

TRACE ESTIMATION AS AN ALGORITHMIC PRIMITIVE

CONNER DIPAOLO*

Abstract. This lecture note details motivation, results, and upper/lower sample complexity bounds for estimating the trace of large matrices when computing elements of \mathbf{A} is about as slow as computing $\langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle$ for vectors \mathbf{x} .

1. Motivation: Computing Schatten Norms. Let \mathbf{A} be any positive definite matrix in \mathbf{M}_n . We're going to assume that \mathbf{A} is real but these results largely carry over directly to the complex case. This matrix might represent the error of a low rank matrix factorization, or might specify something about the curvature of the loss surface of a deep neural network. In both cases, knowing how 'large' \mathbf{A} is, in some norm, would help us determine how good the factorization is or how drastically our loss surface is curving.

For whatever reason, let's suppose we want to measure \mathbf{A} in a Schatten- p norm, where $p = 1, 2, \dots < \infty$ is an integer.

DEFINITION 1.1. Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ be the singular values of \mathbf{A} (which are the same as the eigenvalues since \mathbf{A} is positive definite). The Schatten- p norm of \mathbf{A} is

$$\|\mathbf{A}\|_p^p = \sum_{i=1}^n \lambda_i^p = \text{tr}(\mathbf{A}^p).$$

The normal procedure for computing $\|\mathbf{A}\|_p$ goes as follows.

Algorithm 1.1 Naive Schatten norm computation

Data: A positive semi-definite matrix $\mathbf{A} \in \mathbf{M}_n$ and an integer $p = 1, 2, \dots < \infty$.

Result: The Schatten- p norm $\|\mathbf{A}\|_p$ using $\frac{4}{3}n^3 + np + 17 = \mathcal{O}(n^3 + np)$ flops.

Overwrite \mathbf{A} with $\text{diag}(\lambda_1, \dots, \lambda_n)$ using $\sim 4n^3/3$ flops. (Golub and Van Loan, 2012, Alg 8.3.3)

Compute $\mathbf{A} \leftarrow \mathbf{A}^p$ using $n(p-1)$ flops.

Compute $\text{tr } \mathbf{A} = \sum_{i=1}^n \mathbf{A}_{ii}$ using $n-1$ flops.

Return $\|\mathbf{A}\|_p^p = (\text{tr } \mathbf{A})^{1/p}$ using about 18 flops.

But what happens when n is large, say $n \geq 10,000$? Then the n^3 computation time can become unbearable because the eigendecomposition becomes slow. Indeed, on my 1.4GHz computer the estimated computation time for arbitrary \mathbf{A} is at least a day for $n = 50,000$ and $p = 2$. What can we do to speed this up? In particular, how can we drop the n^3 dependence?

Well, one way to do this is to note that for any \mathbf{x} , we can compute $\mathbf{A}\mathbf{x}$ in $2n^2 - n$ flops, and as a result we can compute $\mathbf{A}^p \mathbf{x}$ in $p(2n^2 - n) = 2n^2 p - np$ flops. In particular, we can compute $\langle \mathbf{A}^p \mathbf{x}, \mathbf{x} \rangle$ in $2n^2 p - np + 2n - 1$ flops. Now

$$\|\mathbf{A}\|_p^p = \text{tr } \mathbf{A}^p = \sum_{i=1}^n \mathbf{A}_{ii}^p = \sum_{i=1}^n \langle \mathbf{A}^p \mathbf{e}_i, \mathbf{e}_i \rangle,$$

so naive computation in this form gives an algorithm for computing the Schatten- p norm in $n(2n^2 p - np + 2n - 1) = 2n^3 p - n^2 p + 2n^2 - n + 18 = \mathcal{O}(n^3 p)$ flops. Immediately, this performs worse than our original Algorithm 1.1, but this provides an interesting interpretation of $\|\mathbf{A}\|_p^p$ as an expectation.

$$\|\mathbf{A}\|_p^p = n \mathbf{E} \langle \mathbf{A}^p \mathbf{e}_i, \mathbf{e}_i \rangle, \text{ where } i \sim \text{Uniform}\{1, 2, \dots, n\}.$$

If we take a hint from Monte-Carlo methods (as well as ideas from the world of concentration inequalities), this means that the following randomized algorithm should return a good *estimate* for $\|\mathbf{A}\|_p$ with high probability, while requiring significantly fewer flops for large n .

*Harvey Mudd College (cdipaolo@hmc.edu)

Algorithm 1.2 Naive randomized Schatten norm computation

Data: Positive semi-definite $\mathbf{A} \in \mathbb{M}_n$, integer $p = 1, 2, \dots < \infty$, and accuracy parameter $m \in \{1, 2, \dots\}$

Result: Estimated Schatten- p norm $S \approx \|\mathbf{A}\|_p$ using $2mn^2p - mnp + m + 19 = \mathcal{O}(mn^2p)$ flops.

Initialize $S \leftarrow 0$.

Initialize $\mathbf{x} \leftarrow 0 \in \mathbb{C}^n$.

for $1 \leq k \leq m$ **do**

 Sample an index $i \sim \text{Uniform}\{1, 2, \dots, n\}$.

 Let $\mathbf{x} \leftarrow \mathbf{e}_i$.

for $1 \leq j \leq p$ **do**

 Compute $\mathbf{x} \leftarrow \mathbf{A}\mathbf{x}$ using $2n^2 - n$

end

 Accumulate $S \leftarrow S + \mathbf{x}_i$ using 1 flop.

end

Let $S \leftarrow \frac{n}{m}S$ using 1 flop.

Return $\|\mathbf{A}\|_p \approx S^{1/p}$ using about 18 flops.

Actual bounds on how good of an estimate Algorithm 1.2 returns will follow from results in the later sections, and are provided explicitly in (Woodruff et al., 2014, Thm 69), but for now just remark how the simple realization that the trace is an expected inner product

$$\text{tr}(\mathbf{A}^p) = n\mathbf{E}\langle \mathbf{A}^p \mathbf{e}_i, \mathbf{e}_i \rangle = \mathbf{E}\langle \mathbf{A}^p \sqrt{n}\mathbf{e}_i, \sqrt{n}\mathbf{e}_i \rangle, \text{ where } i \sim \text{Uniform}\{1, 2, \dots, n\}.$$

allowed us to drop an n^3 dependence in our algorithm to a randomized algorithm with n^2 dependence without bringing in any heavy machinery. One thing we can say without any real work is that Algorithm 1.2 returns an estimate $S_m \rightarrow \|\mathbf{A}\|_p$ as $m \rightarrow \infty$ by the law of large numbers. (Wasserman, 2013, Thm 5.18) This follows since by positive definiteness $\mu = \mathbf{E}|\langle \mathbf{A}^p \sqrt{n}\mathbf{e}_i, \sqrt{n}\mathbf{e}_i \rangle| = \mathbf{E}\langle \mathbf{A}^p \sqrt{n}\mathbf{e}_i, \sqrt{n}\mathbf{e}_i \rangle = \text{tr}(\mathbf{A}^p) < \infty$.

1.1. Trace Estimation as an Algorithmic Primitive. The above progress for Schatten norm estimation made progress by representing our quantity of interest as a trace, and then applying randomized trace estimators to give us our quantity of interest. This idea has become a canonical tool for computing traces of functions of matrices represented by power series. Indeed, if we want to compute some function $\text{tr}(f(\mathbf{A}))$ where

$$f(x) = \sum_{n=0}^{\infty} a_n x^n$$

is analytic, then

$$\text{tr}(f(\mathbf{A})) = \sum_{n=0}^{\infty} a_n \text{tr}(\mathbf{A}^n) \approx \sum_{n=0}^N a_n \text{tr}(\mathbf{A}^n)$$

for some small N depending on f and \mathbf{A} . From here we can approximate further each $\text{tr}(\mathbf{A}^n)$ by a trace estimator to compute an approximation to the original quantity of interest $\text{tr}(f(\mathbf{A}))$. This exact idea has been used to create fast input-sparsity time algorithms for computing

$$\log \det(\mathbf{I} + \mathbf{A}) = \text{tr} \log(\mathbf{I} + \mathbf{A}) = \sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n} \text{tr}(\mathbf{A}^n)$$

in the recent work of Boutsidis et al. (2017). Trace estimation has been used as a primitive for creating better algorithms for counting the number of triangles in large graphs, (Avron, 2010) counting the number of eigenvalues of a matrix inside an interval, (Di Napoli et al., 2016) and even estimating the entire spectrum of matrices. (Adams et al., 2018) In the coming sections, we will go into detail on this primitive, why it works, and why we can't hope to do better by other means.

1.2. Computing Traces. Trace computation is easy if we have full access to a matrix \mathbf{A} . Indeed,

$$\text{tr} \mathbf{A} = \sum_{i=1}^n \mathbf{A}_{ii}$$

can be computed in $n - 1$ flops. But what if we don't have direct access to \mathbf{A} , or in particular if computing \mathbf{A}_{ii} takes about as long as computing $\mathbf{x}^* \mathbf{A} \mathbf{x}$ for any input vector \mathbf{x} ? The above motivation suggests a natural algorithmic process for computing the trace of matrices of this type. We start with some random vector \mathbf{r} distributed so that

$$\mathbf{E}\langle \mathbf{A} \mathbf{r}, \mathbf{r} \rangle = \text{tr } \mathbf{A}$$

for all positive matrices $\mathbf{A} \in \mathbf{M}_n$. From then, we sample m identical and independent copies of \mathbf{r} : \mathbf{r}_i for $i = 1, 2, \dots, m$. Finally, we return

$$T_m = \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A} \mathbf{r}_i, \mathbf{r}_i \rangle \approx \text{tr } \mathbf{A}$$

For which random vectors \mathbf{r} does this informal algorithm work? Well, it must be that $\mathbf{E}\langle \mathbf{A} \mathbf{r}, \mathbf{r} \rangle = \text{tr } \mathbf{A}$, so by the previous section we know $\mathbf{r} \sim \text{Uniform}\{\sqrt{n} \mathbf{e}_i\}$ works. It turns out that a much larger class of vectors works, however.

THEOREM 1.2. $\mathbf{E}\langle \mathbf{A} \mathbf{r}, \mathbf{r} \rangle = \text{tr } \mathbf{A}$ for all positive definite $\mathbf{A} \in \mathbf{M}_n$ if and only if $\mathbf{E} \mathbf{r} \mathbf{r}^* = \mathbf{E} \mathbf{r} \otimes \mathbf{r} = \mathbf{I}$.

Proof. By linearity and the cyclic properties of the trace,

$$\mathbf{E}\langle \mathbf{A} \mathbf{r}, \mathbf{r} \rangle = \mathbf{E} \text{tr}(\mathbf{r}^* \mathbf{A} \mathbf{r}) = \mathbf{E} \text{tr}(\mathbf{A} \mathbf{r} \mathbf{r}^*) = \text{tr}(\mathbf{A} (\mathbf{E} \mathbf{r} \mathbf{r}^*))$$

Let $\Sigma = \mathbf{E} \mathbf{r} \mathbf{r}^*$. If $\Sigma = \mathbf{I}$ then clearly $\mathbf{E}\langle \mathbf{A} \mathbf{r}, \mathbf{r} \rangle = \text{tr}(\mathbf{A} \Sigma) = \text{tr}(\mathbf{A})$ for all $\mathbf{A} \in \mathbf{M}_n$. On the other hand, if

$$\text{tr}(\mathbf{A} \Sigma) = \text{tr}(\mathbf{A})$$

for all positive *semi*-definite $\mathbf{A} \in \mathbf{M}_n$ then in particular

$$\Sigma_{ii} = \mathbf{e}_i^* \Sigma \mathbf{e}_i = \text{tr}(\mathbf{e}_i \mathbf{e}_i^* \Sigma) = \text{tr}(\mathbf{e}_i \mathbf{e}_i^*) = 1.$$

But then

$$2 + 2\Sigma_{ij} = \Sigma_{ii} + \Sigma_{jj} + 2\Sigma_{ij} = \text{tr}((\mathbf{e}_i + \mathbf{e}_j)(\mathbf{e}_i + \mathbf{e}_j)^* \Sigma) = \text{tr}(\mathbf{e}_i + \mathbf{e}_j)(\mathbf{e}_i + \mathbf{e}_j)^* = 2$$

by the same logic so $\Sigma_{ij} = 0$ for $i \neq j$. In sum $\Sigma = \mathbf{I}$ if $\mathbf{E}\langle \mathbf{A} \mathbf{r}, \mathbf{r} \rangle = \text{tr } \mathbf{A}$ for all positive *semi*-definite \mathbf{A} . But by constructing positive semi-definite matrices as a limit of strictly positive definite matrices (e.g. by adding $\epsilon \mathbf{I}$) this also holds if $\mathbf{E}\langle \mathbf{A} \mathbf{r}, \mathbf{r} \rangle = \text{tr } \mathbf{A}$ for only the strictly positive definite matrices. \square

This implies that standard normal vectors $\mathbf{r} \sim \mathcal{N}(0, \mathbf{I})$ as well as Rademacher vectors $\mathbf{r} \sim \text{Uniform}\{-1, 1\}^n$ both work for constructing trace estimators.

2. Upper Bounds on Performance of the Hutchinson Estimator. The above work suggests that if $\mathbf{r}_i \sim \text{Uniform}\{-1, 1\}^n$ are identical and independent Rademacher random vectors then

$$H_m = \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A} \mathbf{r}_i, \mathbf{r}_i \rangle \rightarrow \text{tr } \mathbf{A}$$

almost surely for any fixed matrix \mathbf{A} as $m \rightarrow \infty$. We can get a more quantitative bound on the performance if \mathbf{A} is known to be positive semi-definite, though.

THEOREM 2.1 (Roosta-Khorasani and Ascher Thm 1). *If $\mathbf{A} \in \mathbf{M}_n$ is positive semi-definite, then*

$$\mathbf{P}(|H_m - \text{tr } \mathbf{A}| > \epsilon \text{tr } \mathbf{A}) < 2e^{-\frac{m}{2} \left(\frac{\epsilon^2}{2} - \frac{\epsilon^3}{3} \right)}.$$

In other words, if $m \geq \frac{12 \log(\frac{2}{\delta})}{\epsilon^2(3-2\epsilon)}$ then the relative error $|H_m - \text{tr } \mathbf{A}| \leq \epsilon \text{tr } \mathbf{A}$ with probability at least $1 - \delta$. Moreover, if $\epsilon < \frac{1}{2}$, then $m \geq \frac{6}{\epsilon^2} \log \frac{2}{\delta}$ implies the relative error $|H_m - \text{tr } \mathbf{A}| \leq \epsilon \text{tr } \mathbf{A}$ with probability at least $1 - \delta$.

Note crucially that the necessary number of queries m needed to construct an ϵ -relative error approximation to $\text{tr } \mathbf{A}$ with high probability doesn't depend on the dimensions of \mathbf{A} , which is why this approach lends itself so well to reducing dimension dependence of computational algorithms.

The proof of Theorem 2.1 is surprisingly tedious, and crucially relies on a sequence of results from Achlioptas (2001) which bound the moment generating function of the square of a projection of a Rademacher by reducing to a Gaussian case. To avoid these troubles, we will prove the identical result for an analogous estimator

$$G_n = \frac{1}{n} \sum_{i=1}^n \langle \mathbf{A} \mathbf{g}_i, \mathbf{g}_i \rangle$$

where $\mathbf{g}_i \sim \mathcal{N}(0, I)$ identically and independently. Actually sampling Gaussian vectors and computing G_n takes a bit more work than G_m , but from a statistical perspective these work just as well. Note that the result presented here is stronger than that presented as (Roosta-Khorasani and Ascher, 2015, Thm 3), and brings the current best known sample complexity for trace estimators of Hutchinson type to a tie between the Gaussian and Rademacher cases.

THEOREM 2.2. *If $\mathbf{A} \in \mathbb{M}_n$ is positive semi-definite, then*

$$\mathbf{P}(|G_m - \text{tr } \mathbf{A}| > \epsilon \text{tr } \mathbf{A}) < 2e^{-\frac{m}{2} \left(\frac{\epsilon^2}{2} - \frac{\epsilon^3}{3} \right)}.$$

In other words, if $m \geq \frac{12 \log(\frac{2}{\delta})}{\epsilon^2(3-2\epsilon)}$ then the relative error $|G_m - \text{tr } \mathbf{A}| \leq \epsilon \text{tr } \mathbf{A}$ with probability at least $1 - \delta$. Moreover, if $\epsilon < \frac{1}{2}$, then $m \geq \frac{6}{\epsilon^2} \log \frac{2}{\delta}$ implies the relative error $|G_m - \text{tr } \mathbf{A}| \leq \epsilon \text{tr } \mathbf{A}$ with probability at least $1 - \delta$.

Proof. Diagonalize $\mathbf{A} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^*$ where \mathbf{U} is unitary and $\mathbf{\Lambda}$ is diagonal. Let $\mathbf{z}_i = \mathbf{U}^* \mathbf{g}_i$ so that

$$\begin{aligned} \mathbf{P}(G_m > (1 + \epsilon) \text{tr } \mathbf{A}) &= \mathbf{P}\left(\sum_{i=1}^m \mathbf{g}_i^* \mathbf{A} \mathbf{g}_i > (1 + \epsilon) \text{tr } \mathbf{A}\right) \\ &= \mathbf{P}\left(\sum_{i=1}^m \mathbf{z}_i^* \mathbf{\Lambda} \mathbf{z}_i > m(1 + \epsilon) \text{tr } \mathbf{A}\right) \\ &= \mathbf{P}\left(\sum_{i=1}^m \sum_{j=1}^n \lambda_j \mathbf{z}_{ij}^2 > m(1 + \epsilon) \text{tr } \mathbf{A}\right) \\ &= \mathbf{P}\left(\sum_{j=1}^n \frac{\lambda_j}{\text{tr } \mathbf{A}} \sum_{i=1}^m \mathbf{z}_{ij}^2 > m(1 + \epsilon)\right) \\ &= \mathbf{P}\left(\exp\left(t \sum_{j=1}^n \frac{\lambda_j}{\text{tr } \mathbf{A}} \sum_{i=1}^m \mathbf{z}_{ij}^2\right) > e^{(1+\epsilon)mt}\right) \\ &\leq e^{-(1+\epsilon)mt} \mathbf{E} \exp\left(\sum_{j=1}^n \frac{\lambda_j}{\text{tr } \mathbf{A}} \sum_{i=1}^m t \mathbf{z}_{ij}^2\right) \end{aligned}$$

for $t > 0$ by Markov's inequality. (Wasserman, 2013, Thm 4.1) Now by convexity of the exponential and linearity of expectation

$$\mathbf{E} \exp\left(\sum_{j=1}^n \frac{\lambda_j}{\text{tr } \mathbf{A}} \sum_{i=1}^m t \mathbf{z}_{ij}^2\right) \leq \sum_{j=1}^n \frac{\lambda_j}{\text{tr } \mathbf{A}} \mathbf{E} \exp\left(\sum_{i=1}^m t \mathbf{z}_{ij}^2\right) = \sum_{j=1}^n \frac{\lambda_j}{\text{tr } \mathbf{A}} \prod_{i=1}^m \mathbf{E} e^{t \mathbf{z}_{ij}^2}$$

since $\lambda_i \geq 0$ and $\sum_{i=1}^n \lambda_i = \text{tr } \mathbf{A}$, and for a fixed j , \mathbf{z}_{ij} are independent and identically distributed. Now, to bound the moment generating function, observe that \mathbf{z}_i are just independent standard normal vectors by

the rotation invariance of the Gaussian and unitary nature of \mathbf{U} . This implies that $z_{ij} \sim \mathcal{N}(0, 1)$ and hence z_{ij}^2 is just a χ^2 random variable with one degree of freedom. We then know (see [Wasserman \(2013\)](#)) that

$$\mathbf{E} e^{t z_{ij}^2} = \frac{1}{\sqrt{1-2t}}$$

so long as $0 < t < 1/2$. This implies that

$$\mathbf{E} \exp\left(\sum_{j=1}^n \frac{\lambda_j}{\text{tr } \mathbf{A}} \sum_{i=1}^m t z_{ij}^2\right) \leq \left(\frac{1}{\sqrt{1-2t}}\right)^m$$

for these $0 < t < 1/2$ and hence taking $t = \frac{1}{2} \frac{\epsilon}{1+\epsilon} < 1/2$ we have

$$\mathbf{P}(G_m > (1+\epsilon) \text{tr } \mathbf{A}) \leq \left(\frac{1}{\sqrt{1-2t}}\right)^m e^{-(1+\epsilon)mt} = \left((1+\epsilon)e^{-\epsilon}\right)^{m/2} < e^{-\frac{m}{2}(\frac{\epsilon^2}{2}-\frac{\epsilon^3}{3})},$$

relying on the fact that $(1+\epsilon)e^{-\epsilon} < e^{\frac{\epsilon^3}{3}-\frac{\epsilon^2}{2}}$ for all $\epsilon > 0$, verifiable by simple calculus. Essentially the same argument for the lower tail tells us via a union bound that

$$\mathbf{P}(|G_m - \text{tr } \mathbf{A}| > \epsilon \text{tr } \mathbf{A}) < 2e^{-\frac{m}{2}(\frac{\epsilon^2}{2}-\frac{\epsilon^3}{3})}. \quad \square$$

3. Lower Bounds. We've illustrated extremely simple randomized algorithms for estimating the trace of a large matrix with implicit access. Moreover, we've shown upper bounds on these algorithms that guarantee fast performance with only $\frac{6}{\epsilon^2} \log \frac{2}{\delta}$ queries to \mathbf{A} . Nevertheless, questions still remain. Can we estimate traces with the same number of queries to \mathbf{A} without using any randomness? Perhaps better estimators for the trace exist that only need $\frac{6}{\epsilon} \log \frac{2}{\delta}$ or $\frac{6}{\epsilon^2} \log \log \frac{2}{\delta}$ queries to \mathbf{A} . This section will illustrate why these conjectures are both false.

3.1. Deterministic Algorithms Can't Estimate Traces. Let's assume we have an oracle $\mathcal{O} : \mathbf{x} \mapsto \mathbf{A}\mathbf{x}$ and we'd like to deterministically estimate the trace of \mathbf{A} by sequentially submitting queries m to \mathcal{O} , perhaps performing an uncomputable amount of work in between or rounds or after the querying process. The following theorem essentially says that the only real way to get an estimator T_n from this process with relative error

$$|T_m - \text{tr } \mathbf{A}| < \epsilon \text{tr } \mathbf{A}$$

is to just compute

$$T_n = \sum_{i=1}^n \langle \mathcal{O}(e_i), e_i \rangle = \sum_{i=1}^n \mathbf{A}_{ii} = \text{tr } \mathbf{A}$$

exactly.

THEOREM 3.1. *Any deterministic trace estimator T_m satisfying*

$$|T_m - \text{tr } \mathbf{A}| \leq \epsilon \text{tr } \mathbf{A}$$

for all positive semi-definite matrices $\mathbf{A} \in \mathbf{M}_n$ and some $0 < \epsilon < 1$ requires $m \geq n$.

Proof. We will follow via a resisting oracle. In particular, let's suppose that \mathcal{O} returns $\mathcal{O}(\mathbf{x}_i) = 0$ for the first $n-1$ queries $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n-1}$. At this point, both the zero matrix $\mathbf{A}_1 = 0$ and the orthogonal projection onto the orthogonal complement of the span of $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n-1}\}$ \mathbf{A}_2 would return the 0 for these queries. If we output an estimator $T_m = T_{n-1}$ which satisfied

$$|T_m - \text{tr } \mathbf{A}| \leq \epsilon \text{tr } \mathbf{A}$$

for all positive semi-definite matrices $\mathbf{A} \in \mathbf{M}_n$, then since \mathbf{A}_1 is positive semi-definite this series of responses by \mathcal{O} would ensure

$$|T_m| = |T_m - \text{tr } \mathbf{A}_1| \leq \epsilon \text{tr } \mathbf{A}_1 = 0,$$

so $T_m = 0$. On the other hand, the matrix \mathbf{A}_2 would give the same sequence of responses by \mathcal{O} and so by determinism

$$\text{tr } \mathbf{A}_2 = |\text{tr } \mathbf{A}_2| = |T_m - \text{tr } \mathbf{A}_2| \leq \epsilon \text{tr } \mathbf{A}_2 < \text{tr } \mathbf{A}_2,$$

a contradiction. It follows that at least n queries to \mathcal{O} are needed to guarantee this uniform error bound deterministically. \square

3.2. Simple Lower Bounds for Simple Randomized Trace Estimators. The next section will reference much more general lower bounds for randomized trace estimation. Those proofs are long and tedious. This section will present simple lower bounds due to the author and Weiqing Gu that guarantee the tightness of [Roosta-Khorasani and Ascher \(2015\)](#) without a ton of heavy machinery, at the expense of the full generality of the bounds of [Wimmer et al. \(2014\)](#).

Recall that the Hutchinson estimator of the trace of a matrix $\mathbf{A} \in \mathbb{M}_n$ is

$$H_m = \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A} \mathbf{r}_i, \mathbf{r}_i \rangle$$

where $\mathbf{r}_i \sim \text{Uniform}\{-1, 1\}^n$ independently and identically. For this section, we will call H_n a Hutchinson-type estimator for the trace of a matrix \mathbf{A} if

$$H_m = \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A} \mathbf{r}_i, \mathbf{r}_i \rangle$$

where $\mathbf{E} \mathbf{r}_i \mathbf{r}_i^* = \mathbf{I}$ and \mathbf{r} is real almost surely. Recall from Theorem 1.2 that $H_n \rightarrow \text{tr } \mathbf{A}$ almost surely. While the original Hutchinson estimator returns the trace exactly in one query for diagonal matrices, it happens that every Hutchinson-type estimator has some matrix \mathbf{A} for which $\text{Var}(H_m) > 0$.

LEMMA 3.2. *If the real random vector \mathbf{r} satisfies $\mathbf{E} \mathbf{r} \mathbf{r}^* = \mathbf{I}$, then $\text{Var}(\langle \mathbf{A} \mathbf{r}, \mathbf{r} \rangle) > 0$ for some positive semi-definite $\mathbf{A} \in \mathbb{M}_n$.*

Proof. In light of Theorem 1.2, suppose to the contrary that $\langle \mathbf{A} \mathbf{r}, \mathbf{r} \rangle = \text{tr } \mathbf{A}$ almost surely for every positive semi-definite $\mathbf{A} \in \mathbb{M}_n$. In particular,

$$\mathbf{r}_i^2 = \langle \mathbf{e}_i \mathbf{e}_i^* \mathbf{r}, \mathbf{r} \rangle = \text{tr } \mathbf{e}_i \mathbf{e}_i^* = 1$$

almost surely for all $i, j = 1, 2, \dots, n$, so $\mathbf{r} \in \{-1, 1\}^d$ almost surely. On the other hand, if $\mathbf{r}_i \neq \mathbf{r}_j$ with positive probability for some $i \neq j$ we would have

$$0 = (\mathbf{r}_i + \mathbf{r}_j)^2 = \langle (\mathbf{e}_i + \mathbf{e}_j)(\mathbf{e}_i + \mathbf{e}_j)^* \mathbf{r}, \mathbf{r} \rangle = \text{tr}(\mathbf{e}_i + \mathbf{e}_j)(\mathbf{e}_i + \mathbf{e}_j)^* = 2$$

with positive probability, a contradiction, so almost surely every component of \mathbf{r} is the same. But then $0 = \mathbf{E} \mathbf{r}_i \mathbf{r}_j = \mathbf{E} \mathbf{r}_i^2 = 1$, contradicting our assumption that $\mathbf{E} \mathbf{r} \mathbf{r}^* = \mathbf{I}$. \square

The following theorem relies weakly on the above lemma to say that any Hutchinson-type estimator needs $\Omega(\frac{1}{\epsilon^2} \log \frac{1}{\delta})$ queries to \mathbf{A} in order to achieve an ϵ -relative error approximation to $\text{tr } \mathbf{A}$ with probability at least $1 - \delta$. Recall that W is the Lambert-W function, satisfying $W(x) = \log x - \log \log x + o(1) = \Theta(\log x)$ as $x \rightarrow \infty$. ([Hoorfar and Hassani, 2007](#))

THEOREM 3.3. *Fix $0 < \delta < \frac{1}{10}$. For every Hutchinson-type estimator H_n there exists a positive semi-definite matrix $\mathbf{A} \in \mathbb{M}_n$ and an $\epsilon_0 > 0$ so that*

$$\mathbf{P}(|H_n - \text{tr } \mathbf{A}| > \epsilon \text{tr } \mathbf{A}) > \delta$$

whenever $0 < \epsilon < \epsilon_0$ and $n = \lfloor \frac{\text{Var}(\langle \mathbf{A} \mathbf{x}, \mathbf{x} \rangle)}{\text{tr}(\mathbf{A})^2} \frac{1}{\epsilon^2} W(\frac{2/\pi}{\delta^2}) \rfloor = \Theta(\frac{1}{\epsilon^2} \log \frac{1}{\delta})$.

Proof. Pick \mathbf{A} so that $\text{Var}(\langle \mathbf{A} \mathbf{r}, \mathbf{r} \rangle) > 0$ by the lemma. If $Z \sim \mathcal{N}(0, 1)$ is a standard normal random variable then

$$\mathbf{P}(|Z| > t) = 2\mathbf{P}(Z > t) > \sqrt{\frac{2}{\pi}} \frac{t}{t^2 + 1} e^{-t^2} \geq \frac{1}{\sqrt{2\pi}} \frac{1}{t} e^{-t^2}$$

when $t \geq 1$. (Cook, 2009) Setting the right hand side to 2δ and solving gives

$$\mathbf{P}(|Z| > 2^{-1/2} \sqrt{W(\pi^{-1}\delta^{-2})}) > 2\delta$$

whenever $2^{-1/2} \sqrt{W(\pi^{-1}\delta^{-2})} \geq 1$ or $0 < \delta \leq (e\sqrt{2\pi})^{-1}$. This is satisfied when $0 < \delta < \frac{1}{10}$. If we fix $n = \lfloor \frac{\mathbf{Var}(\langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle)}{\text{tr}(\mathbf{A})^2} \frac{1}{\epsilon^2} W(\frac{2/\pi}{\delta^2}) \rfloor$ and write $\sigma^2 = \mathbf{Var}(\langle \mathbf{A}\mathbf{r}, \mathbf{r} \rangle) > 0$,

$$\mathbf{P}(|H_n - \text{tr } \mathbf{A}| > \epsilon \text{tr } \mathbf{A}) = \mathbf{P}\left(\frac{\sqrt{n}}{\sigma} |H_n - \text{tr } \mathbf{A}| > \frac{\sqrt{n}}{\sigma} \epsilon \text{tr } \mathbf{A}\right) \geq \mathbf{P}\left(\frac{\sqrt{n}}{\sigma} |H_n - \text{tr } \mathbf{A}| > \frac{1}{\sqrt{2}} \sqrt{W(\pi^{-1}\delta^{-2})}\right).$$

Conveniently, the latter expression converges to $\mathbf{P}(|Z| \geq 2^{-1/2} \sqrt{W(\pi^{-1}\delta^{-2})})$ as $n \rightarrow \infty$, (Wasserman, 2013, Thm 5.10) and so

$$\lim_{\epsilon \rightarrow 0} \mathbf{P}(|H_n - \text{tr } \mathbf{A}| > \epsilon \text{tr } \mathbf{A}) > 2\delta.$$

This implies the existence of an $\epsilon_0 > 0$ so that $0 < \epsilon < \epsilon_0$ ensures $\mathbf{P}(|H_n - \text{tr } \mathbf{A}| > \epsilon \text{tr } \mathbf{A}) > \delta$ under our relation defining n . \square

3.3. Generic Lower Bounds. Somewhat recent work by Wimmer et al. (2014) has shown that the upper bounds given by Roosta-Khorasani and Ascher (2015) are tight: no algorithm that can sequentially submit vectors \mathbf{x} to an oracle $\mathcal{O} : \mathbf{x} \mapsto \langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle$ and possibly estimate the trace of \mathbf{A} unless $\Omega(\frac{1}{\epsilon^2} \log \frac{1}{\delta})$ queries to \mathbf{A} are made. In particular, Wimmer, Wu, and Zhang prove the following sequence of Theorems.

THEOREM 3.4 (Wimmer et al. Thm 1). *If we consider estimators for $\text{tr } \mathbf{A}$ that pre-decide a distribution over queries $(\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_m)$ as well as weights (w_1, w_2, \dots, w_m) and output*

$$T_m = \sum_{i=1}^m w_i \mathcal{O}(\mathbf{r}_i),$$

the minimum variance estimator for which $\mathbf{E}T_m = \text{tr } \mathbf{A}$ uniformly on \mathbf{M}_n is achieved by sampling $\{\mathbf{r}_i\}$ as a collection of m orthogonal unit vectors and outputting

$$T_m^* = \frac{n}{m} \sum_{i=1}^m \mathcal{O}(\mathbf{r}_i)$$

THEOREM 3.5 (Wimmer et al. Thm 2). *Any possibly nonlinear or adaptive estimator for the trace of a matrix \mathbf{A} that sequentially submits random queries \mathbf{r}_i to \mathcal{O} after seeing the previous $i - 1$ queries needs $\Omega(1/\epsilon)$ queries to achieve ϵ mean squared error uniformly across all matrices with Frobenius norm 1.*

THEOREM 3.6 (Wimmer et al. Thm 3). *Any possibly nonlinear or adaptive estimator for the trace of a matrix \mathbf{A} that sequentially submits random queries \mathbf{r}_i to \mathcal{O} after seeing the previous $i - 1$ queries needs $\Omega(\frac{1}{\epsilon^2} \log \frac{1}{\delta})$ queries to output an estimator T_m that satisfies*

$$\mathbf{P}(|T - \text{tr } \mathbf{A}| > \epsilon \text{tr } \mathbf{A}) \leq \delta$$

for any rank-one positive definite matrix \mathbf{A} .

The proofs of these theorems are long and beyond the scope of this lecture. See Wimmer et al. (2014) for more details.

As we saw in the case of Schatten norm estimation, the common algorithmic setting we have isn't an oracle mapping $\mathbf{x} \mapsto \langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle$, but rather an oracle which can tell us $\mathbf{A}\mathbf{x}$ for any given \mathbf{x} . Thus, while the lower bounds of Wimmer et al. (2014) are interesting, correct, and relevant, the following open problem really gets to the heart of the fundamental problem of trace estimation as it is actually used in practice.

OPEN PROBLEM 3.7. *If we have access to an oracle $\mathcal{O} : \mathbf{x} \mapsto \mathbf{A}\mathbf{x}$, is it true that $\Omega(\frac{1}{\epsilon^2} \log \frac{1}{\delta})$ queries to \mathcal{O} are needed to return a (possibly adaptive and nonlinear) estimator T for the trace of \mathbf{A} with*

$$\mathbf{P}(|T - \text{tr } \mathbf{A}| > \epsilon \text{tr } \mathbf{A}) \leq \delta$$

for every positive semi-definite matrix \mathbf{A} .

References.

- Dimitris Achlioptas. Database-friendly random projections. In *Proceedings of the twentieth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 274–281. ACM, 2001.
- Ryan P Adams, Jeffrey Pennington, Matthew J Johnson, Jamie Smith, Yaniv Ovadia, Brian Patton, and James Saunderson. Estimating the spectral density of large implicit matrices. *arXiv preprint arXiv:1802.03451*, 2018.
- Haim Avron. Counting triangles in large graphs using randomized matrix trace estimation. In *Workshop on Large-scale Data Mining: Theory and Applications*, volume 10, pages 10–9, 2010.
- Christos Boutsidis, Petros Drineas, Prabhanjan Kambadur, Eugenia-Maria Kontopoulou, and Anastasios Zouzias. A randomized algorithm for approximating the log determinant of a symmetric positive definite matrix. *Linear Algebra and its Applications*, 533:95–117, 2017.
- John D Cook. Upper and lower bounds for the normal distribution function, 2009.
- Edoardo Di Napoli, Eric Polizzi, and Yousef Saad. Efficient estimation of eigenvalue counts in an interval. *Numerical Linear Algebra with Applications*, 23(4):674–692, 2016.
- Gene H Golub and Charles F Van Loan. *Matrix computations*, volume 3. JHU Press, 2012.
- Abdolhossein Hoorfar and Mehdi Hassani. Approximation of the lambert w function and hyperpower function. *Research report collection*, 10(2), 2007.
- Farbod Roosta-Khorasani and Uri Ascher. Improved bounds on sample size for implicit matrix trace estimators. *Foundations of Computational Mathematics*, 15(5):1187–1212, 2015.
- Larry Wasserman. *All of statistics: a concise course in statistical inference*. Springer Science & Business Media, 2013.
- Karl Wimmer, Yi Wu, and Peng Zhang. Optimal query complexity for estimating the trace of a matrix. In *International Colloquium on Automata, Languages, and Programming*, pages 1051–1062. Springer, 2014.
- David P Woodruff et al. Sketching as a tool for numerical linear algebra. *Foundations and Trends® in Theoretical Computer Science*, 10(1–2):1–157, 2014.