

Hutch++: Optimal Stochastic Trace Estimation

Raphael A. Meyer¹ Cameron Musco² Christopher Musco¹ David P. Woodruff³

¹New York University

²University of Massachusetts Amherst

³Carnegie Mellon University

Introduction and Overview

Trace Estimation

Trace Estimation

- Goal: Estimate trace of $n \times n$ matrix A :

$$\text{tr}(A) = \sum_{i=1}^n A_{ii} = \sum_{i=1}^n \lambda_i$$

- In Downstream Applications, A is not stored in memory.
- Instead, B is in memory and $A = f(B)$:

No. Triangles	Estrada Index	Log-Determinant
$\text{tr}(\frac{1}{6}B^3)$	$\text{tr}(e^B)$	$\text{tr}(\ln(B))$

- Computing $A = \frac{1}{6}B^3$ takes $O(n^3)$ time
- Computing $Ax = \frac{1}{6}B(Bx)$ takes $O(n^2)$ time
- If $A = f(B)$, then we can often compute Ax quickly

Matrix-Vector Oracle Model

Idea: Matrix-Vector Product as a Computational Primitive

- Given access to a $n \times n$ matrix A only through a Matrix-Vector Multiplication Oracle:

$$x \xrightarrow{\text{input}} \text{ORACLE} \xrightarrow{\text{output}} Ax$$

- e.g. Krylov Methods, Sketching, Streaming, ...

Implicit Matrix Trace Estimation:

Estimate $\text{tr}(A)$ with as few Matrix-Vector products Ax_1, \dots, Ax_m as possible:

$$|\tilde{\text{tr}}(A) - \text{tr}(A)| \leq \varepsilon \text{tr}(A)$$

Our Contributions

For PSD A , we show $\Theta(\frac{1}{\varepsilon})$ products are necessary and sufficient. Prior work used $O(\frac{1}{\varepsilon^2})$ products.

Hutch++ Algorithm:

- Input: Number of matrix-vector queries m , matrix A
- 1. Sample $S \in \mathbb{R}^{d \times \frac{m}{3}}$ and $G \in \mathbb{R}^{d \times \frac{m}{3}}$ with i.i.d. $\{+1, -1\}$ entries
- 2. Compute $Q = \text{qr}(AS)$
- 3. Return $\text{tr}(Q^T A Q) + \frac{3}{m} \text{tr}(G^T (I - Q Q^T) A (I - Q Q^T) G)$

Main Theorem and Context

Prior Work: Hutchinson's Estimator

Hutchinson's Estimator:

- If $x \sim \mathcal{N}(0, I)$, then

$$\mathbb{E}[x^T A x] = \text{tr}(A) \quad \text{Var}[x^T A x] = \|A\|_F^2$$

- Hutchinson's Estimator: $H_\ell(A) := \frac{1}{\ell} \sum_{i=1}^{\ell} x_i^T A x_i$

$$\mathbb{E}[H_\ell(A)] = \text{tr}(A) \quad \text{Var}[H_\ell(A)] = \frac{1}{\ell} \|A\|_F^2$$

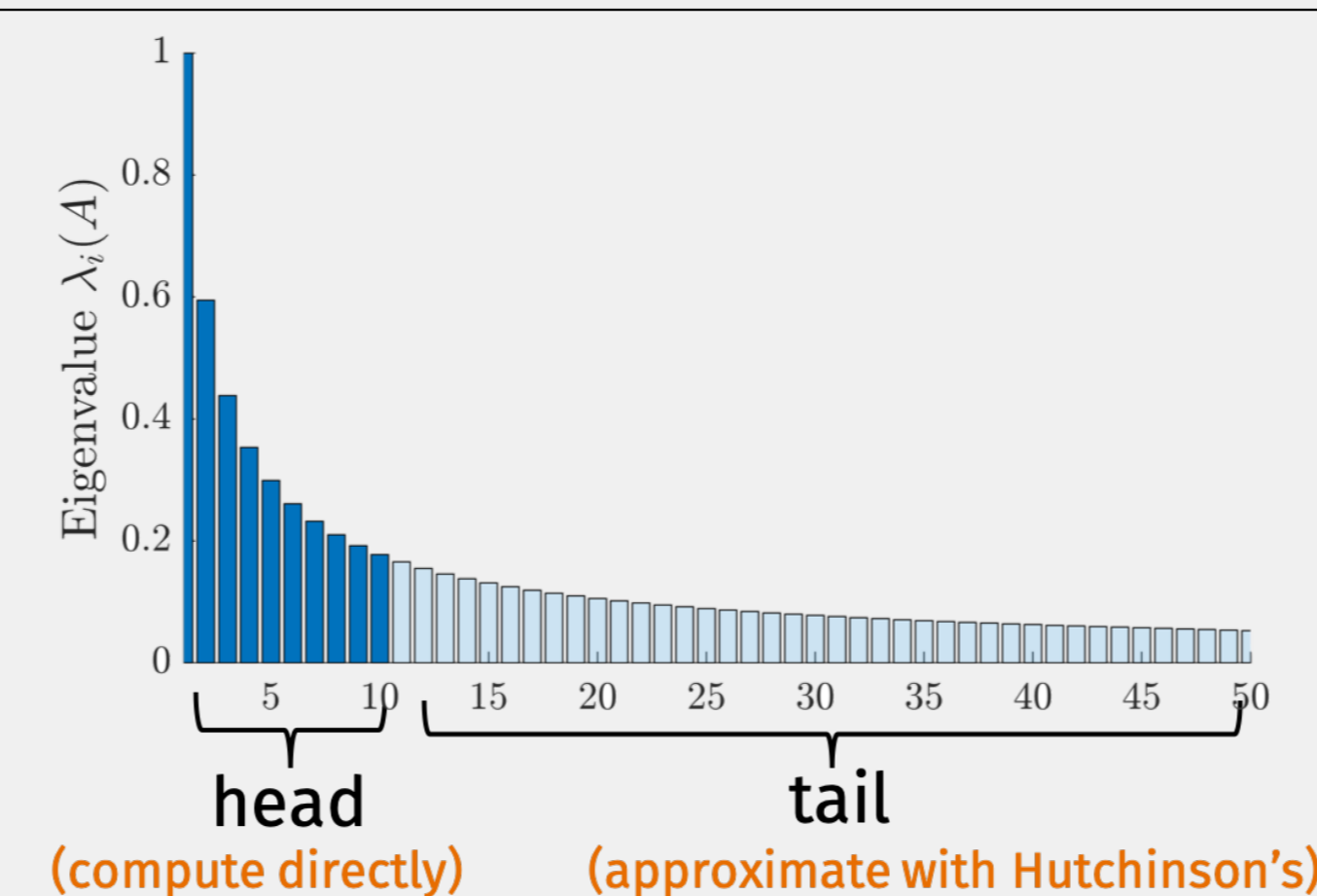
- For $A \succeq 0$, we have $\|A\|_F \leq \text{tr}(A)$, so that

$$\begin{aligned} |H_\ell(A) - \text{tr}(A)| &\leq O\left(\frac{1}{\sqrt{\ell}}\right) \|A\|_F && \text{(Chebyshev Ineq.)} \\ &\leq O\left(\frac{1}{\sqrt{\ell}}\right) \text{tr}(A) && (\|A\|_F \leq \text{tr}(A)) \\ &= \varepsilon \text{tr}(A) && (\ell = O(\frac{1}{\varepsilon^2})) \end{aligned}$$

- This analysis is only tight if $\|A\|_F \approx \text{tr}(A)$!
- $\|A\|_F \approx \text{tr}(A)$ only if A 's eigenvalues are nearly sparse!

If A is hard for the Hutchinson estimator, then the sum of the top k eigenvalues represents most of $\text{tr}(A)$.

Core Intuition



- Find a good rank- k approximation \tilde{A}_k
- Notice that $\text{tr}(A) = \text{tr}(\tilde{A}_k) + \text{tr}(A - \tilde{A}_k)$
- Compute $\text{tr}(\tilde{A}_k)$ exactly
- Return $\text{Hutch++}(A) := \text{tr}(\tilde{A}_k) + H_\ell(A - \tilde{A}_k)$

Theorem: If $k = \ell = O(\frac{1}{\varepsilon})$, then $|\text{Hutch++}(A) - \text{tr}(A)| \leq \varepsilon \text{tr}(A)$

Fundamental Rate: $|\text{Hutch++}(A) - \text{tr}(A)| \leq O(\frac{1}{\ell}) \|A - A_k\|_F$

Conclusions

Experiments

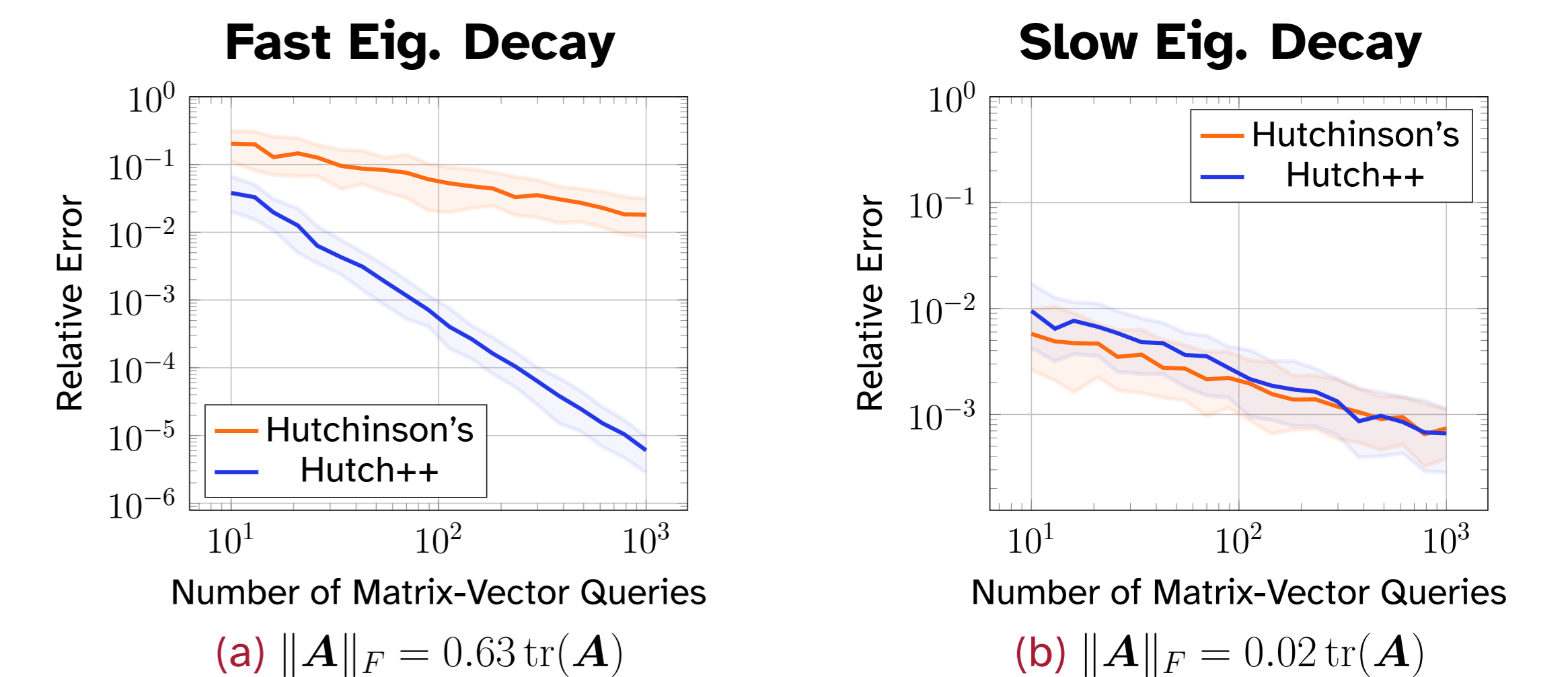


Figure 1. When $\|A\|_F \approx \text{tr}(A)$, A has quickly decaying eigenvalues, and Hutch++ is much faster than H_ℓ .

Extensions Et Cetera

Lower Bounds

- All algorithms must use $\Omega(\frac{1}{\varepsilon})$ queries.

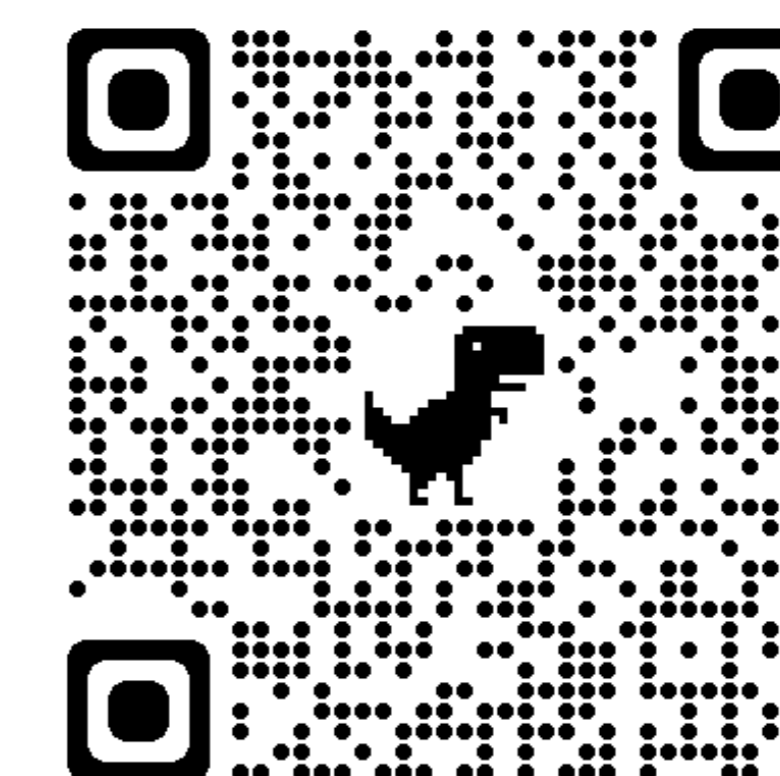
Indefinite Matrices

- We instead achieve $|\text{tr}(A) - \text{Hutch++}(A)| \leq \varepsilon \|A\|_*$.

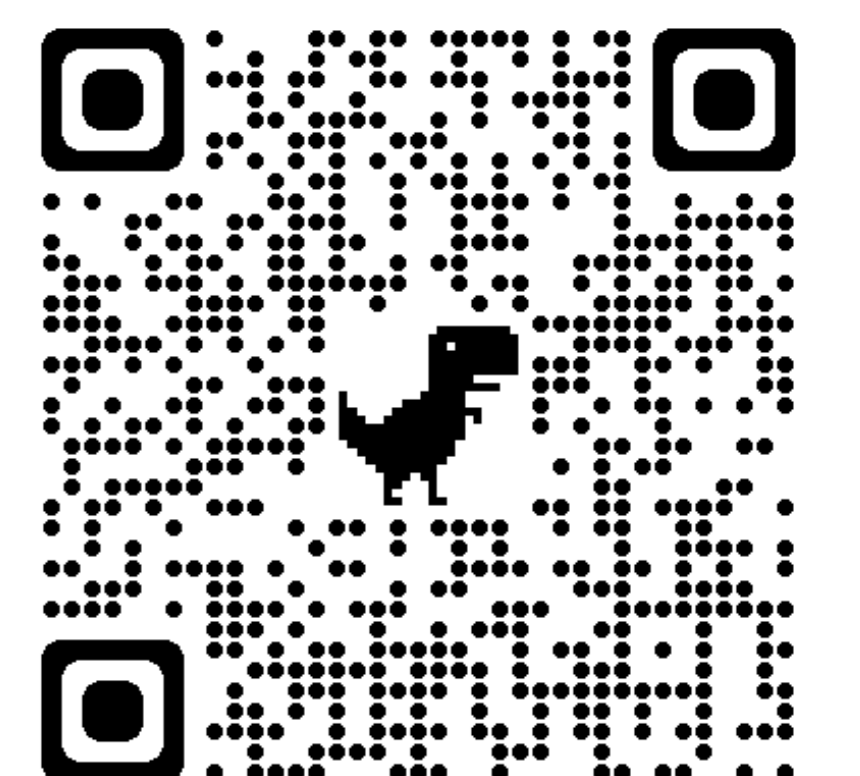
Non-Adaptive Algorithms

- We submit all queries in parallel. x_2 cannot depend on Ax_1 .
- NA-Hutch++: Non-Adaptive variant of Hutch++, still $O(\frac{1}{\varepsilon})$
- No Adaptivity Gap

QR Codes & Links for More Details



(a) Hutch++ for Undergrads



(b) Hutch++ Full Paper