

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)

Implicit Trace Estimation

Basic problem in linear algebra:

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

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

- Goal is to approximate $\text{tr}(\mathbf{A}) = \sum_{i=1}^n \mathbf{A}_{ii} = \sum_{i=1}^n \lambda_i$

Main Question: How many matrix-vector multiplication queries $\mathbf{Ax}_1, \dots, \mathbf{Ax}_m$ are required to compute $\text{tr}(\mathbf{A})$?¹

¹ \mathbf{x}_j can be chosen *adaptively*, based on the results $\mathbf{Ax}_1, \dots, \mathbf{Ax}_{j-1}$

Application: Trace of a Function of a Matrix

- ⊙ Suppose \mathbf{B} is the adjacency matrix for graph G . Then $\frac{1}{6} \text{tr}(\mathbf{B}^3)$ counts the number of triangles in G .
 - Computing \mathbf{B}^3 directly takes $O(n^3)$ time
 - Computing $\mathbf{B}^3\mathbf{x}$ takes $O(n^2)$ time

Application: Trace of a Function of a Matrix

- ⊙ Suppose \mathbf{B} is the adjacency matrix for graph G . Then $\frac{1}{6} \text{tr}(\mathbf{B}^3)$ counts the number of triangles in G .
 - Computing \mathbf{B}^3 directly takes $O(n^3)$ time
 - Computing $\mathbf{B}^3\mathbf{x}$ takes $O(n^2)$ time
- ⊙ Other functions of interest: $\text{tr}(e^{\mathbf{B}})$, $\text{tr}(\ln(\Sigma))$, etc.

Application: Trace of a Function of a Matrix

- ⊙ Suppose \mathbf{B} is the adjacency matrix for graph G . Then $\frac{1}{6} \text{tr}(\mathbf{B}^3)$ counts the number of triangles in G .
 - Computing \mathbf{B}^3 directly takes $O(n^3)$ time
 - Computing $\mathbf{B}^3\mathbf{x}$ takes $O(n^2)$ time
- ⊙ Other functions of interest: $\text{tr}(e^{\mathbf{B}})$, $\text{tr}(\ln(\Sigma))$, etc.
- ⊙ Computing $f(\mathbf{B})\mathbf{x}$ is often much faster than computing $f(\mathbf{B})$ directly
 - Especially if we only need very few \mathbf{x} vectors

Algorithms:

- ⊙ Krylov Methods, Sketching Methods, Streaming Methods, etc.
- ⊙ See also: *Implicit Matrix Methods*, *Matrix-Free Methods*
- ⊙ Useful framework for algorithmic lower bounds
 - Allows us to prove optimality in a very general setting

Background: Hutchinson's Estimator

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

If $m = O\left(\frac{\log(1/\delta)}{\varepsilon^2}\right)$, then with probability $1 - \delta$,

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

Background: Hutchinson's Estimator

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

If $m = O\left(\frac{\log(1/\delta)}{\varepsilon^2}\right)$, then with probability $1 - \delta$,

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

⊙ If \mathbf{A} is PSD, then $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$, so that

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

Contribution: $O(1/\varepsilon)$ vectors is optimal

Theorems

1. For PSD \mathbf{A} and $m = O(\frac{\log(1/\delta)}{\varepsilon})$, with probability $1 - \delta$,

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

```
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
```

For the rest of the talk, \mathbf{A} is always PSD.

Contribution: $O(1/\varepsilon)$ vectors is optimal

Theorems

1. For PSD \mathbf{A} and $m = O(\frac{\log(1/\delta)}{\varepsilon})$, with probability $1 - \delta$,

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

2. For any b -bit precision oracle, $\tilde{\Omega}(\frac{1}{\varepsilon b})$ possibly adaptive queries are necessary.

```
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
```

For the rest of the talk, \mathbf{A} is always PSD.

Contribution: $O(1/\varepsilon)$ vectors is optimal

Theorems

1. For PSD \mathbf{A} and $m = O(\frac{\log(1/\delta)}{\varepsilon})$, with probability $1 - \delta$,

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

2. For any b -bit precision oracle, $\tilde{\Omega}(\frac{1}{\varepsilon b})$ possibly adaptive queries are necessary.
3. For any infinite precision oracle, $\Omega(\frac{1}{\varepsilon})$ non-adaptive queries are necessary.

```
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
```

For the rest of the talk, \mathbf{A} is always PSD.

Hutchinson's Estimator Versus the Top Few Eigenvalues

Hutchinson Analysis

Let's return to the result for Hutchinson's Estimator:

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

Hutchinson Analysis

Let's return to the result for Hutchinson's Estimator:

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

Hutchinson Analysis

Let's return to the result for Hutchinson's Estimator:

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

Hutchinson Analysis

Let's return to the result for Hutchinson's Estimator:

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

Let's return to the result for Hutchinson's Estimator:

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

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

Let's return to the result for Hutchinson's Estimator:

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

- ⊙ 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

Let's return to the result for Hutchinson's Estimator:

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

- ⊙ 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

Let's return to the result for Hutchinson's Estimator:

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

- ⊙ 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

Let's return to the result for Hutchinson's Estimator:

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

- ⊙ 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

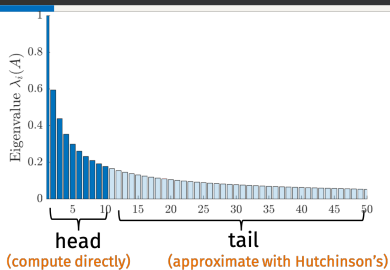
Hutchinson Analysis

Let's return to the result for Hutchinson's Estimator:

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

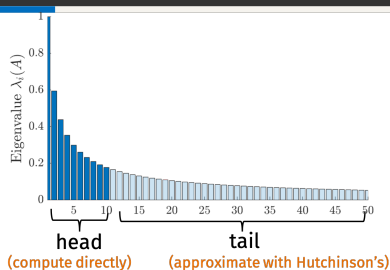
- ⊙ When does Hutchinson's Estimator truly need $O\left(\frac{1}{\varepsilon^2}\right)$ queries?
- ⊙ When is the bound $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$ tight?
- ⊙ Let $\mathbf{v} = [\lambda_1 \ \dots \ \lambda_n]$ be the eigenvalues of PSD \mathbf{A}
- ⊙ When is the bound $\|\mathbf{v}\|_2 \leq \|\mathbf{v}\|_1$ 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 Hutchinson's Estimator
5. Return $\text{Hutch++}(\mathbf{A}) = \text{tr}(\tilde{\mathbf{A}}_k) + \tilde{T}$

Finding a Good Low-Rank Approximation

Let \mathbf{A}_k be the best rank- k approximation of \mathbf{A} .

Lemma (Sarlos 2006, Woodruff 2014)

Let $\mathbf{S} \in \mathbb{R}^{n \times m}$ have i.i.d. uniform ± 1 entries, $\mathbf{Q} = \text{orth}(\mathbf{AS})$, and $\tilde{\mathbf{A}}_k = \mathbf{AQQ}^\top$. Then, with probability $1 - \delta$,

$$\|\mathbf{A} - \tilde{\mathbf{A}}_k\|_F \leq 2\|\mathbf{A} - \mathbf{A}_k\|_F$$

so long as \mathbf{S} has $m = O(k + \log(1/\delta))$ columns.

Finding a Good Low-Rank Approximation

Let \mathbf{A}_k be the best rank- k approximation of \mathbf{A} .

Lemma (Sarlos 2006, Woodruff 2014)

Let $\mathbf{S} \in \mathbb{R}^{n \times m}$ have i.i.d. uniform ± 1 entries, $\mathbf{Q} = \text{orth}(\mathbf{AS})$, and $\tilde{\mathbf{A}}_k = \mathbf{AQQ}^\top$. Then, with probability $1 - \delta$,

$$\|\mathbf{A} - \tilde{\mathbf{A}}_k\|_F \leq 2\|\mathbf{A} - \mathbf{A}_k\|_F$$

so long as \mathbf{S} has $m = O(k + \log(1/\delta))$ columns.

We can compute the trace of $\tilde{\mathbf{A}}_k$ with m queries and $O(mn)$ space:

$$\text{tr}(\tilde{\mathbf{A}}_k) = \text{tr}(\mathbf{AQQ}^\top) = \text{tr}(\mathbf{Q}^\top(\mathbf{A}\mathbf{Q}))$$

Lemma: $\|\mathbf{A} - \mathbf{A}_k\|_F \leq \frac{1}{\sqrt{k}} \text{tr}(\mathbf{A})$

Proof. Note that $\lambda_{k+1} \leq \frac{1}{k} \sum_{i=1}^k \lambda_i \leq \frac{1}{k} \text{tr}(\mathbf{A})$.

Lemma: $\|\mathbf{A} - \mathbf{A}_k\|_F \leq \frac{1}{\sqrt{k}} \text{tr}(\mathbf{A})$

Proof. Note that $\lambda_{k+1} \leq \frac{1}{k} \sum_{i=1}^k \lambda_i \leq \frac{1}{k} \text{tr}(\mathbf{A})$. Then,

$$\|\mathbf{A} - \mathbf{A}_k\|_F^2 = \sum_{i=k+1}^n \lambda_i^2 \leq \lambda_{k+1} \sum_{i=k+1}^n \lambda_i \leq \left(\frac{1}{k} \text{tr}(\mathbf{A})\right) \cdot \text{tr}(\mathbf{A})$$

Lemma: $\|\mathbf{A} - \mathbf{A}_k\|_F \leq \frac{1}{\sqrt{k}} \text{tr}(\mathbf{A})$

Proof. Note that $\lambda_{k+1} \leq \frac{1}{k} \sum_{i=1}^k \lambda_i \leq \frac{1}{k} \text{tr}(\mathbf{A})$. Then,

$$\|\mathbf{A} - \mathbf{A}_k\|_F^2 = \sum_{i=k+1}^n \lambda_i^2 \leq \lambda_{k+1} \sum_{i=k+1}^n \lambda_i \leq \left(\frac{1}{k} \text{tr}(\mathbf{A})\right) \cdot \text{tr}(\mathbf{A})$$

- ⊙ Formalizes our earlier intuition
- ⊙ Replaces the earlier bound $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$
- ⊙ Similar to standard compressed sensing result:

For all $\mathbf{v} \in \mathbb{R}^d$, there exists k -sparse $\tilde{\mathbf{v}}$ such that

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

Complete Analysis

Using rank- k approximation and ℓ sample for Hutchinson's.

Complete Analysis

Using rank- k approximation and ℓ sample for Hutchinson's.

1. We can only make an error in the Hutchinson's step:

$$|\text{tr}(\mathbf{A}) - \text{Hutch}^{++}(\mathbf{A})| = |\text{tr}(\mathbf{A} - \tilde{\mathbf{A}}_k) - \tilde{T}|$$

Using rank- k approximation and ℓ sample for Hutchinson's.

1. We can only make an error in the Hutchinson's step:

$$|\text{tr}(\mathbf{A}) - \text{Hutch}^{++}(\mathbf{A})| = |\text{tr}(\mathbf{A} - \tilde{\mathbf{A}}_k) - \tilde{T}|$$

2. Guarantees for Hutchinson's and Low-Rank Approximation:

$$|\text{tr}(\mathbf{A} - \tilde{\mathbf{A}}_k) - \tilde{T}| \leq O\left(\frac{1}{\sqrt{\ell}}\right) \|\mathbf{A} - \tilde{\mathbf{A}}_k\|_F \leq O\left(\frac{1}{\sqrt{\ell}}\right) \cdot 2 \|\mathbf{A} - \mathbf{A}_k\|_F$$

Using rank- k approximation and ℓ sample for Hutchinson's.

1. We can only make an error in the Hutchinson's step:

$$|\text{tr}(\mathbf{A}) - \text{Hutch}^{++}(\mathbf{A})| = |\text{tr}(\mathbf{A} - \tilde{\mathbf{A}}_k) - \tilde{T}|$$

2. Guarantees for Hutchinson's and Low-Rank Approximation:

$$|\text{tr}(\mathbf{A} - \tilde{\mathbf{A}}_k) - \tilde{T}| \leq O\left(\frac{1}{\sqrt{\ell}}\right) \|\mathbf{A} - \tilde{\mathbf{A}}_k\|_F \leq O\left(\frac{1}{\sqrt{\ell}}\right) \cdot 2\|\mathbf{A} - \mathbf{A}_k\|_F$$

3. Use the lemma from the last slide:

$$|\text{tr}(\mathbf{A}) - \text{Hutch}^{++}(\mathbf{A})| \leq O\left(\frac{1}{\sqrt{k\ell}}\right) \text{tr}(\mathbf{A})$$

Using rank- k approximation and ℓ sample for Hutchinson's.

1. We can only make an error in the Hutchinson's step:

$$|\text{tr}(\mathbf{A}) - \text{Hutch}^{++}(\mathbf{A})| = |\text{tr}(\mathbf{A} - \tilde{\mathbf{A}}_k) - \tilde{T}|$$

2. Guarantees for Hutchinson's and Low-Rank Approximation:

$$|\text{tr}(\mathbf{A} - \tilde{\mathbf{A}}_k) - \tilde{T}| \leq O\left(\frac{1}{\sqrt{\ell}}\right) \|\mathbf{A} - \tilde{\mathbf{A}}_k\|_F \leq O\left(\frac{1}{\sqrt{\ell}}\right) \cdot 2\|\mathbf{A} - \mathbf{A}_k\|_F$$

3. Use the lemma from the last slide:

$$|\text{tr}(\mathbf{A}) - \text{Hutch}^{++}(\mathbf{A})| \leq O\left(\frac{1}{\sqrt{k\ell}}\right) \text{tr}(\mathbf{A})$$

4. If $k = \ell = O\left(\frac{1}{\varepsilon}\right)$, then $|\text{tr}(\mathbf{A}) - \text{Hutch}^{++}(\mathbf{A})| \leq \varepsilon \text{tr}(\mathbf{A}) \square$

Lower Bound:
Communication Complexity

- ⊙ Really rich area of theoretical computing

Gap-Hamming Problem

Let Alice and Bob each have vectors $\mathbf{s}, \mathbf{t} \in \{+1, -1\}^n$. Using as few bits of communication as possible, they must decide if

$$\langle \mathbf{s}, \mathbf{t} \rangle \geq \sqrt{n} \quad \text{or if} \quad \langle \mathbf{s}, \mathbf{t} \rangle \leq -\sqrt{n}$$

Chakrabarti, Regev 2012

Any (possibly adaptive) protocol between Alice and Bob must use $\Omega(n)$ bits to solve the Gap-Hamming problem with probability $\geq \frac{2}{3}$.

A Reduction from Gap-Hamming

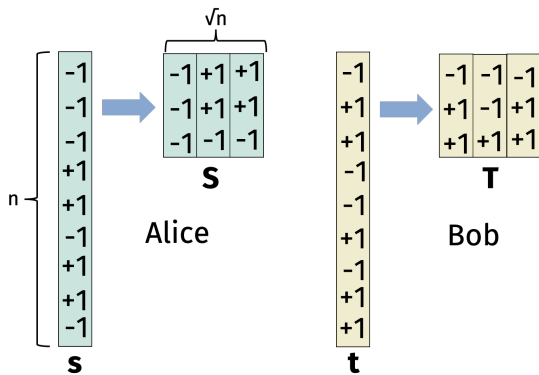
- ⊙ Suppose the Matrix-Vector Oracle for \mathbf{A} only accepts queries with entries that use b bits of precision
 - (e.g. the entries of \mathbf{x} are integers between -2^b and 2^b).

Theorem

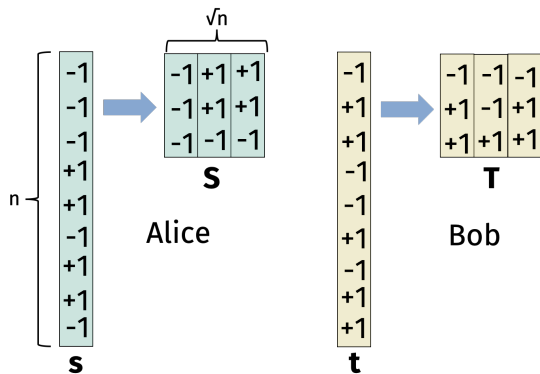
Any (possibly adaptive) algorithm that estimates $\text{tr}(\mathbf{A})$ to relative error ε with probability $\geq \frac{2}{3}$ must use $\Omega\left(\frac{1}{\varepsilon(b+\log(1/\varepsilon))}\right)$ queries.

Proof Idea: Simulate a m -query trace-estimation algorithm to solve a n -bit Gap-Hamming problem

A Reduction to Trace Estimation



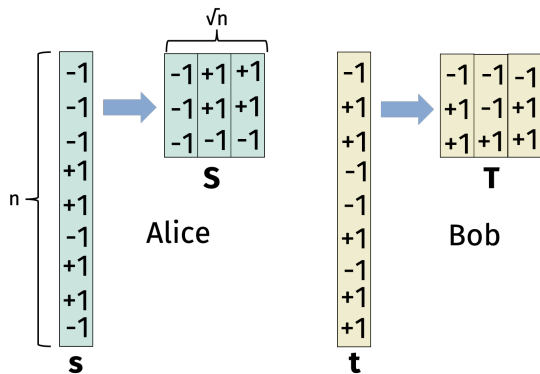
A Reduction to Trace Estimation



Let $\mathbf{Z} = \mathbf{S} + \mathbf{T}$ and $\mathbf{A} = \mathbf{Z}^T \mathbf{Z}$, so that

$$\text{tr}(\mathbf{A}) = \|\mathbf{Z}\|_F^2 = \|\mathbf{s} + \mathbf{t}\|_2^2 = 2n - 2\langle \mathbf{s}, \mathbf{t} \rangle$$

A Reduction to Trace Estimation



Let $\mathbf{Z} = \mathbf{S} + \mathbf{T}$ and $\mathbf{A} = \mathbf{Z}^T \mathbf{Z}$, so that

$$\text{tr}(\mathbf{A}) = \|\mathbf{Z}\|_F^2 = \|\mathbf{s} + \mathbf{t}\|_2^2 = 2n - 2\langle \mathbf{s}, \mathbf{t} \rangle$$

If Alice and Bob can estimate $\text{tr}(\mathbf{A})$ to error $(1 \pm \frac{1}{\sqrt{n}})$, they can solve the Gap-Hamming problem (so $\varepsilon = \frac{1}{\sqrt{n}}$).

- ⊙ For any precision b vector \mathbf{x} , Alice and Bob can compute \mathbf{Ax} with $O(\sqrt{n}(\log(n) + b))$ bits of communication

A Reduction to Trace Estimation

- ⊙ For any precision b vector \mathbf{x} , Alice and Bob can compute \mathbf{Ax} with $O(\sqrt{n}(\log(n) + b))$ bits of communication
- ⊙ They can simulate any m -query trace estimation algorithm with $O(m \cdot \sqrt{n}(\log(n) + b))$ bits of communication

A Reduction to Trace Estimation

- ⊙ For any precision b vector \mathbf{x} , Alice and Bob can compute \mathbf{Ax} with $O(\sqrt{n}(\log(n) + b))$ bits of communication
- ⊙ They can simulate any m -query trace estimation algorithm with $O(m \cdot \sqrt{n}(\log(n) + b))$ bits of communication
- ⊙ Gap-Hamming Lower bound: $m \geq \Omega\left(\frac{n}{\sqrt{n}(\log(n)+b)}\right)$

A Reduction to Trace Estimation

- ⊙ For any precision b vector \mathbf{x} , Alice and Bob can compute \mathbf{Ax} with $O(\sqrt{n}(\log(n) + b))$ bits of communication
- ⊙ They can simulate any m -query trace estimation algorithm with $O(m \cdot \sqrt{n}(\log(n) + b))$ bits of communication
- ⊙ Gap-Hamming Lower bound: $m \geq \Omega\left(\frac{n}{\sqrt{n}(\log(n)+b)}\right)$
- ⊙ Substitute $\varepsilon = \frac{1}{\sqrt{n}}$: $m \geq \Omega\left(\frac{1}{\varepsilon(b+\log(1/\varepsilon))}\right)$

Lower Bound:
Statistical Hypothesis Testing

Design distributions \mathcal{P}_0 and \mathcal{P}_1 over PSD matrices such that

1. A trace estimator can distinguish \mathcal{P}_0 from \mathcal{P}_1
 - If $\mathbf{A}_0 \sim \mathcal{P}_0$ and $\mathbf{A}_1 \sim \mathcal{P}_1$
 - With high probability, $\text{tr}(\mathbf{A}_0) \leq (1 - 2\varepsilon) \text{tr}(\mathbf{A}_1)$

Statistical Hypothesis Testing

Design distributions \mathcal{P}_0 and \mathcal{P}_1 over PSD matrices such that

1. A trace estimator can distinguish \mathcal{P}_0 from \mathcal{P}_1
 - If $\mathbf{A}_0 \sim \mathcal{P}_0$ and $\mathbf{A}_1 \sim \mathcal{P}_1$
 - With high probability, $\text{tr}(\mathbf{A}_0) \leq (1 - 2\varepsilon) \text{tr}(\mathbf{A}_1)$
2. No estimator can distinguish \mathcal{P}_0 from \mathcal{P}_1 with $\Omega(\frac{1}{\varepsilon})$ queries
 - Nature samples $i \sim \{0, 1\}$, and $\mathbf{A} \sim \mathcal{P}_i$
 - Any estimator that correctly guesses i with probability $\geq \frac{3}{4}$ must use $\Omega(\frac{1}{\varepsilon})$ queries

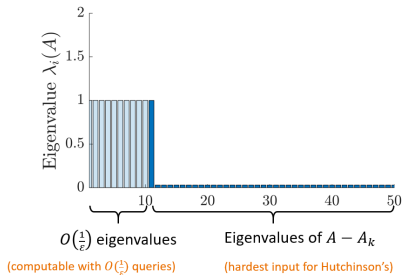
Statistical Hypothesis Testing

Design distributions \mathcal{P}_0 and \mathcal{P}_1 over PSD matrices such that

1. A trace estimator can distinguish \mathcal{P}_0 from \mathcal{P}_1
 - If $\mathbf{A}_0 \sim \mathcal{P}_0$ and $\mathbf{A}_1 \sim \mathcal{P}_1$
 - With high probability, $\text{tr}(\mathbf{A}_0) \leq (1 - 2\varepsilon) \text{tr}(\mathbf{A}_1)$
2. No estimator can distinguish \mathcal{P}_0 from \mathcal{P}_1 with $\Omega(\frac{1}{\varepsilon})$ queries
 - Nature samples $i \sim \{0, 1\}$, and $\mathbf{A} \sim \mathcal{P}_i$
 - Any estimator that correctly guesses i with probability $\geq \frac{3}{4}$ must use $\Omega(\frac{1}{\varepsilon})$ queries

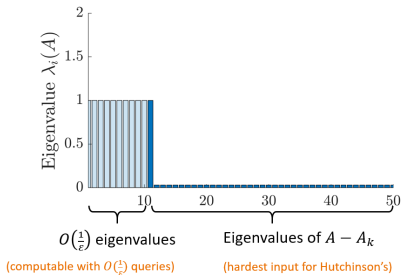
The design of \mathcal{P}_0 and \mathcal{P}_1 should reflect what structure makes trace estimation hard!

Designing a Hard Instance



What would the hardest input for Hutch++ be?

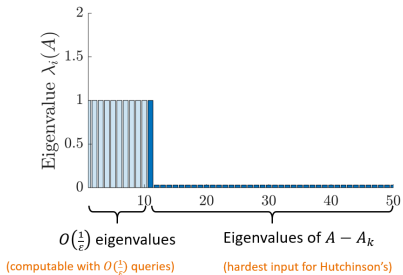
Designing a Hard Instance



What would the hardest input for Hutch++ be?

- ⊙ Hutch++ only makes errors with Hutchinson's estimator on $\text{tr}(\mathbf{A} - \tilde{\mathbf{A}}_k)$

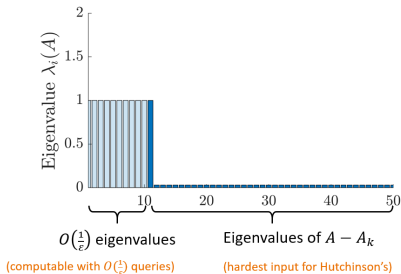
Designing a Hard Instance



What would the hardest input for Hutch++ be?

- ⊙ Hutch++ only makes errors with Hutchinson's estimator on $\text{tr}(\mathbf{A} - \tilde{\mathbf{A}}_k)$
- ⊙ For what \mathbf{A} would Hutchinson's estimator have difficulty estimating $\text{tr}(\mathbf{A} - \mathbf{A}_k)$?

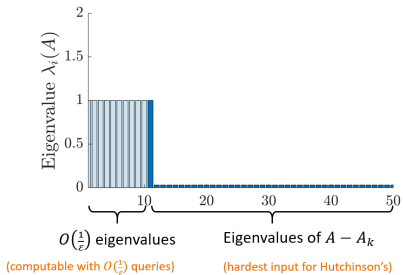
Designing a Hard Instance



What would the hardest input for Hutch++ be?

- ⊙ Hutch++ only makes errors with Hutchinson's estimator on $\text{tr}(\mathbf{A} - \tilde{\mathbf{A}}_k)$
- ⊙ For what \mathbf{A} would Hutchinson's estimator have difficulty estimating $\text{tr}(\mathbf{A} - \mathbf{A}_k)$?
 - Hutchinson's estimator needs many samples when $\mathbf{A} - \mathbf{A}_k$ has concentrated eigenvalues

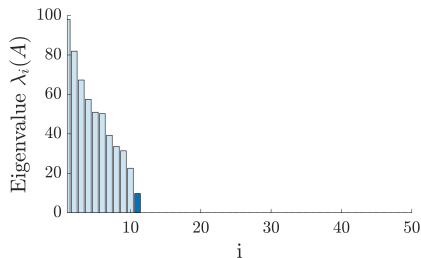
Designing a Hard Instance



What would the hardest input for Hutch++ be?

- ⊙ Hutch++ only makes errors with Hutchinson's estimator on $\text{tr}(\mathbf{A} - \tilde{\mathbf{A}}_k)$
- ⊙ For what \mathbf{A} would Hutchinson's estimator have difficulty estimating $\text{tr}(\mathbf{A} - \mathbf{A}_k)$?
 - Hutchinson's estimator needs many samples when $\mathbf{A} - \mathbf{A}_k$ has concentrated eigenvalues
- ⊙ \mathbf{A} has $k = O(\frac{1}{\epsilon})$ large eigenvalues. The rest are zero.

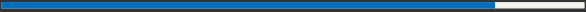
Designing a Hard Instance



Formally, for large enough integer d ,

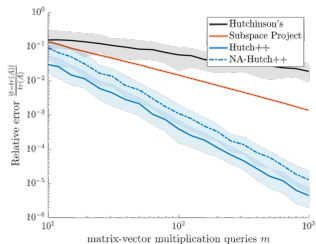
$$\begin{array}{l|l} \mathcal{P}_0 & \mathbf{A} = \mathbf{G}^T \mathbf{G} \text{ for } \mathbf{G} \in \mathbb{R}^{d \times (\frac{1}{\epsilon})} \text{ Gaussian} \\ \hline \mathcal{P}_1 & \mathbf{A} = \mathbf{G}^T \mathbf{G} \text{ for } \mathbf{G} \in \mathbb{R}^{d \times (\frac{1}{\epsilon} + 1)} \text{ Gaussian} \end{array}$$

Experiments

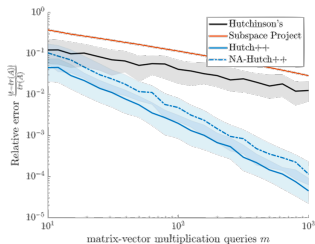


Synthetic Experiments

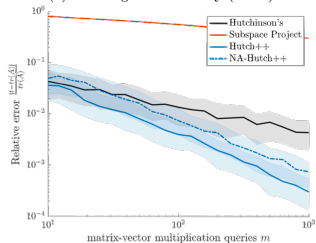
Results on synthetic matrix \mathbf{A} with spectrum $\lambda_i = i^{-c}$ for different values of c :



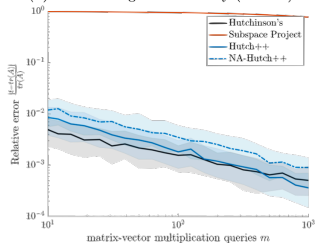
(a) Fast Eigenvalue Decay ($c = 2$)



(b) Medium Eigenvalue Decay ($c = 1.5$)



(c) Slow Eigenvalue Decay ($c = 1$)



(d) Very Slow Eigenvalue Decay ($c = .5$)

Non-PSD Experiments

Hutch++ works well empirically for many non-PSD matrices.

Let \mathbf{B} be the (indefinite) adjacency matrix of an undirected graph G , $\frac{1}{6} \text{tr}(\mathbf{B}^3)$ is exactly equal to the number of *triangles* in G .

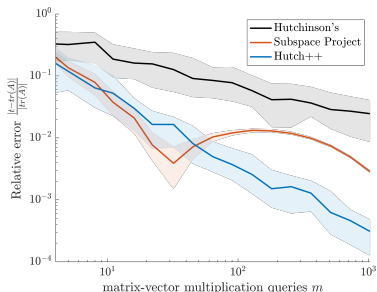
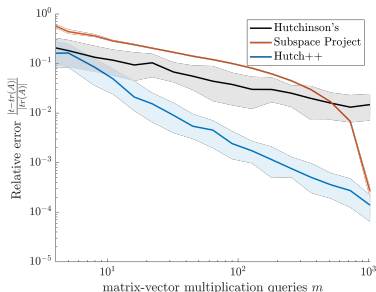


Figure: $\mathbf{A} = \mathbf{B}^3$ for arXiv.org citation network and Wikipedia voting network.

- ⊙ **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