

Hutch++

Optimal Stochastic Trace Estimation

Raphael A. Meyer (New York University)

With Christopher Musco (New York University), Cameron Musco (University of Massachusetts Amherst), and David P. Woodruff (Carnegie Mellon University)

Collaborators



Christopher Musco
(NYU)



Cameron Musco
(UMass. Amherst)



David P. Woodruff
(CMU)

Trace Estimation

- ⊙ Goal: Estimate trace of $n \times n$ matrix \mathbf{A} :

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

Trace Estimation

- ⊙ Goal: Estimate trace of $n \times n$ matrix \mathbf{A} :

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

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

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

Trace Estimation

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

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

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

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

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

Matrix-Vector Oracle Model

Idea: Matrix-Vector Product as a Computational Primitive

Matrix-Vector Oracle Model

Idea: Matrix-Vector Product as a Computational Primitive

- Given access to a $n \times n$ matrix \mathbf{A} only through a **Matrix-Vector Multiplication Oracle**

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

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

Matrix-Vector Oracle Model

Idea: Matrix-Vector Product as a Computational Primitive

- Given access to a $n \times n$ matrix \mathbf{A} only through a **Matrix-Vector Multiplication Oracle**

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

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

Implicit Matrix Trace Estimation: Estimate $\text{tr}(\mathbf{A})$ with as few Matrix-Vector products $\mathbf{Ax}_1, \dots, \mathbf{Ax}_m$ as possible.

$$(1 - \varepsilon) \text{tr}(\mathbf{A}) \leq \tilde{\text{tr}}(\mathbf{A}) \leq (1 + \varepsilon) \text{tr}(\mathbf{A})$$

3 Core Contributions

For PSD matrix trace estimation,

1. Hutch++ algorithm, which uses $\tilde{O}(\frac{1}{\epsilon})$ matrix-vector products.
 - Improves prior rate of $\tilde{O}(\frac{1}{\epsilon^2})$
 - Empirically works well

¹ \tilde{O} notation only hide logarithmic dependence on the failure probability.

3 Core Contributions

For PSD matrix trace estimation,

1. Hutch++ algorithm, which uses $\tilde{O}(\frac{1}{\epsilon})$ matrix-vector products.
 - Improves prior rate of $\tilde{O}(\frac{1}{\epsilon^2})$
 - Empirically works well
2. All adaptive algorithms with finite-precision oracles use $\Omega(\frac{1}{\epsilon \log(1/\epsilon)})$ queries.
3. All nonadaptive algorithms with infinite-precision oracles use $\Omega(\frac{1}{\epsilon})$ queries.

¹ \tilde{O} notation only hide logarithmic dependence on the failure probability.

The classical approach to trace estimation:

Hutchinson 1991, Girard 1987

1. Draw $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$ with i.i.d. uniform $\{+1, -1\}$ entries
2. Return $\tilde{T} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^\top \mathbf{A} \mathbf{x}_i$

Avron, Toledo 2011, Roosta, Ascher 2015

With probability $1 - \delta$,

$$|\tilde{T} - \text{tr}(\mathbf{A})| \leq \tilde{O}\left(\frac{1}{\sqrt{m}}\right) \|\mathbf{A}\|_F$$

Hutchinson Analysis

For PSD \mathbf{A} , $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$, so that:

$$|\tilde{T} - \text{tr}(\mathbf{A})| \leq O\left(\frac{1}{\sqrt{m}}\right) \|\mathbf{A}\|_F$$

Hutchinson Analysis

For PSD \mathbf{A} , $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$, so that:

$$\begin{aligned} |\tilde{T} - \text{tr}(\mathbf{A})| &\leq O\left(\frac{1}{\sqrt{m}}\right) \|\mathbf{A}\|_F \\ &\leq O\left(\frac{1}{\sqrt{m}}\right) \text{tr}(\mathbf{A}) \end{aligned}$$

Hutchinson Analysis

For PSD \mathbf{A} , $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$, so that:

$$\begin{aligned} |\tilde{T} - \text{tr}(\mathbf{A})| &\leq O\left(\frac{1}{\sqrt{m}}\right) \|\mathbf{A}\|_F \\ &\leq O\left(\frac{1}{\sqrt{m}}\right) \text{tr}(\mathbf{A}) \\ &= \varepsilon \text{tr}(\mathbf{A}) \end{aligned}$$

Hutchinson Analysis

For PSD \mathbf{A} , $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$, so that:

$$\begin{aligned} |\tilde{T} - \text{tr}(\mathbf{A})| &\leq O\left(\frac{1}{\sqrt{m}}\right) \|\mathbf{A}\|_F \\ &\leq O\left(\frac{1}{\sqrt{m}}\right) \text{tr}(\mathbf{A}) \\ &= \varepsilon \text{tr}(\mathbf{A}) \end{aligned}$$

- ⊙ When does Hutchinson's Estimator truly need $O\left(\frac{1}{\varepsilon^2}\right)$ queries?

Hutchinson Analysis

For PSD \mathbf{A} , $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$, so that:

$$\begin{aligned} |\tilde{T} - \text{tr}(\mathbf{A})| &\approx O\left(\frac{1}{\sqrt{m}}\right) \|\mathbf{A}\|_F \\ &\leq O\left(\frac{1}{\sqrt{m}}\right) \text{tr}(\mathbf{A}) \\ &= \varepsilon \text{tr}(\mathbf{A}) \end{aligned}$$

- ⊙ When does Hutchinson's Estimator truly need $O\left(\frac{1}{\varepsilon^2}\right)$ queries?

Hutchinson Analysis

For PSD \mathbf{A} , $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$, so that:

$$\begin{aligned} |\tilde{T} - \text{tr}(\mathbf{A})| &\approx O\left(\frac{1}{\sqrt{m}}\right) \|\mathbf{A}\|_F \\ &\leq O\left(\frac{1}{\sqrt{m}}\right) \text{tr}(\mathbf{A}) \\ &= \varepsilon \text{tr}(\mathbf{A}) \end{aligned}$$

- ⊙ When does Hutchinson's Estimator truly need $O\left(\frac{1}{\varepsilon^2}\right)$ queries?
- ⊙ When is the bound $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$ tight?

Hutchinson Analysis

For PSD \mathbf{A} , $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$, so that:

$$\begin{aligned} |\tilde{T} - \text{tr}(\mathbf{A})| &\approx O\left(\frac{1}{\sqrt{m}}\right) \|\mathbf{A}\|_F \\ &\leq O\left(\frac{1}{\sqrt{m}}\right) \text{tr}(\mathbf{A}) \\ &= \varepsilon \text{tr}(\mathbf{A}) \end{aligned}$$

- ⊙ When does Hutchinson's Estimator truly need $O\left(\frac{1}{\varepsilon^2}\right)$ queries?
- ⊙ When is the bound $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$ tight?
- ⊙ Let $\mathbf{v} = [\lambda_1 \ \dots \ \lambda_n]$ be the eigenvalues of PSD \mathbf{A}

Hutchinson Analysis

For PSD \mathbf{A} , $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$, so that:

$$\begin{aligned} |\tilde{T} - \text{tr}(\mathbf{A})| &\approx O\left(\frac{1}{\sqrt{m}}\right) \|\mathbf{A}\|_F \\ &\leq O\left(\frac{1}{\sqrt{m}}\right) \text{tr}(\mathbf{A}) \\ &= \varepsilon \text{tr}(\mathbf{A}) \end{aligned}$$

- ⊙ When does Hutchinson's Estimator truly need $O\left(\frac{1}{\varepsilon^2}\right)$ queries?
- ⊙ When is the bound $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$ tight?
- ⊙ Let $\mathbf{v} = [\lambda_1 \ \dots \ \lambda_n]$ be the eigenvalues of PSD \mathbf{A}
- ⊙ When is the bound $\|\mathbf{v}\|_2 \leq \|\mathbf{v}\|_1$ tight?

Hutchinson Analysis

For PSD \mathbf{A} , $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$, so that:

$$\begin{aligned} |\tilde{T} - \text{tr}(\mathbf{A})| &\approx O\left(\frac{1}{\sqrt{m}}\right) \|\mathbf{A}\|_F \\ &\leq O\left(\frac{1}{\sqrt{m}}\right) \text{tr}(\mathbf{A}) \\ &= \varepsilon \text{tr}(\mathbf{A}) \end{aligned}$$

- ⊙ When does Hutchinson's Estimator truly need $O\left(\frac{1}{\varepsilon^2}\right)$ queries?
- ⊙ When is the bound $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$ tight?
- ⊙ Let $\mathbf{v} = [\lambda_1 \ \dots \ \lambda_n]$ be the eigenvalues of PSD \mathbf{A}
- ⊙ When is the bound $\|\mathbf{v}\|_2 \leq \|\mathbf{v}\|_1$ tight?
 - Property of norms: $\|\mathbf{v}\|_2 \approx \|\mathbf{v}\|_1$ only if \mathbf{v} is nearly sparse

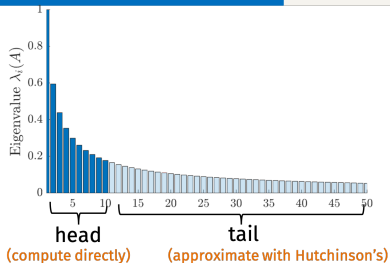
Hutchinson Analysis

For PSD \mathbf{A} , $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$, so that:

$$\begin{aligned} |\tilde{T} - \text{tr}(\mathbf{A})| &\approx O\left(\frac{1}{\sqrt{m}}\right) \|\mathbf{A}\|_F \\ &\leq O\left(\frac{1}{\sqrt{m}}\right) \text{tr}(\mathbf{A}) \\ &= \varepsilon \text{tr}(\mathbf{A}) \end{aligned}$$

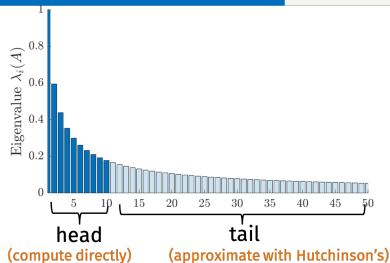
- ⊙ When does Hutchinson's Estimator truly need $O\left(\frac{1}{\varepsilon^2}\right)$ queries?
- ⊙ When is the bound $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$ tight?
- ⊙ Let $\mathbf{v} = [\lambda_1 \ \dots \ \lambda_n]$ be the eigenvalues of PSD \mathbf{A}
- ⊙ When is the bound $\|\mathbf{v}\|_2 \leq \|\mathbf{v}\|_1$ tight?
 - Property of norms: $\|\mathbf{v}\|_2 \approx \|\mathbf{v}\|_1$ only if \mathbf{v} is nearly sparse
- ⊙ Hutchinson only requires $O\left(\frac{1}{\varepsilon^2}\right)$ queries if \mathbf{A} has a few large eigenvalues

Helping Hutchinson's Estimator



Idea: Explicitly estimate the top few eigenvalues of \mathbf{A} . Use Hutchinson's for the rest.

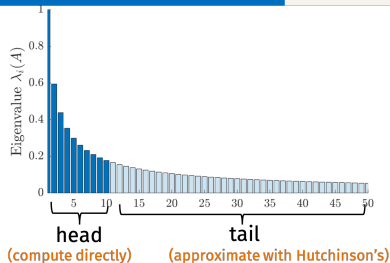
Helping Hutchinson's Estimator



Idea: Explicitly estimate the top few eigenvalues of \mathbf{A} . Use Hutchinson's for the rest.

1. Find a good rank- k approximation $\tilde{\mathbf{A}}_k$
2. Notice that $\text{tr}(\mathbf{A}) = \text{tr}(\tilde{\mathbf{A}}_k) + \text{tr}(\mathbf{A} - \tilde{\mathbf{A}}_k)$
3. Compute $\text{tr}(\tilde{\mathbf{A}}_k)$ exactly
4. Compute $\tilde{T} \approx \text{tr}(\mathbf{A} - \tilde{\mathbf{A}}_k)$ with m steps of Hutchinson's
5. Return $\text{Hutch++}(\mathbf{A}) = \text{tr}(\tilde{\mathbf{A}}_k) + \tilde{T}$

Helping Hutchinson's Estimator



Idea: Explicitly estimate the top few eigenvalues of \mathbf{A} . Use Hutchinson's for the rest.

1. Find a good rank- k approximation $\tilde{\mathbf{A}}_k$
2. Notice that $\text{tr}(\mathbf{A}) = \text{tr}(\tilde{\mathbf{A}}_k) + \text{tr}(\mathbf{A} - \tilde{\mathbf{A}}_k)$
3. Compute $\text{tr}(\tilde{\mathbf{A}}_k)$ exactly
4. Compute $\tilde{T} \approx \text{tr}(\mathbf{A} - \tilde{\mathbf{A}}_k)$ with m steps of Hutchinson's
5. Return $\text{Hutch++}(\mathbf{A}) = \text{tr}(\tilde{\mathbf{A}}_k) + \tilde{T}$

- ⊙ Lemma: $\|\mathbf{A} - \tilde{\mathbf{A}}_k\|_F \leq \frac{2}{\sqrt{k}} \text{tr}(\mathbf{A})$
 - Replaces earlier bound $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$
 - For all \mathbf{v} , there exists k -sparse $\tilde{\mathbf{v}}$ such that

$$\|\mathbf{v} - \tilde{\mathbf{v}}\|_2 \leq \frac{1}{\sqrt{k}} \|\mathbf{v}\|_1$$

- ⊙ Final Theorem:
 - Using rank- k approximation and m samples in Hutchinson's
 - $|\text{tr}(\mathbf{A}) - \text{Hutch++}(\mathbf{A})| \leq O\left(\frac{1}{\sqrt{km}}\right) \text{tr}(\mathbf{A})$
 - Set $k = m = \tilde{O}\left(\frac{1}{\epsilon}\right)$

Implimentation

⊙ Input: Number of matrix-vector queries m

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

```
1 function T = hutchplusplus(A, m)
2     S = 2*randi(2, size(A,1), m/3);
3     G = 2*randi(2, size(A,1), m/3);
4     [Q,~] = qr(A*S,0);
5     G = G - Q*(Q'*G);
6     T = trace(Q'*A*Q) + 1/size(G,2)*trace(G'*A*G);
7 end
```

If you want to learn more

25 Minute Version of this Talk: More Details

- ⊙ Full proof of Hutch++ Correctness
- ⊙ Intuitions for both lower bounds
- ⊙ Discussion of some experiments

In the full paper: Even more details

- ⊙ Non-Adaptive Algorithm
- ⊙ Minor Optimizations
- ⊙ Full Proofs
- ⊙ Richer discussion of experiments

Code: github.com/RaphaelArkadyMeyerNYU/hutchplusplus

- ⊙ **In progress:** Lower bounds for e.g. $\text{tr}(\mathbf{A}^3)$, $\text{tr}(e^{\mathbf{A}})$, $\text{tr}(\mathbf{A}^{-1})$
- ⊙ What about inexact oracles? We often approximate $f(\mathbf{A})\mathbf{x}$ with iterative methods. How accurate do these computations need to be?
- ⊙ Extend to include row/column sampling? This would encapsulate e.g. SGD/SCD.

THANK
YOU

Code available at
github.com/RaphaelArkadyMeyerNYU/hutchplusplus