \mu)$ converges to the normal distribution $\mathcal{N}(0, \sigma^2)$. Therefore, for large N , h_0 approximately follows $\mathcal{N}(\mu, \sigma^2/N)$. Because σ^2/N is nothing but the variance of h_0 , we have for any $\beta > 0$,

$$\Pr \left(|h_0 - \mu| \leq \beta \sqrt{\text{Var}(h_0)} \right) \approx \text{erf} \left(\frac{\beta}{\sqrt{2}} \right), \quad (1.3)$$

where erf is the error function. Although $\text{Var}(h_0)$ in (1.3) could be replaced by the sample variance of $h(y_i)$ scaled by N , in this paper we consider the case when the $h(y_i)$'s are not computed accurately. Then, instead, we suppose that an unbiased estimator $h_1 > 0$ of $\text{Var}(h_0)$ is available; that is,

$$\mathbb{E}[h_1] = \text{Var}(h_0).$$

Let \tilde{h}_0 be the computed result of h_0 with error Δh_0 (i.e., $h_0 = \tilde{h}_0 + \Delta h_0$). If the error can be controlled to within an α portion of the estimate of the standard deviation of h_0 :

$$|\Delta h_0| \leq \alpha \sqrt{h_1}, \quad \alpha > 0, \quad (1.4)$$

then, we can maintain a confidence interval for \tilde{h}_0 as in

$$\Pr\left(|\tilde{h}_0 - \mu| \leq \beta \sqrt{\text{Var}(h_0)} + \alpha \sqrt{h_1}\right) \geq \Pr\left(|h_0 - \mu| \leq \beta \sqrt{\text{Var}(h_0)}\right) \approx \text{erf}\left(\frac{\beta}{\sqrt{2}}\right). \quad (1.5)$$

Further, if the estimator h_1 is asymptotically normal (akin to the fact that h_0 is an average of iid samples), then by invoking the normal density again we have for all $\gamma > 0$,

$$\Pr\left(|h_1 - \text{Var}(h_0)| \leq \gamma \sqrt{\text{Var}(h_1)}\right) \approx \text{erf}\left(\frac{\gamma}{\sqrt{2}}\right).$$

This approximate equality translates to

$$\Pr\left(\sqrt{h_1} \leq \sqrt{\text{Var}(h_0)} + \gamma^{\frac{1}{2}} \text{Var}(h_1)^{\frac{1}{4}}\right) \gtrsim \frac{1}{2} \left(1 + \text{erf}\left(\frac{\gamma}{\sqrt{2}}\right)\right),$$

which, combined with (1.5), yields

$$\Pr\left(|\tilde{h}_0 - \mu| \leq (\beta + \alpha) \sqrt{\text{Var}(h_0)} + \alpha \gamma^{\frac{1}{2}} \text{Var}(h_1)^{\frac{1}{4}}\right) \gtrsim \frac{1}{2} \left(1 + \text{erf}\left(\frac{\gamma}{\sqrt{2}}\right)\right) \text{erf}\left(\frac{\beta}{\sqrt{2}}\right). \quad (1.6)$$

The numerical interpretation of the quantity $s = (\beta + \alpha) \sqrt{\text{Var}(h_0)} + \alpha \gamma^{\frac{1}{2}} \text{Var}(h_1)^{\frac{1}{4}}$ is that with probability at least $\rho \approx \frac{1}{2}(1 + \text{erf}(\gamma/\sqrt{2})) \text{erf}(\beta/\sqrt{2})$, the relative error of the actually obtainable numerical result \tilde{h}_0 is bounded by

$$\frac{|\tilde{h}_0 - \mu|}{|\mu|} \leq \frac{s}{|\mu|}.$$

Take the common three-sigma rule, for example, where $\beta = 3$ and $\gamma = 3$. We have a sufficiently high probability $\rho = 99.6\%$. Omitting the $\alpha \gamma^{\frac{1}{2}} \text{Var}(h_1)^{\frac{1}{4}}$ term in s , we also have the relative error bound

$$\frac{s}{|\mu|} \approx \frac{\beta + \alpha}{|\mu|} \sqrt{\text{Var}(h_0)} = \frac{\sigma}{|\mu|} \cdot \frac{3 + \alpha}{\sqrt{N}}.$$

This bound means that the asymptotics with respect to N is key to the accuracy of \tilde{h}_0 , whereas the magnitude of α (which controls the part of numerical error) plays a less significant role. Thus, one may safely consider α as large as 1.0. We note that the term $\alpha \gamma^{\frac{1}{2}} \text{Var}(h_1)^{\frac{1}{4}}$ omitted in this discussion is assumed to decay at least as fast as the other term $O(N^{-\frac{1}{2}})$, a subject we will return to shortly.

1.2. Organization. The central contribution of this work is the establishment of conditions that ensure (1.4). To materialize the estimators h_0 and h_1 , we first define the random vector Y in Theorem 1.1. As oppose to the common use of independent ± 1 's (symmetric Bernoulli variables) as the elements of Y , in Section 2, we justify that normal variables are often as effective. In fact, the invention of the use of normal variables was attributed to Girard [15] by Hutchinson [19], before the symmetric Bernoulli variables became popular. One advantage of symmetric Bernoulli variables is that they minimize the variance in (1.2); however, we show two examples, motivated by electronic structure calculations with matrices of structural decay [6], that demonstrate that symmetric Bernoulli variables often cannot improve the estimation accuracy over normal variables by even one digit. Normal variables, on the other hand, allow a simultaneous estimation of the variance with almost negligible cost, which is challenging for symmetric Bernoulli variables to achieve. A further advantage of working with normal variables is that the estimator h_1 fulfills asymptotic normality and $O(N^{-\frac{1}{2}})$ decay of the one-fourth power of variance, as required in the preceding discussions.

In Section 3, we establish the condition ensuring (1.4) for the case $M = A^{-1}$. The condition is with respect to the *absolute* residual in the solution of linear equations $Ax_i = y_i$. This condition can be straightforwardly used as the *absolute* residual tolerance in an iterative solver. Additionally, an example with symmetric tridiagonal matrices is shown. Similar to the example in the section that follows, these are “toy” matrices because many properties (e.g., $\text{tr}(A^{-1})$ and the estimator variance) can be analytically derived. The purpose of the toy examples, however, is to show the asymptotics with respect to the matrix size n and the sample size N and to give a flavor of the numerical results under randomness. A numerical example with matrices from applications is given two sections later.

In Section 4, we establish the condition ensuring (1.4) in a probabilistic manner for the case $M = f(A)$. A general approach for computing implicit matrix-vector products of the form $f(A)y$ is to replace f by an approximate function p (e.g., a polynomial or a rational function) such that the evaluation of $p(A)y$ renders to matrix-vector multiplications with A . This approach should bare no surprise since for the case of linear systems, a Krylov solver can be interpreted as building a polynomial that interpolates $f(x) = x^{-1}$ at the approximate eigenvalues of A . Different from the condition for $M = A^{-1}$, however, the condition here is with respect to the *relative* error of the approximant p in the uniform norm. This condition is straightforwardly applicable when p is a polynomial, such as in the approach proposed by Chen et al. [12], because the approximation error can be monitored without the knowledge of y and because $p(A)y$ can be evaluated accurately to machine precision. For rational or other approximations (see, e.g., [18]), one must take into account the numerical error in evaluating $p(A)y$ in addition to the approximation error of p . As before, we show a numerical example with Toeplitz matrices with structural decay to illustrate the use of the condition. These matrices are model matrices for electronic structures and the attainable relative error scales as $\Theta(n^{-\frac{1}{2}}N^{-\frac{1}{2}})$ with high probability.

We show further computational experiences in Section 5, by using the PARSEC collection¹ of matrices arising from density functional theory [10, 9]. The function f therein is defined based on the Fermi–Dirac function (1.1). In this case, the number N of samples is chosen to be 1,000 and the trace estimate is generally two to four

¹Available from the University of Florida Sparse Matrix Collection <https://www.cise.ufl.edu/research/sparse/matrices/>

digits accurate (with sufficiently high probability). Interestingly, the general trend of results suggests that the relative accuracy improves as the matrix size increases, even though the same N is used throughout. This observation agrees with that of the model matrices with structural decay in Section 4.

Related work, discussions, and concluding remarks are given in Section 6.

Note: In this paper we use the Bachmann–Landau notations $O(\cdot)$, $\Omega(\cdot)$, and $\Theta(\cdot)$ to mean “bounded above”, “bounded below”, and “bounded both above and below” for expressing orders of decay or growth.

2. Hutchinson estimator with normal variables. In this section, we justify the use of normal variables in the Hutchinson estimator and study the properties of the estimator, the variance of the estimator, and the estimator of the variance.

2.1. Normal v.s. symmetric Bernoulli. A symmetric Bernoulli variable is a discrete random variable that takes the values ± 1 with equal probabilities. The following result is straightforward.

COROLLARY 2.1 (Hutchinson [19]). *Under the condition of Theorem 1.1,*

1. *If the entries of Y are independent symmetric Bernoulli variables, then*

$$\text{Var}(y^T M y) = 2 \text{tr}(M^2) - 2 \sum_{i=1}^n M_{ii}^2.$$

This variance is the minimum among all possible distributions of Y .

2. *If $Y \sim \mathcal{N}(0, I)$, then*

$$\text{Var}(y^T M y) = 2 \text{tr}(M^2).$$

A direct consequence of the corollary is that the variance may vanish (when M is diagonal) for symmetric Bernoulli variables; but for normal variables, the standard-deviation-to-mean ratio admits a lower bound $\Omega(n^{-\frac{1}{2}})$.

PROPOSITION 2.2. *Let $M \in \mathbb{R}^{n \times n}$ be symmetric and $Y \sim \mathcal{N}(0, I_n)$. Then, for a sample y of Y ,*

$$\frac{\sqrt{\text{Var}(y^T M y)}}{|\mathbb{E}[y^T M y]|} \geq \sqrt{\frac{2}{n}}.$$

Proof. Let λ be the vector of eigenvalues of M . One obtains the inequality by noting that

$$\text{Var}(y^T M y) = 2 \text{tr}(M^2) = 2 \|\lambda\|_2^2, \quad |\mathbb{E}[y^T M y]| = |\text{tr}(M)| \leq \|\lambda\|_1,$$

and that $\|\lambda\|_1 \leq \sqrt{n} \|\lambda\|_2$. \square

As such, it may appear that normal variables are inferior to symmetric Bernoulli variables, because if the energy of the matrix (in the sense of Frobenius norm $\|M\|_F^2 = \text{tr}(M^2)$) is concentrated on the diagonal, then the latter will yield highly accurate estimates. In many scenarios, however, this is an impractical assumption. In the following, we show two examples, both of which entail a decaying structure, and demonstrate that the standard-deviation-to-mean ratio attains the rate $\Theta(n^{-\frac{1}{2}})$ in both estimators. To maintain clarity, we use subscript “N” to mean normal and “B” to mean symmetric Bernoulli.

Example: Toeplitz matrix with exponential decay. Consider $M_{ij} = \theta^{|i-j|}$ where $0 < \theta < 1$. Then, $\text{tr}(M) = n$ and

$$\text{tr}(M^2) = n \frac{1 + \theta^2}{1 - \theta^2} - 2 \frac{\theta^2(1 - \theta^{2n})}{(1 - \theta^2)^2}, \quad \sum_{i=1}^n M_{ii}^2 = n.$$

Thus,

$$\frac{\sqrt{\text{Var}_{\text{N}}(y^T M y)}}{|\mathbb{E}[y^T M y]|} = \sqrt{\frac{2}{n} \frac{1 + \theta^2}{1 - \theta^2}} + O\left(\frac{1}{n}\right),$$

and

$$\frac{\sqrt{\text{Var}_{\text{B}}(y^T M y)}}{|\mathbb{E}[y^T M y]|} = \sqrt{\frac{2}{n} \frac{2\theta^2}{1 - \theta^2}} + O\left(\frac{1}{n}\right).$$

Asymptotically, θ needs to be $\leq 1/\sqrt{199} \approx 0.07$ in order that $\sqrt{\text{Var}_{\text{B}}}$ is a factor of 10 smaller than $\sqrt{\text{Var}_{\text{N}}}$ (i.e., one more digit accurate under the same probability).

Example: Toeplitz matrix with algebraic decay. Consider $M_{ij} = |i - j + 1|^{-1}$. Then, $\text{tr}(M) = n$ and

$$\text{tr}(M^2) = -n + 2(n+1) \sum_{i=1}^n \frac{1}{i^2} - 2 \sum_{i=1}^n \frac{1}{i}, \quad \sum_{i=1}^n M_{ii}^2 = n.$$

By applying the inequalities

$$\frac{\pi^2}{6} - \frac{1}{n} < \sum_{i=1}^n \frac{1}{i^2} < \frac{\pi^2}{6} - \frac{1}{n+1} \quad \text{and} \quad \ln(n+1) < \sum_{i=1}^n \frac{1}{i} \leq \ln n + 1,$$

we obtain

$$\left(\frac{\pi^2}{3} - 1\right)n - \frac{2}{n} - 2 \ln n + \frac{\pi^2}{3} - 4 < \text{tr}(M^2) < \left(\frac{\pi^2}{3} - 1\right)n - 2 \ln(n+1) + \frac{\pi^2}{3} - 2.$$

Therefore,

$$\frac{\sqrt{\text{Var}_{\text{N}}(y^T M y)}}{|\mathbb{E}[y^T M y]|} = \sqrt{\frac{2}{n} \left(\frac{\pi^2}{3} - 1\right)} + O\left(\frac{\ln n}{n}\right),$$

and

$$\frac{\sqrt{\text{Var}_{\text{B}}(y^T M y)}}{|\mathbb{E}[y^T M y]|} = \sqrt{\frac{2}{n} \left(\frac{\pi^2}{3} - 2\right)} + O\left(\frac{\ln n}{n}\right).$$

Asymptotically, the ratio between $\sqrt{\text{Var}_{\text{B}}}$ and $\sqrt{\text{Var}_{\text{N}}}$ is approximately 0.75, which means that the relative error resulting from the use of symmetric Bernoulli variables is only slightly better than that of normal variables. Improvement on the number of accurate digits is impossible.

2.2. Estimator, variance, and estimator of variance. Structural decay is an important property in electronic structures [6]. As demonstrated above, for model matrices with such a property, normal variables are generally as effective as symmetric Bernoulli variables. An advantage of the former is that it allows for a simultaneous estimation of the variance with negligible costs. Thus, in this subsection, we define the variance estimator, which itself has a variance that in turn admits an estimator. The recurrence of estimator and variance interestingly repeats endlessly. From now on, we use the sample average in place of a single sample in the estimator.

PROPOSITION 2.3. *Let $M \in \mathbb{R}^{n \times n}$ be symmetric and $y_i, i = 1, \dots, N$, be random iid vectors from $\mathcal{N}(0, I_n)$. Then, for all $j = 0, 1, \dots$*

$$\text{Var} \left(\frac{2^{2^j-1}}{N^{2^j}} \sum_{i=1}^N y_i^T M^{2^j} y_i \right) = \mathbb{E} \left[\frac{2^{2^{j+1}-1}}{N^{2^{j+1}}} \sum_{i=1}^N y_i^T M^{2^{j+1}} y_i \right] = \frac{2^{2^{j+1}-1}}{N^{2^{j+1}-1}} \text{tr}(M^{2^{j+1}}).$$

Proof. With basic properties of the variance,

$$\begin{aligned} \text{Var} \left(\frac{2^{2^j-1}}{N^{2^j}} \sum_{i=1}^N y_i^T M^{2^j} y_i \right) &= \left(\frac{2^{2^j-1}}{N^{2^j-1}} \right)^2 \text{Var} \left(\frac{1}{N} \sum_{i=1}^N y_i^T M^{2^j} y_i \right) \\ &= \frac{1}{N} \left(\frac{2^{2^j-1}}{N^{2^j-1}} \right)^2 \text{Var} \left(y_1^T M^{2^j} y_1 \right). \end{aligned}$$

Invoking Corollary 2.1 followed by Theorem 1.1, we obtain

$$\text{Var} \left(y_1^T M^{2^j} y_1 \right) = 2 \text{tr} \left(M^{2^{j+1}} \right) = 2 \mathbb{E} \left[y_1^T M^{2^{j+1}} y_1 \right].$$

Then, together with

$$\mathbb{E} \left[y_1^T M^{2^{j+1}} y_1 \right] = \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N y_i^T M^{2^{j+1}} y_i \right],$$

we conclude the proposition. \square

Based on the proposition, we define

$$h_j := \frac{2^{2^j-1}}{N^{2^j}} \sum_{i=1}^N y_i^T M^{2^j} y_i, \quad j = 0, 1, \dots \quad (2.1)$$

This definition is consistent with the notation h_0 and h_1 introduced earlier in the introduction. Then, Equation (1.2) together with Proposition 2.3 state that h_0 is an estimator of $\text{tr}(M)$, h_1 is an estimator of the variance of h_0 , and generally, h_{j+1} is an estimator of the variance of h_j . Clearly, all estimators are unbiased. Moreover, h_1 is a sample average and thus fulfills asymptotic normality as $N \rightarrow \infty$. Its variance $\text{Var}(h_1) = 8N^{-3} \text{tr}(M^4)$ decays as $\Theta(N^{-3})$.

Regarding the decay, we have a further result.

THEOREM 2.4. *Denote by σ_{\min} and σ_{\max} the smallest the largest singular value of a symmetric nonzero matrix M , respectively. For all $j > 0$,*

$$\left(\frac{2\sigma_{\min}}{N} \right)^{2^j} \leq \frac{h_{j+1}}{|h_j|} \leq \left(\frac{2\sigma_{\max}}{N} \right)^{2^j}. \quad (2.2)$$

Additionally, when $j = 0$,

1. (2.2) holds when M is semidefinite.
2. The left half of (2.2) holds when M is indefinite.

Proof. When $j > 0$, $h_j > 0$. Write $z_i = M^{2^{j-1}} y_i$. Then,

$$h_j = \frac{2^{2^j-1}}{N^{2^j}} \sum_{i=1}^N \|z_i\|^2 \quad \text{and} \quad h_{j+1} = \frac{2^{2^{j+1}-1}}{N^{2^{j+1}}} \sum_{i=1}^N z_i^T M^{2^j} z_i.$$

Because for each i ,

$$\sigma_{\min}^{2^j} \|z_i\|^2 \leq z_i^T M^{2^j} z_i \leq \sigma_{\max}^{2^j} \|z_i\|^2,$$

summing over i we conclude (2.2).

When $j = 0$ and M is positive semidefinite, the same argument can be used to establish (2.2), by noting that z_i is well defined. When $j = 0$ and M is negative semidefinite, $|h_0| = -h_0$. Then, replacing M by $-M$ will prove (2.2).

We now proceed to the remaining case: $j = 0$ and M is indefinite. We concatenate all the vectors y_i to form a vector z , and duplicate M diagonally to form a matrix \widetilde{M} (i.e., \widetilde{M} is block diagonal and the diagonal blocks are M). Then,

$$|h_0| = \frac{1}{N} |z^T \widetilde{M} z| \quad \text{and} \quad h_1 = \frac{2}{N^2} z^T \widetilde{M}^2 z.$$

Because \widetilde{M} can be diagonalized and the unitary factor can be absorbed to z , we treat \widetilde{M} a diagonal matrix where the diagonal elements are the eigenvalues of M . Then,

$$\left(\frac{2\sigma_{\min}}{N} \right) |h_0| = \left(\frac{2\sigma_{\min}}{N^2} \right) \left| \sum_k \lambda_k z_k^2 \right| \leq \frac{2}{N^2} \sum_k \sigma_{\min} |\lambda_k z_k^2| \leq \frac{2}{N^2} \sum_k \lambda_k^2 z_k^2 = h_1,$$

where z_k are the elements of z and λ_k are the eigenvalues of M . This shows the left half of (2.2). \square

Remark. The right half of (2.2) may not hold when $j = 0$ and M is indefinite, because h_0 may attain 0.

An immediate consequence of Theorem 2.4 is that

$$\left(\frac{2\sigma_{\min}}{N} \right)^{2^j-2} h_1 \leq h_j \leq \left(\frac{2\sigma_{\max}}{N} \right)^{2^j-2} h_1 \quad \text{and} \quad \left(\frac{2\sigma_{\min}}{N} \right) |h_0| \leq h_1.$$

The significance of this result is three fold. First, when N is considered fixed and sufficiently large, h_j decreases doubly exponentially with respect to j ; such a decrease is faster than the exponential. Second, when N varies, h_j decreases as $\Theta(N^{-2^j+1})$ if $h_1 = \Theta(N^{-1})$. Such a decrease is a very-high-order algebraic decrease. Third, for all $j > 0$, the ratio $\sqrt{h_{j+1}}/h_j = \Theta(N^{-\frac{1}{2}})$ if h_1 is $\Theta(N^{-1})$. Moreover, $\sqrt{h_1}/|h_0|$ is at least $\Omega(N^{-\frac{1}{2}})$. Since the ratio quantifies the relative error of using h_j as an estimator of $\text{Var}(h_{j-1})$ when $j > 0$, and of $\text{tr}(M)$ when $j = 0$, the j -independent decrease rate implies that the quality of estimators h_j is similar across j , in the relative term. Then, in the absolute term, h_j is more and more accurate as j becomes large.

We illustrate an example for the last point. Suppose $\text{tr}(M) = 316.22$ and the estimate $h_0 = 314.70$ is two-digit accurate. The true standard deviation $\sqrt{\text{Var}(h_0)} = 1.99$. In practice, when we use h_1 to estimate $\text{Var}(h_0)$, we may obtain $\sqrt{h_1} = 1.90$, again two-digit accurate. In such a case, we may safely treat h_1 the ‘‘same’’ as

$\text{Var}(h_0)$ when establishing the confidence interval for h_0 , because the absolute difference 1.99–1.90 is very small compared with the true trace 316.22.

On closing this section, we note that later we will frequently refer to the quantity

$$\frac{\sqrt{\text{tr}(M^2)}}{|\text{tr}(M)|} = \frac{\sqrt{\text{Var}(h_0)}}{|h_0|} \sqrt{\frac{N}{2}}, \quad (2.3)$$

which is $\sqrt{N/2}$ times the standard-deviation-to-mean ratio. If this quantity scales as $\Theta(n^{-\frac{1}{2}})$ (cf. the lower bound in Proposition 2.2), then the standard-deviation-to-mean ratio $\sqrt{\text{Var}(h_0)}/|h_0|$ scales as $\Theta(n^{-\frac{1}{2}}N^{-\frac{1}{2}})$. This means that the relative error of the estimate of the trace decreases not only with the sample size N but also with the matrix size n .

3. Estimating $\text{tr}(A^{-1})$. In this section, we consider the case $M = A^{-1}$, where the My_i 's are computed through solving linear systems $Ax_i = y_i$. The following is the main result.

THEOREM 3.1. *Let $A \in \mathbb{R}^{n \times n}$ be symmetric nonsingular and $y_i, i = 1, 2, \dots, N$, be independent vectors from $\mathcal{N}(0, I_n)$. For each i , denote by $r_i = y_i - Ax_i$ the residual of the linear system with matrix A and right-hand side y_i , where x_i is an approximate solution. Decompose the trace estimator $h_0 = \tilde{h}_0 + \Delta h_0$, where*

$$h_0 = \frac{1}{N} \sum_{i=1}^N y_i^T A^{-1} y_i \quad \text{and} \quad \tilde{h}_0 = \frac{1}{N} \sum_{i=1}^N y_i^T x_i,$$

and let h_1 be the estimator of the variance of h_0 , i.e.,

$$h_1 = \frac{2}{N^2} \sum_{i=1}^N y_i^T A^{-2} y_i.$$

For any $\alpha > 0$, if

$$\|r_i\| \leq \alpha \sqrt{\frac{2}{N}} \quad \text{for all } i, \quad (3.1)$$

then $|\Delta h_0| \leq \alpha \sqrt{h_1}$.

Proof. We express Δh_0 in terms of r_i :

$$\Delta h_0 = \frac{1}{N} \sum_{i=1}^N y_i^T A^{-1} r_i.$$

Let z be the column concatenation of the vectors $A^{-1}y_i$ and similarly let r be the concatenation of the r_i 's. Then,

$$\Delta h_0 = \frac{1}{N} z^T r \quad \text{and} \quad h_1 = \frac{2}{N^2} z^T z.$$

By Cauchy–Schwarz, $|z^T r| \leq \|z\| \|r\|$. Therefore, if all vectors r_i satisfy (3.1), then $\|r\| \leq \alpha \sqrt{2}$. Thus,

$$|\Delta h_0| = \frac{1}{N} |z^T r| \leq \frac{1}{N} \|z\| \|r\| \leq \frac{\alpha \sqrt{2}}{N} \|z\| = \alpha \sqrt{h_1},$$

which concludes the theorem. \square

We note that the condition (3.1) concerns the absolute residual, but in software implementations, the tolerance of a Krylov solver is typically applied on the relative residual. Hence, care is called for when one uses a software. The following result indicates that the absolute residual is approximately \sqrt{n} times the relative one.

PROPOSITION 3.2. *If $y \sim \mathcal{N}(0, I_n)$, then*

$$\mathbb{E}[\|y\|] = \frac{n\Gamma\left(\frac{n+1}{2}\right)}{\sqrt{2}\Gamma\left(\frac{n}{2}+1\right)} \quad \text{and} \quad \lim_{n \rightarrow \infty} \frac{\mathbb{E}[\|y\|]}{\sqrt{n}} = 1.$$

Proof. The first equality is a straightforward calculation:

$$\begin{aligned} \mathbb{E}[\|y\|] &= \int_{\mathbb{R}^n} \frac{\|y\|}{(2\pi)^{n/2}} \exp\left(-\frac{\|y\|^2}{2}\right) dy && \text{by definition} \\ &= \frac{n\pi^{n/2}}{\Gamma\left(\frac{n}{2}+1\right)} \int_0^\infty \frac{r}{(2\pi)^{n/2}} \exp\left(-\frac{r^2}{2}\right) r^{n-1} dr && \text{spherical coordinate } r = \|y\| \\ &= \frac{n\sqrt{2}}{\Gamma\left(\frac{n}{2}+1\right)} \int_0^\infty \exp(-r^2) r^n dr && \text{change of variable } r/\sqrt{2} \rightarrow r \\ &= \frac{n\sqrt{2}}{\Gamma\left(\frac{n}{2}+1\right)} \cdot \frac{1}{2}\Gamma\left(\frac{n+1}{2}\right). \end{aligned}$$

The second equality follows from

$$\lim_{m \rightarrow \infty} \frac{\Gamma\left(m - \frac{1}{2}\right) \sqrt{m}}{\Gamma(m)} = 1. \quad \square$$

Remark. Clearly, because $\|y\|^2 \sim \chi_n^2$, we have $\mathbb{E}[\|y\|^2] = n$.

The condition (3.1) established in Theorem 3.1 carries a dependence on N . The following result gives a sufficient condition for N to achieve an (ϵ, η) -type of bound. The probability $1 - \eta$ therein comes from the asymptotic normality of h_0 and h_1 and is well approximated for a sufficiently large sample size N .

COROLLARY 3.3. *Under the conditions of Theorem 3.1, let*

$$\rho = \Pr\left(|h_0 - \text{tr}(A^{-1})| \leq \beta\sqrt{\text{Var}(h_0)}\right) \Pr\left(h_1 \leq \text{Var}(h_0) + \gamma\sqrt{\text{Var}(h_1)}\right).$$

Then, for any $\epsilon > 0$, we have

$$\Pr\left(|\tilde{h}_0 - \text{tr}(A^{-1})| \leq \epsilon |\text{tr}(A^{-1})|\right) \geq \rho, \quad (3.2)$$

if $\text{tr}(A^{-1}) \neq 0$ and if

$$N \geq \max\left\{\frac{2\alpha^4\gamma^2}{(\beta + \alpha)^4} \frac{\text{tr}(A^{-4})}{\text{tr}(A^{-2})^2}, \frac{8(\beta + \alpha)^2}{\epsilon^2} \frac{\text{tr}(A^{-2})}{\text{tr}(A^{-1})^2}\right\}. \quad (3.3)$$

Moreover, for any $\tau \in (0, 1)$ and $\eta \in (0, 1 - 2^{-1/\tau})$, we have

$$\rho \rightarrow 1 - \eta \quad \text{as} \quad N \rightarrow \infty,$$

if $\beta = \sqrt{2} \text{erf}^{-1}((1 - \eta)^{1-\tau})$ and $\gamma = \sqrt{2} \text{erf}^{-1}(2(1 - \eta)^\tau - 1)$.

Proof. When N holds as required, we have

$$(\beta + \alpha)\sqrt{\text{Var}(h_0)} = (\beta + \alpha)\sqrt{\frac{2\text{tr}(A^{-2})}{N}} \leq \frac{\epsilon}{2}|\text{tr}(A^{-1})| \quad (3.4)$$

and

$$\alpha\gamma^{\frac{1}{2}}\text{Var}(h_1)^{\frac{1}{4}} = \alpha\gamma^{\frac{1}{2}}\left(\frac{8\text{tr}(A^{-4})}{N^3}\right)^{\frac{1}{4}} \leq (\beta + \alpha)\sqrt{\text{Var}(h_0)}.$$

Therefore,

$$(\beta + \alpha)\sqrt{\text{Var}(h_0)} + \alpha\gamma^{\frac{1}{2}}\text{Var}(h_1)^{\frac{1}{4}} \leq \epsilon|\text{tr}(A^{-1})|.$$

Then, with basic properties in probability,

$$\begin{aligned} & \Pr\left(|\tilde{h}_0 - \text{tr}(A^{-1})| \leq \epsilon|\text{tr}(A^{-1})|\right) \\ & \geq \Pr\left(|\tilde{h}_0 - \text{tr}(A^{-1})| \leq (\beta + \alpha)\sqrt{\text{Var}(h_0)} + \alpha\gamma^{\frac{1}{2}}\text{Var}(h_1)^{\frac{1}{4}}\right) \\ & \geq \Pr\left(|\tilde{h}_0 - \text{tr}(A^{-1})| \leq \beta\sqrt{\text{Var}(h_0)} + \alpha\sqrt{h_1}\right) \Pr\left(\sqrt{h_1} \leq \sqrt{\text{Var}(h_0)} + \gamma^{\frac{1}{2}}\text{Var}(h_1)^{\frac{1}{4}}\right). \end{aligned} \quad (3.5)$$

Thus, substituting

$$\Pr\left(|\tilde{h}_0 - \text{tr}(A^{-1})| \leq \beta\sqrt{\text{Var}(h_0)} + \alpha\sqrt{h_1}\right) \geq \Pr\left(|h_0 - \text{tr}(A^{-1})| \leq \beta\sqrt{\text{Var}(h_0)}\right)$$

(because of the result $|\Delta h_0| \leq \alpha\sqrt{h_1}$ in Theorem 3.1) and

$$\Pr\left(\sqrt{h_1} \leq \sqrt{\text{Var}(h_0)} + \gamma^{\frac{1}{2}}\text{Var}(h_1)^{\frac{1}{4}}\right) \geq \Pr\left(h_1 \leq \text{Var}(h_0) + \gamma\sqrt{\text{Var}(h_1)}\right)$$

to (3.5), we obtain (3.2).

Moreover, based on Proposition 2.3 and definition (2.1), h_0 and h_1 are both asymptotically normal; that is,

$$\lim_{N \rightarrow \infty} \Pr\left(|h_0 - \text{tr}(A^{-1})| \leq \beta\sqrt{\text{Var}(h_0)}\right) = \text{erf}\left(\frac{\beta}{\sqrt{2}}\right) \quad (3.6)$$

and

$$\lim_{N \rightarrow \infty} \Pr\left(h_1 \leq \text{Var}(h_0) + \gamma\sqrt{\text{Var}(h_1)}\right) = \frac{1}{2}\left(1 + \text{erf}\left(\frac{\gamma}{\sqrt{2}}\right)\right). \quad (3.7)$$

Note that β and γ are defined such that

$$1 - \eta = \frac{1}{2}\left(1 + \text{erf}\left(\frac{\gamma}{\sqrt{2}}\right)\right)\text{erf}\left(\frac{\beta}{\sqrt{2}}\right). \quad (3.8)$$

Then, (3.6)–(3.8) lead to $\rho \rightarrow 1 - \eta$. \square

Remark. By the central limit theorem, the convergence $\rho \rightarrow 1 - \eta$ may be read as “ ρ is approximately $1 - \eta$ when N is sufficiently large.” To gain insight from the corollary, let us consider the simplified setting of no numerical error. In this case,

Remark. In the inequalities of the third part, the right half is generally tight but not the left half. The rationale follows a technical detail in the proof, which is given in the appendix. As a consequence,

$$\frac{\sqrt{\text{tr}(A^{-2})}}{|\text{tr}(A^{-1})|} \approx \sqrt{\frac{1}{n} \frac{\zeta^2 + 1}{\zeta^2 - 1}}.$$

With the analytic understanding of $\text{tr}(A^{-1})$ and $\text{tr}(A^{-2})$, we now show numerical results for the following six quantities:

- (a) $\text{tr}(A^{-1})$ (b) h_0 (c) h_0 , iterative solve (`tol`)
 (d) $\sqrt{\frac{2}{N} \text{tr}(A^{-2})}$ (e) $\sqrt{h_1}$ (f) $\sqrt{h_1}$, iterative solve (`tol`).

See Table 3.1.

TABLE 3.1

Computational results for the tridiagonal matrix A defined in (3.9). Parameters: $n = 1000$, $N = 100$, $\alpha = 1$. Residual tolerance computed from (3.1): $\text{tol} = 1.41\text{e-}01$, $E[\text{rtol}] = 4.47\text{e-}03$.

Case $a = 2$, average residual = 1.23e-01,			
average relative residual = 3.90e-03			
	Truth	Estim. (full solve)	Estim. (<code>tol</code>)
$\text{tr}(A^{-1})$	167000	175960	175708
<code>stddev(estim)</code>	14937	15507	15470
Case $a = 1.7$, average residual = 9.75e-02,			
average relative residual = 3.08e-03			
	Truth	Estim. (full solve)	Estim. (<code>tol</code>)
$\text{tr}(A^{-1})$	1138.61	1030.96	1029.46
<code>stddev(estim)</code>	209.47	201.93	201.72
Case $a = 2.6$, average residual = 9.19e-02,			
average relative residual = 2.91e-03			
	Truth	Estim. (full solve)	Estim. (<code>tol</code>)
$\text{tr}(A^{-1})$	601.589	600.135	600.088
<code>stddev(estim)</code>	3.365	3.358	3.358

Let us recall the meaning of these quantities. Term (a) is the truth. Term (b) is the estimate, with standard deviation in Term (d). In the context of this paper, we approximate Term (b) by using Term (c), which is computed approximately based on the condition (3.1) in Theorem 3.1; and Term (e) is the estimator of Term (d). We approximate Term (e) by using Term (f), a byproduct of the calculation of Term (c).

We pick three values of a for demonstrating the numerical results, each corresponding to one part of Proposition 3.4. Because the case $a = 2$ corresponds to the standard 1D Laplacian, whose eigenvalues lie in the interval $(0, 4)$, we easily see that A is positive definite when $a = 2$ and 2.6, but is indefinite when $a = 1.7$. In light of the indefiniteness, we use GMRES as the linear solver and block Jacobi as the preconditioner.

When the sample size N is 100, the condition (3.1) indicates that we need only an absolute residual tolerance $1.41\text{e-}01$, which corresponds to an average relative

residual tolerance $4.47\text{e-}03$ for a matrix of size $n = 1,000$. In such a setting, the results in Table 3.1 indicate that the estimates are one- (close to two-) digit accurate for $a = 2$ and $a = 1.7$, but are two- (close to three-) digit accurate for $a = 2.6$. The absolute errors all happen to be within one times the standard deviation.

4. Estimating $\text{tr}(f(A))$. In this section, we consider the case $M = f(A)$, where the function f is approximated by p . The following is the main result.

THEOREM 4.1. *Let $A \in \mathbb{R}^{n \times n}$ be symmetric, f and p be functions well defined on the spectrum of A , and $y_i, i = 1, 2, \dots, N$, be independent vectors from $\mathcal{N}(0, I_n)$. Decompose the trace estimator $h_0 = \tilde{h}_0 + \Delta h_0$, where*

$$h_0 = \frac{1}{N} \sum_{i=1}^N y_i^T f(A) y_i \quad \text{and} \quad \tilde{h}_0 = \frac{1}{N} \sum_{i=1}^N y_i^T p(A) y_i,$$

and let h_1 be the estimator of the variance of h_0 , i.e.,

$$h_1 = \frac{2}{N^2} \sum_{i=1}^N y_i^T f(A)^2 y_i.$$

For any $\alpha > 0$ and $\delta \in (0, 1/2)$, if

$$|1 - p(\lambda)/f(\lambda)| \leq \alpha \sqrt{\frac{2}{(1 + \delta)nN}} \quad (4.1)$$

for all eigenvalues λ of A , then $|\Delta h_0| \leq \alpha \sqrt{h_1}$ with probability at least $1 - e^{-\delta^2 nN/6}$.

The proof of the theorem relies on the following lemma, whose result appears in various slightly different forms; see, e.g., [14, 21].

LEMMA 4.2. *Let $z \sim \mathcal{N}(0, I_d)$. Then, for any $\delta \in (0, 1/2)$,*

$$\Pr \left(\|z\|^2 \leq (1 + \delta)d \right) \geq 1 - e^{-\delta^2 d/6}.$$

Proof of Theorem 4.1. When $f(A)$ is nonsingular, write

$$\Delta h_0 = \frac{1}{N} \sum_{i=1}^N y_i^T f(A) (I - f(A)^{-1} p(A)) y_i = \frac{1}{N} \sum_{i=1}^N y_i^T f(A) r_i,$$

where $r_i = (I - f(A)^{-1} p(A)) y_i$. Let z be the column concatenation of the vectors $f(A) y_i$ and similarly let r be the concatenation of r_i . Then,

$$\Delta h_0 = \frac{1}{N} z^T r \quad \text{and} \quad h_1 = \frac{2}{N^2} z^T z.$$

If (4.1) is satisfied, we have

$$\|r\| \leq \alpha \sqrt{\frac{2}{(1 + \delta)nN}} \|w\|,$$

where w is the column concatenation of the vectors y_i . By Lemma 4.2, with probability at least $1 - e^{-\delta^2 nN/6}$, we have $\|w\| \leq \sqrt{(1 + \delta)nN}$. Therefore, with at least this

probability, $\|r\| \leq \alpha\sqrt{2}$. Then, by Cauchy–Schwarz $|z^T r| \leq \|z\|\|r\|$, we immediately conclude the theorem. \square

Remark. The Theorem suffers no loss of generality when $f(\lambda) = 0$ for some eigenvalue λ . In such a case, because f is bounded within the spectrum interval, one may define $\tilde{f} = f + c$ for some constant c so that $\tilde{f}(\lambda) \neq 0$ for all eigenvalues. Then, the theorem applies to the new function \tilde{f} .

Different from the condition in the preceding section, the condition (4.1) here ensures $|\Delta h_0| \leq \alpha\sqrt{h_1}$ only in the probabilistic sense. Hence, the establishment of the confidence interval (1.5) for \tilde{h}_0 needs a modification:

$$\begin{aligned} & \Pr\left(|\tilde{h}_0 - \mu| \leq \beta\sqrt{\text{Var}(h_0)} + \alpha\sqrt{h_1}\right) \\ & \geq \Pr\left(|h_0 - \mu| \leq \beta\sqrt{\text{Var}(h_0)}\right) \Pr\left(|\Delta h_0| \leq \alpha\sqrt{h_1}\right) \approx (1 - e^{-\delta^2 nN/6}) \text{erf}\left(\frac{\beta}{\sqrt{2}}\right). \end{aligned}$$

This modification, however, is minor. For example, for a probability ρ , the quantity $(1 - e^{-\delta^2 nN/6}) \times \rho$ agrees with ρ in the first three digits as long as $\delta^2 nN/6 > 10$.

With the modification in mind, we provide an (ϵ, η) -type of bound similar to that in the preceding section.

COROLLARY 4.3. *Under the conditions of Theorem 4.1, let*

$$\rho = \Pr\left(|h_0 - \text{tr}(f(A))| \leq \beta\sqrt{\text{Var}(h_0)}\right) \Pr\left(h_1 \leq \text{Var}(h_0) + \gamma\sqrt{\text{Var}(h_1)}\right).$$

Then, for any $\epsilon > 0$ and $c \in (0, 1)$, we have

$$\Pr\left(|\tilde{h}_0 - \text{tr}(f(A))| \leq \epsilon|\text{tr}(f(A))|\right) \geq c\rho,$$

if $\text{tr}(f(A)) \neq 0$ and if

$$N \geq \max\left\{\frac{2\alpha^4\gamma^2}{(\beta + \alpha)^4} \frac{\text{tr}(f(A)^4)}{\text{tr}(f(A)^2)^2}, \frac{8(\beta + \alpha)^2}{\epsilon^2} \frac{\text{tr}(f(A)^2)}{\text{tr}(f(A))^2}, \frac{-6 \log(1 - c)}{n\delta^2}\right\}. \quad (4.2)$$

Moreover, for any $\tau \in (0, 1)$ and $\eta \in (0, 1 - 2^{-1/\tau})$, we have

$$c\rho \rightarrow 1 - \eta \quad \text{as } N \rightarrow \infty,$$

if $\beta = \sqrt{2} \text{erf}^{-1}(((1 - \eta)/c)^{1-\tau})$ and $\gamma = \sqrt{2} \text{erf}^{-1}(2((1 - \eta)/c)^\tau - 1)$.

Remark. In (4.2), recall that n is the size of the matrix and δ , which does not need to be too small, is a quantity introduced in Theorem 4.1.

Proof. The proof is analogous to that of Corollary 3.3; hence, we mention only the difference. When N holds as required, we have

$$\begin{aligned} & \Pr\left(|\tilde{h}_0 - \text{tr}(f(A))| \leq \epsilon|\text{tr}(f(A))|\right) \\ & \geq \Pr\left(|h_0 - \text{tr}(f(A))| \leq \beta\sqrt{\text{Var}(h_0)}\right) (1 - e^{-\delta^2 nN/6}) \\ & \quad \times \Pr\left(h_1 \leq \text{Var}(h_0) + \gamma\sqrt{\text{Var}(h_1)}\right) \geq c\rho. \end{aligned}$$

Moreover, β and γ are defined such that

$$1 - \eta = \frac{c}{2} \left(1 + \text{erf}\left(\frac{\gamma}{\sqrt{2}}\right)\right) \text{erf}\left(\frac{\beta}{\sqrt{2}}\right).$$

Hence, $c\rho \rightarrow 1 - \eta$ by the central limit theorem. \square

Remark. Similar to the second remark after Corollary 3.3, here we may relax the sufficient sample size N in (4.2) to

$$N \geq \max \left\{ \frac{2\alpha^4\gamma^2}{(\beta + \alpha)^4}, \frac{8(\beta + \alpha)^2}{\epsilon^2}, \frac{-6 \log(1 - c)}{n\delta^2} \right\},$$

if $f(A)$ is definite.

4.1. Example: Toeplitz matrices with decay. To demonstrate the use of Theorem 4.1, here we consider Toeplitz matrices with structural decay (exponential or algebraic). Similar to the example in the preceding section, decaying Toeplitz matrices enjoy interesting properties. As we will show later, the standard-deviation-to-mean ratio decreases as $\Theta(n^{-\frac{1}{2}}N^{-\frac{1}{2}})$ for any continuous function f . Hence, one expects that the estimate can be more accurate when the matrix becomes larger. As model matrices for electronic structures, they hint on the effective use of the trace estimator in this application.

Let $A_{ij} = t_{i-j}$, where for symmetry $t_k = t_{-k}$. Assume that the infinite sequence $\dots, t_{-2}, t_{-1}, t_0, t_1, t_2, \dots$ consists of the coefficients of the Fourier series of a 2π -periodic function $q(\omega)$ in that

$$q(\omega) = \sum_{k=-\infty}^{\infty} t_k e^{ik\omega} \quad \text{with} \quad t_k = \frac{1}{2\pi} \int_0^{2\pi} q(\omega) e^{-ik\omega} d\omega. \quad (4.3)$$

Then, for any length- n vector x ,

$$x^T A x = \frac{1}{2\pi} \int_0^{2\pi} q(\omega) \left| \sum_{j=1}^n x_j e^{-ij\omega} \right|^2 d\omega.$$

Because

$$\frac{1}{2\pi} \int_0^{2\pi} \left| \sum_{j=1}^n x_j e^{-ij\omega} \right|^2 d\omega = \|x\|_2^2,$$

if we choose an x with a unit norm, then

$$\inf_{\omega \in [0, 2\pi]} q(\omega) \leq x^T A x \leq \sup_{\omega \in [0, 2\pi]} q(\omega),$$

which means that the eigenvalues of A are bounded within the range of q .

We will add a subscript n to A when asymptotics is in concern. A useful result for Toeplitz matrices is that the eigenvalues of A_n , denoted as $\lambda_j^{(n)}$, $j = 0, \dots, n-1$, are close to equally-spaced samples of q . For this, we need the definition of *equal distribution*.

DEFINITION 4.4. *Two sets of real numbers $\{a_j^{(n)}\}_{j=0, \dots, n-1}$ and $\{b_j^{(n)}\}_{j=0, \dots, n-1}$ are equally distributed in the interval $[M_1, M_2]$ if for any continuous function $F : [M_1, M_2] \rightarrow \mathbb{R}$,*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=0}^{n-1} [F(a_j^{(n)}) - F(b_j^{(n)})] = 0.$$

It is well known that the eigenvalues $\{\lambda_j^{(n)}\}$ of A_n and the set $\{q(2\pi j/n)\}_{j=0, \dots, n-1}$ are equally distributed (see, e.g, [16, 11]). An immediate consequence is that for a matrix function f , if it is continuous on the range of q , then the trace of $f(A_n)$ and that of $f^2(A_n)$ can be characterized by

$$\lim_{n \rightarrow \infty} \frac{1}{n} \operatorname{tr}(f(A_n)) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=0}^{n-1} f(\lambda_j^{(n)}) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=0}^{n-1} f(q(2\pi j/n)) = \frac{1}{2\pi} \int_0^{2\pi} f(q(\omega)) d\omega$$

and

$$\lim_{n \rightarrow \infty} \frac{1}{n} \operatorname{tr}(f^2(A_n)) = \frac{1}{2\pi} \int_0^{2\pi} f^2(q(\omega)) d\omega.$$

Therefore, we have the following result.

THEOREM 4.5. *Given an infinite sequence $\dots, t_{-2}, t_{-1}, t_0, t_1, t_2, \dots$ that satisfies the symmetry $t_k = t_{-k}$ and the assumption (4.3) for some 2π -periodic function $q(\omega)$, define a sequence of matrices A_n , $n = 1, 2, \dots$, where $(A_n)_{ij} = t_{i-j}$. Then, for any f continuous on the range of q ,*

$$\lim_{n \rightarrow \infty} n^{\frac{1}{2}} \frac{\sqrt{\operatorname{tr}(f^2(A_n))}}{\operatorname{tr}(f(A_n))} = \frac{\left(\int_0^{2\pi} 2\pi f^2(q(\omega)) d\omega \right)^{\frac{1}{2}}}{\int_0^{2\pi} f(q(\omega)) d\omega}.$$

We now consider two examples.

Exponential decay. Let

$$t_k = \theta^{|k|}, \quad 0 < \theta < 1. \quad (4.4)$$

We have

$$\begin{aligned} q(\omega) &= \sum_{k=-\infty}^{\infty} t_k e^{ik\omega} = -1 + 2\Re \left(\sum_{k=0}^{\infty} e^{k \ln \theta} e^{ik\omega} \right) \\ &= -1 + 2\Re \left(\frac{1}{1 - e^{\ln \theta + i\omega}} \right) = \frac{1 - \theta^2}{1 - 2\theta \cos \omega + \theta^2}. \end{aligned}$$

Therefore,

$$q_{\max} = q(0) = \frac{1 + \theta}{1 - \theta} \quad \text{and} \quad q_{\min} = q(\pi) = \frac{1 - \theta}{1 + \theta}.$$

Algebraic decay. Let

$$t_k = (k^2 + 1)^{-1}. \quad (4.5)$$

Note that this decay has a different order from that in Section 2. We make use of the well known Fourier transform

$$\frac{2a}{a^2 + \omega^2} = \int_{-\infty}^{+\infty} e^{-a|t|} e^{-i\omega t} dt, \quad a > 0$$

to write

$$\begin{aligned} q(\omega) &= \sum_{k=-\infty}^{+\infty} \frac{e^{ik\omega}}{1+k^2} = \frac{1}{2} \sum_{k=-\infty}^{+\infty} \int_{-\infty}^{+\infty} e^{-|t|} e^{-ikt} dt e^{ik\omega} \\ &= \frac{1}{2} \int_{-\infty}^{+\infty} e^{-|t|} \left[\sum_{k=-\infty}^{+\infty} e^{-ikt} e^{ik\omega} \right] dt \\ &= \pi \int_{-\infty}^{+\infty} e^{-|t|} \text{III}(t - \omega) dt = \pi \sum_{j=-\infty}^{+\infty} e^{-|\omega+2\pi j|}, \end{aligned}$$

where III denotes the Dirac comb. Then, in the interval $[0, 2\pi]$,

$$q(\omega) = \frac{\pi e^{2\pi}}{e^{2\pi} - 1} e^{-\omega} + \frac{\pi}{e^{2\pi} - 1} e^{\omega}.$$

Therefore,

$$q_{\max} = q(0) = \pi \frac{e^{2\pi} + 1}{e^{2\pi} - 1} \quad \text{and} \quad q_{\min} = q(\pi) = \frac{2\pi e^{\pi}}{e^{2\pi} - 1}.$$

In both examples, q is positive and hence A is positive definite. We use the square root function $f(x) = \sqrt{x}$ to demonstrate numerical results of the following eight quantities:

- (a) $\text{tr}(f(A))$ (b) $\frac{n}{2\pi} \int f(q(\omega))$ (c) h_0 , exact f (d) h_0 , using p
 (e) $\sqrt{\frac{2}{N} \text{tr}(f^2(A))}$ (f) $\sqrt{\frac{n}{N\pi} \int f^2(q(\omega))}$ (g) $\sqrt{h_1}$, exact f (h) $\sqrt{h_1}$, using p .

See Table 4.1.

TABLE 4.1

Computational results for Toeplitz matrices defined in (4.4) and (4.5). Parameters: $n = 1000$, $N = 100$, $\alpha = 1$, $\delta = 0.1$. Tolerance on relative approximation error computed from (4.1): $\text{rtol} = 4.26\text{e-}03$.

Case: exponential decay, $\theta = 0.7$, interpolation interval $[0.176, 5.667]$,
 polynomial degree $k = 10$, relative approximation error = $3.40\text{e-}03$

	Truth	Approx. Truth	Estim. (exact f)	Estim. (rtol)
$\text{tr}(A^{1/2})$	839.299	839.122	837.889	837.944
$\text{stddev}(\text{estim})$	4.472	4.472	4.464	4.464

Case: algebraic decay, interpolation interval $[0.272, 3.154]$,
 polynomial degree $k = 6$, relative approximation error = $2.35\text{e-}03$

	Truth	Approx. Truth	Estim. (exact f)	Estim. (rtol)
$\text{tr}(A^{1/2})$	930.241	930.169	929.038	929.046
$\text{stddev}(\text{estim})$	4.472	4.472	4.465	4.465

Let us recall the meaning of these quantities. Term (a) is the truth. Term (c) is the estimate, with standard deviation in Term (e). Terms (b) and (f) are approximations of Terms (a) and (e), respectively, based on equal distributions at large n . In the context of this paper, we approximate Term (c) by using Term (d), where

the approximation $p \approx f$ achieves a relative error dictated by (4.1) of Theorem 4.1; and Term (g) is the estimator of Term (e). Moreover, we approximate Term (g) by using Term (h), a byproduct of the calculation of Term (d).

Because the matrix A is well conditioned (see the values of q_{\min} and q_{\max} previously analyzed), the square root function f can be well approximated by using the simple Chebyshev interpolation. Let k be the degree of the interpolating polynomial p that interpolates f at the (shifted) Chebyshev nodes in the interval $[q_{\min}, q_{\max}]$. These nodes are the roots of the (shifted) Chebyshev polynomial T_{k+1} of degree $k+1$. Then, the relative approximation error $\max |1 - p/f|$, which is required in (4.1), can be well approximated by the maximum of $|1 - p/f|$ evaluated at the extrema of T_{k+1} , because these nodes interleave with the Chebyshev nodes.

When the sample size N is 100, the condition (4.1) indicates that we need only a relative error 4.26e-03. In such a setting, the results in Table 4.1 indicate that the estimates yield a relative error between 10^{-2} and 10^{-3} . The absolute errors are all within the standard deviation.

5. Further numerical examples in electronic structures. In this section, we show further computational results by using the PARSEC collection of matrices arising from density functional theory. The function f is a simple scaling and shift of the Fermi-Dirac function f_{FD} in (1.1):

$$f(x) = 2f_{\text{FD}}(x) - 1 = \frac{1 - \exp[\beta(x - \mu)]}{1 + \exp[\beta(x - \mu)]}, \quad \beta = (kT)^{-1} > 0. \quad (5.1)$$

The reason of performing this transformation is that f_{FD} approaches 0 as $x \rightarrow \infty$. Then, the relative approximation error is hard to control were f_{FD} used as the function. According to the remark of Theorem 4.1, we perform the transformation to resolve this challenge. Clearly, the function f is nothing but the negative hyperbolic tangent: $f(x) = -\tanh[\beta(x - \mu)/2]$.

We set β to be a large number 100 so that f appears close to the negative sign function. We set the chemical potential μ to lie at 1/3 of the spectrum interval, and scale the shifted matrix $A - \mu I$ such that its spectral radius is 1. Hence, function approximation is carried out in the unit interval $[-1, 1]$. The extreme eigenvalues of A are computed by using the Lanczos method. Numerical results are shown in Table 5.1.

Because of the steep slope of f at the origin, Chebyshev interpolation is no longer the best choice. We use the approach proposed by Chen et al. [12] to perform the spline/polynomial approximation. In this approach, f is first approximated by a cubic spline, where the knots are at geometrically progressing locations $\pm\theta^k$ together with the origin (for Table 5.1 we set $\theta = 9/10$ and $k = 0, 1, \dots, 99$). Then, the spline is in turn approximated by a least squares polynomial, where the inner product is defined as the sum of the $(1 - x^2)^{-\frac{1}{2}}$ -weighted inner products in each subinterval. The relative approximation error is approximated by the maximum of $|1 - p/f|$ evaluated at the mid-points of the spline knots. Under such a scheme, the degree of the polynomial p which satisfies the relative error tolerance dictated by (4.1) is on the order of several hundreds (see the columns “Degree” and “rtol” in Table 5.1).

Unlike the examples in the preceding sections, where a number of quantities can be computed because of the known expressions and the small size of the matrix, here we compute only the (approximate) h_0 and h_1 . We increase the sample size N to 1,000 for a more accurate estimation. One observes from the table that overall, the

TABLE 5.1

Computational results for the PARSEC collection of matrices. The function f is defined in (5.1). Parameters: $N = 1000$, $\alpha = 1$, $\delta = 0.1$.

Matrix	n	rtol	Degree	h_0	$\sqrt{h_1}$
Si2	769	1.54e-03	241	-269.096	1.210
SiH4	5,041	6.01e-04	273	-2049.98	3.09
SiNa	5,743	5.63e-04	275	-2324.90	3.30
Na5	5,832	5.58e-04	275	-2006.35	3.32
benzene	8,219	4.70e-04	281	-2817.46	3.94
Si10H16	17,077	3.26e-04	291	-6694.64	5.70
Si5H12	19,896	3.02e-04	295	-7290.60	6.14
SiO	33,401	2.33e-04	303	-12601.3	7.9
Ga3As3H12	61,349	1.72e-04	313	61321.1	11.0
GaAsH6	61,349	1.72e-04	313	61334.8	11.0
H2O	67,024	1.65e-04	315	-22269.1	11.2
Si34H36	97,569	1.37e-04	321	-36723.1	13.6
Ge87H76	112,985	1.27e-04	323	-43704.8	14.6
Ge99H100	112,985	1.27e-04	323	-43371.0	14.6
Ga10As10H30	113,081	1.27e-04	323	113020.	15.
Ga19As19H42	133,123	1.17e-04	327	133034.	16.
SiO2	155,331	1.08e-04	329	-52984.2	17.1
Si41Ge41H72	185,639	9.90e-05	333	-67395.9	18.7
CO	221,119	9.07e-05	335	-81894.4	20.4
Si87H76	240,369	8.70e-05	337	-89993.4	21.3
Ga41As41H72	268,096	8.24e-05	339	267859.	23.

estimates are two to four digits accurate and the accuracy improves when the matrix size becomes larger.

6. Concluding remarks. Computing the trace of a large, implicit matrix M has diverse interesting applications in scientific computing. The Hutchinson trace estimator is a matrix-free approach that makes use of efficient M -vector multiplications to remedy the expensive cost of the explicit construction of M . We have studied the effect of the numerical error in the evaluation of M -vector products. In particular, we derive conditions (3.1) and (4.1) that ensure that the numerical error is comparable to the uncertainty of the stochastic estimation. These conditions are readily applicable as a computational guidance for the approximate evaluation of M -vector products. Combined with these conditions, we also give in Corollary 3.3 and 4.3 the sufficient sample size that guarantees a relative error bound given any desired probability. Several examples with special matrices and matrices from an application demonstrate the effective use of the conditions. Particularly, in the experiments with the PARSEC collection of matrices from density functional theory, we observe that the relative error of the estimation decreases as the matrix size n increases, a qualitative agreement with the theoretical order $\Theta(n^{-\frac{1}{2}}N^{-\frac{1}{2}})$ for model matrices with structural decay.

6.1. Other trace estimators. Many additional methods and variants exist for the trace computation. Avron and Toledo [2] proposed and analyzed several distributions where the random vectors are drawn. The effectiveness of the resulting estimators, as long as they are unbiased, may be measured by the estimator variance, because the central limit theorem ensures that confidence intervals can be established

by treating the sample average approximately Gaussian for large N . In that vein, variance reduction is a valuable resort. For example, Stein et al. [29] proposed using dependent random vectors in groups so that the variance of any diagonal block of M corresponding to the group is eliminated. The grouping of the rows and columns of M , in this case, respects the closeness of spatial data in order that the eliminated variance majorizes the remained covariance in the off-diagonal blocks.

Interestingly, an opposite approach for grouping, coined “probing,” was proposed as well. In this approach, rows and columns far apart are grouped together. The rationale is simple. Consider, for the moment, that M is tridiagonal. It suffices to use three deterministic vectors

$$y_1 = \sum_{i=0}^{\lfloor (n-1)/3 \rfloor} e_{3i+1}, \quad y_2 = \sum_{i=0}^{\lfloor (n-2)/3 \rfloor} e_{3i+2}, \quad y_3 = \sum_{i=0}^{\lfloor (n-3)/3 \rfloor} e_{3i+3}$$

to exactly recover the trace:

$$\text{tr}(M) = y_1^T M y_1 + y_2^T M y_2 + y_3^T M y_3,$$

because $e_i^T M e_j$ vanishes whenever $|i - j| \geq 3$. In other words, every other three columns (and rows) of M are grouped together. Then, the diagonal blocks of the permuted matrix is diagonal. Hence, each deterministic vector is used to compute the trace of one block. Such a technique can be generalized for a sparse matrix by coloring the graph representation of the matrix, so that no two same-colored nodes share neighbors [5]. In practice, however, the sparsity pattern of an implicit matrix M is unknown; but if M is the inverse of A whose sparsity is given, it is often a reasonable assumption that the magnitude of M_{ij} decreases as the distance between nodes i and j in the graph of A increases. Hence, the heuristic is to group graph nodes that are a certain distance apart [28]. In all these methods, if the variance of the resulting estimator is formulated, the technique for deriving computational conditions similar to those in this paper is possibly transferable.

6.2. Bounds. In addition to the computable conditions, in this paper we also offer a sufficient sample size N that, together with these conditions, guarantees (ϵ, η) -type of probability bounds (see Corollary 3.3 and 4.3). Such bounds are generally framed in the form [2, 25, 17, 31]

$$\Pr \left(|\tilde{h}_0 - \mu| \leq \epsilon |\mu| \right) \geq 1 - \eta,$$

where \tilde{h}_0 is an estimation (possibly with numerical error) of the truth μ . For this bound to hold, the required sample size N grows as $\Theta(\epsilon^{-2} \log(1/\eta))$. Our results for N are similar, but the inequality of the bounds holds only approximately due to the use of asymptotic normality arguments. Nevertheless, the approximation becomes highly accurate when N is sufficiently large. Moreover, an advantage of our analysis is that it applies to any nonsingular symmetric matrix M , whereas other work generally relies on the presumption of definiteness. Such an advantage becomes a necessity in applications (e.g., electronic structures) where the implicit matrix M is almost always indefinite.

6.3. Bias caused by numerical error. Although the estimator h_0 is unbiased, the actually obtainable result \tilde{h}_0 may not be, unless the numerical error $\Delta h_0 = h_0 - \tilde{h}_0$ has a zero mean. In some applications, e.g., lattice quantum chromodynamics [28],

the traces of a sequence of A^{-1} are computed. If we allow the linear solve errors to be at the same order as the statistical uncertainty, averaging over all these traces creates a systematic bias, because the numerical errors could be in similar directions. Hence, it might be preferable to reduce α in (1.4) to, say, the order of $\Theta(m^{-1})$ or $\Theta(m^{-\frac{1}{2}})$, where m is the number of traces computed. The theory developed in this paper fully covers such a case, because in principle α can be any positive number, although a moderate value (such as 1.0) appears to be sufficient were there no concern on the numerical bias.

6.4. The lower bound of Proposition 2.2. The Toeplitz examples in Section 2 state that for the Hutchinson estimator with normal variables, the asymptotic lower bound $\Omega(n^{-\frac{1}{2}})$ of the standard-deviation-to-mean ratio is attainable. One of the referees inspired us to formulate another example that demonstrates the same bound. This example concerns definite matrices with a uniformly bounded condition number.

PROPOSITION 6.1. *Let $M \in \mathbb{R}^{n \times n}$ be an element of a sequence of symmetric definite matrices of increasing sizes, whose condition numbers are upper bounded by κ independent of the matrix size, and let $Y \sim \mathcal{N}(0, I_n)$. Then, for a sample y of Y ,*

$$\sqrt{\frac{2}{n}} \leq \frac{\sqrt{\text{Var}(y^T M y)}}{|\mathbb{E}[y^T M y]|} \leq \kappa \sqrt{\frac{2}{n}}.$$

Proof. The left inequality is a restatement of Proposition 2.2. For the right inequality, let σ_{\max} and σ_{\min} be the largest and smallest singular values of M , respectively. One obtains the inequality by noting that

$$\text{Var}(y^T M y) = 2 \text{tr}(M^2) \leq 2n\sigma_{\max}^2 \quad \text{and} \quad |\mathbb{E}[y^T M y]| = |\text{tr}(M)| \geq n\sigma_{\min}. \quad \square$$

Acknowledgments. The author is indebted to the editor Andreas Stathopoulos and two anonymous referees for their constructive comments that have substantially improved the paper. This work was supported by the XDATA program of the Defense Advanced Research Projects Agency (DARPA), administered through Air Force Research Laboratory contract FA8750-12-C-0323.

Appendix. Proof of Proposition 3.4. We split the proof into four cases.

Case $a = 2$. The recurrence relation (3.10) solves to $u_i = i$ for all i . Then,

$$(A^{-1})_{ij} = \min(i, j) - \frac{ij}{n+1}, \quad i, j = 1, \dots, n.$$

Summing over i (and j), we obtain

$$\text{tr}(A^{-1}) = \frac{n^2 + 2n}{6} \quad \text{and} \quad \text{tr}(A^{-2}) = \frac{2n^4 + 8n^3 + 17n^2 + 18n}{180}.$$

Case $a = -2$. The recurrence relation (3.10) solves to $u_i = (-1)^{i+1}i$ for all i . Then,

$$(A^{-1})_{ij} = (-1)^{i-j+1} \left[\min(i, j) - \frac{ij}{n+1} \right], \quad i, j = 1, \dots, n.$$

Summing over i (and j), we obtain

$$\text{tr}(A^{-1}) = -\frac{n^2 + 2n}{6} \quad \text{and} \quad \text{tr}(A^{-2}) = \frac{2n^4 + 8n^3 + 17n^2 + 18n}{180}.$$

Case $|a| < 2$. Based on the definition of θ , if $(n+1)\theta$ is a multiple of π , then A is singular; otherwise, the recurrence relation (3.10) solves to

$$u_i = \frac{2 \sin(i\theta)}{\sqrt{4-a^2}} \quad \text{for all } i.$$

Hence, the elements of the lower triangular part of A^{-1} are

$$(A^{-1})_{ij} = \frac{2 \sin[(n+1-i)\theta] \sin(j\theta)}{\sqrt{4-a^2} \sin[(n+1)\theta]}, \quad i = 1, \dots, n, \quad j = 1, \dots, i.$$

With tedious algebraic calculations, we obtain

$$\text{tr}(A^{-1}) = \frac{-(n+1) \cot[(n+1)\theta] + \cot \theta}{\sqrt{4-a^2}}$$

and

$$\text{tr}(A^{-2}) = \frac{(n+1)^2 \csc^2[(n+1)\theta] + (n+1) \cot[(n+1)\theta] \cot \theta + 1 - 2 \csc^2 \theta}{4-a^2}.$$

Case $|a| > 2$. The recurrence relation (3.10) solves to

$$u_i = \frac{1}{\sqrt{a^2-4}} \left[\left(\frac{a + \sqrt{a^2-4}}{2} \right)^i - \left(\frac{a - \sqrt{a^2-4}}{2} \right)^i \right] \quad \text{for all } i.$$

In what follows, we will give the details for only the scenario $a > 2$. The details of the other scenario $a < -2$ are analogous. Based on the definition of ζ and δ , we note that $\zeta > 1$, $0 < \delta < 1$ and $u_i = (1 - \delta^i) \zeta^i / \sqrt{a^2 - 4}$. Then, when $i \geq j$,

$$\frac{u_{n+1-i} u_j \sqrt{a^2-4}}{u_{n+1}} = \frac{(1 - \delta^{n+1-i})(1 - \delta^j)}{1 - \delta^{n+1}} \zeta^{j-i},$$

which is bounded on both sides by $(1 - \delta)^2 \zeta^{j-i}$ and ζ^{j-i} . Therefore,

$$(1 - \delta)^2 \zeta^{-|i-j|} \leq (A^{-1})_{ij} \sqrt{a^2-4} \leq \zeta^{-|i-j|}.$$

Summing over $i = j$, we obtain

$$\frac{(1 - \delta)^2 n}{\sqrt{a^2-4}} \leq \text{tr}(A^{-1}) \leq \frac{n}{\sqrt{a^2-4}};$$

and summing over i and j , we obtain

$$\frac{(1 - \delta)^4}{a^2 - 4} \left[n \frac{\zeta^2 + 1}{\zeta^2 - 1} - 2 \frac{\zeta^2(\zeta^{2n} - 1)}{(\zeta^2 - 1)^2 \zeta^{2n}} \right] \leq \text{tr}(A^{-2}) \leq \frac{1}{a^2 - 4} \left[n \frac{\zeta^2 + 1}{\zeta^2 - 1} - 2 \frac{\zeta^2(\zeta^{2n} - 1)}{(\zeta^2 - 1)^2 \zeta^{2n}} \right].$$

□

Remark. In the inequalities of the case $|a| > 2$, the right half is generally tight but not the left half. This is because when $\delta < 1$ is away from 1,

$$\frac{(1 - \delta^{n+1-i})(1 - \delta^j)}{1 - \delta^{n+1}}$$

is close to 1, unless when $i \approx n$ and $j \approx 1$, which makes it close to $(1 - \delta)^2$. However, the number of cases when $i \approx n$ and $j \approx 1$ is limited and thus when summing over i (and j), these cases contribute little to the overall sum. Therefore, the trace terms leaned toward the upper bounds. As a result,

$$\frac{\sqrt{\operatorname{tr}(A^{-2})}}{|\operatorname{tr}(A^{-1})|} \approx \sqrt{\frac{1}{n} \frac{\zeta^2 + 1}{\zeta^2 - 1}}.$$

REFERENCES

- [1] M. ANITESCU, J. CHEN, AND L. WANG, *A matrix-free approach for solving the parametric Gaussian process maximum likelihood problem*, SIAM J. Sci. Comput., 34 (2012), pp. A240–A262.
- [2] H. AVRON AND S. TOLEDO, *Randomized algorithms for estimating the trace of an implicit symmetric positive semi-definite matrix*, J. Assoc. Comput. Mach., 58 (2011).
- [3] Z. BAI, G. FAHEY, AND G. GOLUB, *Some large-scale matrix computation problems*, J. Comput. Appl. Math., 74 (1996), pp. 71–89.
- [4] C. BEKAS, A. CURIONI, AND I. FEDULOVA, *Low cost high performance uncertainty quantification*, in Proceedings of the 2nd Workshop on High Performance Computational Finance, 2009.
- [5] C. BEKAS, E. KOKIOPOULOU, AND Y. SAAD, *An estimator for the diagonal of a matrix*, Appl. Numer. Math., 57 (2007), pp. 1214–1229.
- [6] M. BENZI, P. BOITO, AND N. RAZOUK, *Decay properties of spectral projectors with applications to electronic structure*, SIAM Rev., 55 (2013), pp. 3–64.
- [7] S. BOYD AND L. VANDENBERGHE, *Convex Optimization*, Cambridge University Press, 2004.
- [8] G. CASELLA AND R. L. BERGER, *Statistical Inference*, Duxbury Press, 2nd ed., 2001.
- [9] J. CHELIKOWSKY, *The pseudopotential-density functional method applied to nanostructures*, J. Phys. D: Appl. Phys., 33 (2000), pp. R33–R50.
- [10] J. R. CHELIKOWSKY, N. TROULLIER, AND Y. SAAD, *Finite-difference-pseudopotential method: Electronic structure calculations without a basis*, Phys. Rev. Lett., 72 (1994), pp. 1240–1243.
- [11] J. CHEN, *On the use of discrete Laplace operator for preconditioning kernel matrices*, SIAM J. Sci. Comput., 36 (2013), pp. A289–A309.
- [12] J. CHEN, M. ANITESCU, AND Y. SAAD, *Computing $f(A)b$ via least squares polynomial approximations*, SIAM J. Sci. Comput., 33 (2011), pp. 195–222.
- [13] J. CHEN, T. L. H. LI, AND M. ANITESCU, *A parallel linear solver for multilevel Toeplitz systems with possibly several right-hand sides*, Parallel Comput., 40 (2014), pp. 408–424.
- [14] S. DASGUPTA AND A. GUPTA, *An elementary proof of a theorem of Johnson and Lindenstrauss*, Random Structures & Algorithms, 22 (2003), pp. 60–65.
- [15] D. GIRARD, *Un algorithme simple et rapide pour la validation croisé généralisée sur des problèmes de grande taille*, Tech. Rep. RR 669-M, Inf. et Math. Appl. de Grenoble, Grenoble, France, 1987.
- [16] R. M. GRAY, *Toeplitz and Circulant Matrices: A Review*, Now Publishers Inc, 2006.
- [17] I. HAN, D. MALIOUTOV, AND J. SHIN, *Large-scale log-determinant computation through stochastic Chebyshev expansions*, in Proceedings of the 32nd International Conference on Machine Learning, 2015.
- [18] N. J. HIGHAM, *Functions of Matrices: Theory and Computation*, SIAM, 2008.
- [19] M. F. HUTCHINSON, *A stochastic estimator of the trace of the influence matrix for Laplacian smoothing splines*, Communications in Statistics – Simulation and Computation, 19 (1990), pp. 433–450.
- [20] V. KALANTZIS, C. BEKAS, A. CURIONI, AND E. GALLOPOULOS, *Accelerating data uncertainty quantification by solving linear systems with multiple right-hand sides*, Numerical Algorithms, 62 (2013), pp. 637–653.
- [21] P. LI, T. J. HASTIE, AND K. W. CHURCH, *Nonlinear estimators and tail bounds for dimension reduction in l_1 using Cauchy random projections*, in Learning Theory, vol. 4539 of Lecture Notes in Computer Science, 2007, Springer Berlin Heidelberg, pp. 514–529.
- [22] L. LIN, *Randomized estimation of spectral densities of large matrices made accurate*. arXiv:1504.07690 [math.NA], 2015.

- [23] L. LIN, Y. SAAD, AND C. YANG, *Approximating spectral densities of large matrices*, SIAM Rev., (in press).
- [24] D. P. O'LEARY, *The block conjugate gradient algorithm and related methods*, Linear Algebra Appl., 29 (1980), pp. 293–322.
- [25] F. ROOSTA-KHORASANI AND U. ASCHER, *Improved bounds on sample size for implicit matrix trace estimators*, Foundations of Computational Mathematics, 15 (2015), pp. 1187–1212.
- [26] Y. SAAD, *Iterative Methods for Sparse Linear Systems*, SIAM, 2nd ed., 2003.
- [27] Y. SAAD, J. R. CHELIKOWSKY, AND S. M. SHONTZ, *Numerical methods for electronic structure calculations of materials*, SIAM Rev., 52 (2010), pp. 3–54.
- [28] A. STATHOPOULOS, J. LAEUCHLI, AND K. ORGINOS, *Hierarchical probing for estimating the trace of the matrix inverse on toroidal lattices*, SIAM J. Sci. Comput., 35 (2013), pp. S299–S322.
- [29] M. L. STEIN, J. CHEN, AND M. ANITESCU, *Stochastic approximation of score functions for Gaussian processes*, Annals of Applied Statistics, 7 (2013), pp. 1162–1191.
- [30] L. VANDENBERGHE AND S. BOYD, *Semidefinite programming*, SIAM Rev., 38 (1996), pp. 49–95.
- [31] K. WIMMER, Y. WU, AND P. ZHANG, *Optimal query complexity for estimating the trace of a matrix*, in Proceedings of the 41st International Colloquium on Automata, Languages and Programming, 2014.
- [32] L. WU, A. STATHOPOULOS, J. LAEUCHLI, V. KALANTZIS, AND E. GALLOPOULOS, *Estimating the trace of the matrix inverse by interpolating from the diagonal of an approximate inverse*, Journal of Computational Physics, (to appear).