# 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)

1. Introduction
    - What problems am I solving?
    - Why are these problems interesting?
    - How am I solving them?
2. Trace Estimation *(SOSA 2021)*
3. Trace Monomial Estimation *(Ongoing Research)*

- ⊙ Scientific Computing relies on Numerical Linear Algebra
- ⊙ We spent decades building better algorithms

## Numerical Linear Algebra

⊙ Scientific Computing relies on Numerical Linear Algebra

⊙ We spent decades building better algorithms

⊙ We don't know which algorithms are *optimal*

  ○ Krylov Iteration is optimal for top eigenvalue
  ○ Hutchinson's Estimator is suboptimal for trace estimation

# Numerical Linear Algebra

- Scientific Computing relies on Numerical Linear Algebra
- We spent decades building better algorithms
- We don't know which algorithms are *optimal*
  - Krylov Iteration is optimal for top eigenvalue
  - Hutchinson's Estimator is suboptimal for trace estimation
- My goal: Prove the optimality of linear algebra algorithms
  - Emphasis on building lower bounds

⊙ Goal: Estimate trace of $d \times d$ matrix $\boldsymbol{A}$:

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

# Trace Estimation

◎ Goal: Estimate trace of $d \times d$ matrix $\boldsymbol{A}$:

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

◎ In Downstream Applications, $\boldsymbol{A}$ is not stored in memory.

◎ Instead, $\boldsymbol{B}$ is in memory and $\boldsymbol{A} = f(\boldsymbol{B})$:

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

- Goal: Estimate trace of $d \times d$ matrix $\boldsymbol{A}$:

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

- In Downstream Applications, $\boldsymbol{A}$ is not stored in memory.

- Instead, $\boldsymbol{B}$ is in memory and $\boldsymbol{A} = f(\boldsymbol{B})$:

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

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

## Trace Estimation

◎ Goal: Estimate trace of $d \times d$ matrix $\boldsymbol{A}$:

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

◎ In Downstream Applications, $\boldsymbol{A}$ is not stored in memory.

◎ Instead, $\boldsymbol{B}$ is in memory and $\boldsymbol{A} = f(\boldsymbol{B})$:

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

◎ Computing $\boldsymbol{A} = \frac{1}{6}\boldsymbol{B}^3$ takes $O(n^3)$ time

◎ Computing $\boldsymbol{A}\mathbf{x} = \frac{1}{6}\boldsymbol{B}(\boldsymbol{B}(\boldsymbol{B}\mathbf{x}))$ takes $O(n^2)$ time

◎ If $\boldsymbol{A} = f(\boldsymbol{B})$, then we can often compute $\boldsymbol{A}\mathbf{x}$ quickly

◎ Goal: Estimate $\text{tr}(\boldsymbol{A})$ by computing $\boldsymbol{A}\mathbf{x}_1, \ldots \boldsymbol{A}\mathbf{x}_k$

# Matrix-Vector Oracle Model

Formally: Matrix-Vector Product as a Computational Primitive

## Matrix-Vector Oracle Model

Formally: Matrix-Vector Product as a Computational Primitive

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

$$\mathbf{x} \quad \overset{input}{\Longrightarrow} \quad \text{ORACLE} \quad \overset{output}{\Longrightarrow} \quad \mathbf{A}\mathbf{x}$$

- ◉ e.g. Krylov Methods, Sketching, Streaming, . . .
- ◉ Very few existing lower bounds

# Matrix-Vector Oracle Model

Formally: Matrix-Vector Product as a Computational Primitive

- ◉ Given access to a $d \times d$ matrix $\boldsymbol{A}$ only through a
  <span style="color:red">Matrix-Vector Multiplication Oracle</span>

$$\mathbf{x} \xrightarrow{\textit{input}} \text{ORACLE} \xrightarrow{\textit{output}} \boldsymbol{A}\mathbf{x}$$

- ◉ e.g. Krylov Methods, Sketching, Streaming, ...
- ◉ Very few existing lower bounds

**Trace Estimation:** Estimate $\text{tr}(\boldsymbol{A})$ with as few Matrix-Vector products $\boldsymbol{A}\mathbf{x}_1, \ldots, \boldsymbol{A}\mathbf{x}_k$ as possible.

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

## Our Contributions

Prior Work:

- ⊙ Hutchinson's Estimator: $O(\frac{1}{\varepsilon^2})$ products suffice [AT11]
  - ○ 2 Lines of MATLAB code
- ⊙ Lower Bound: Hutchinson's Estimator needs $\Omega(\frac{1}{\varepsilon^2})$ products [WWZ14]

Our Results:

- ⊙ Hutch++ Estimator: $O(\frac{1}{\varepsilon})$ products suffice
  - ○ 5 Lines of MATLAB code
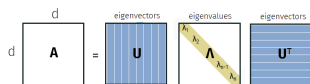- ⊙ Lower Bound: Any estimator needs $\Omega(\frac{1}{\varepsilon})$ products

# Linear Algebra Review



- Symmetric $\boldsymbol{A} \in \mathbb{R}^{d \times d}$ has $\boldsymbol{A} = \boldsymbol{U} \boldsymbol{\Lambda} \boldsymbol{U}^{\mathsf{T}}$

- $\boldsymbol{U}$ is a rotation matrix: $\boldsymbol{U}^{\mathsf{T}} \boldsymbol{U} = \boldsymbol{I}$

- Eigenvalues $\lambda_1 \geq \lambda_2 \geq \ldots \geq \lambda_d$

# Linear Algebra Review



- ⊙ Symmetric $\boldsymbol{A} \in \mathbb{R}^{d \times d}$ has $\boldsymbol{A} = \boldsymbol{U}\boldsymbol{\Lambda}\boldsymbol{U}^\mathsf{T}$

- ⊙ $\boldsymbol{U}$ is a rotation matrix: $\boldsymbol{U}^\mathsf{T}\boldsymbol{U} = \boldsymbol{I}$

- ⊙ Eigenvalues $\lambda_1 \geq \lambda_2 \geq \ldots \geq \lambda_d$

- ⊙ $\|\boldsymbol{A}\|_F^2 = \sum_{i,j} \boldsymbol{A}_{i,j}^2 = \sum_i \lambda_i^2$

- ⊙ $\mathrm{tr}(\boldsymbol{A}) = \sum_i \boldsymbol{A}_{i,i} = \sum_i \lambda_i$

# Linear Algebra Review



- ⊙ Symmetric $\boldsymbol{A} \in \mathbb{R}^{d \times d}$ has $\boldsymbol{A} = \boldsymbol{U} \boldsymbol{\Lambda} \boldsymbol{U}^\mathsf{T}$
- ⊙ $\boldsymbol{U}$ is a rotation matrix: $\boldsymbol{U}^\mathsf{T} \boldsymbol{U} = \boldsymbol{I}$
- ⊙ Eigenvalues $\lambda_1 \geq \lambda_2 \geq \ldots \geq \lambda_d$
- ⊙ $\|\boldsymbol{A}\|_F^2 = \sum_{i,j} \boldsymbol{A}_{i,j}^2 = \sum_i \lambda_i^2$
- ⊙ $\mathrm{tr}(\boldsymbol{A}) = \sum_i \boldsymbol{A}_{i,i} = \sum_i \lambda_i$
- ⊙ Positive Semi-Definite (PSD) $\boldsymbol{A}$ has $\lambda_i \geq 0$ for all $i$
  - ○ $\|\boldsymbol{A}\|_F = \|\boldsymbol{\lambda}\|_2 \leq \|\boldsymbol{\lambda}\|_1 = \mathrm{tr}(\boldsymbol{A})$

# Linear Algebra Review



- ⊙ Symmetric $\boldsymbol{A} \in \mathbb{R}^{d \times d}$ has $\boldsymbol{A} = \boldsymbol{U\Lambda U}^\mathsf{T}$

- ⊙ $\boldsymbol{U}$ is a rotation matrix: $\boldsymbol{U}^\mathsf{T}\boldsymbol{U} = \boldsymbol{I}$

- ⊙ Eigenvalues $\lambda_1 \geq \lambda_2 \geq \ldots \geq \lambda_d$

- ⊙ $\|\boldsymbol{A}\|_F^2 = \sum_{i,j} \boldsymbol{A}_{i,j}^2 = \sum_i \lambda_i^2$

- ⊙ $\mathrm{tr}(\boldsymbol{A}) = \sum_i \boldsymbol{A}_{i,i} = \sum_i \lambda_i$

- ⊙ Positive Semi-Definite (PSD) $\boldsymbol{A}$ has $\lambda_i \geq 0$ for all $i$
  - ○ $\|\boldsymbol{A}\|_F = \|\boldsymbol{\lambda}\|_2 \leq \|\boldsymbol{\lambda}\|_1 = \mathrm{tr}(\boldsymbol{A})$

- ⊙ Low Rank Approximation:
  $\boldsymbol{A}_k = \boldsymbol{U}_k \boldsymbol{\Lambda}_k \boldsymbol{U}_k^\mathsf{T} = \mathrm{argmin}_{rank(\boldsymbol{B})=k} \|\boldsymbol{A} - \boldsymbol{B}\|_F$

⊙ If $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, then $\mathbf{A}\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{A}\mathbf{A}^\mathsf{T})$

⊙ If $X_1, \ldots, X_n \sim \mathcal{N}(0, 1)$, then $S := \sum_i X_i^2 \sim \chi_n^2$, $\mathbb{E}[S] = n$, $\mathrm{Var}[S] = 2n$

⊚ If $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{I})$, then $\boldsymbol{A}\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{A}\boldsymbol{A}^\mathsf{T})$

⊚ If $X_1, \ldots, X_n \sim \mathcal{N}(0, 1)$, then $S := \sum_i X_i^2 \sim \chi_n^2$, $\mathbb{E}[S] = n$, $\mathrm{Var}[S] = 2n$

⊚ Chebyshev's Ineq: $|X - \mathbb{E}[X]| \leq \frac{1}{\sqrt{\delta}} \sqrt{\mathrm{Var}[X]}$ w.p. $\geq 1 - \delta$

⊙ If $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{I})$, then $\boldsymbol{A}\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{A}\boldsymbol{A}^\mathsf{T})$

⊙ If $X_1, \ldots, X_n \sim \mathcal{N}(0, 1)$, then $S := \sum_i X_i^2 \sim \chi_n^2$, $\mathbb{E}[S] = n$, $\mathrm{Var}[S] = 2n$

⊙ Chebyshev's Ineq: $|X - \mathbb{E}[X]| \leq \frac{1}{\sqrt{\delta}}\sqrt{\mathrm{Var}[X]}$ w.p. $\geq 1 - \delta$

⊙ Chebyshev's Ineq: $|X - \mathbb{E}[X]| \leq O(\sqrt{\mathrm{Var}[X]})$ w.p. $\geq \frac{2}{3}$

Towards Optimal

Trace Estimation in the

Matrix-Vector Oracle Model

## Hutchinson's Estimator

⊙ If $\mathbf{x} \sim \mathcal{N}(0, \boldsymbol{I})$, then

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

⊙ If $\mathbf{x} \sim \mathcal{N}(0, \boldsymbol{I})$, then

$$\mathbb{E}[\mathbf{x}^\mathsf{T} \boldsymbol{A} \mathbf{x}] = \mathrm{tr}(\boldsymbol{A}) \qquad\qquad \mathsf{Var}[\mathbf{x}^\mathsf{T} \boldsymbol{A} \mathbf{x}] = 2\|\boldsymbol{A}\|_F^2$$

⊙ Hutchinson's Estimator: $\mathsf{H}_\ell(\boldsymbol{A}) := \frac{1}{\ell} \sum_{i=1}^{\ell} \mathbf{x}_i^\mathsf{T} \boldsymbol{A} \mathbf{x}_i$

$$\mathbb{E}[\mathsf{H}_\ell(\boldsymbol{A})] = \mathrm{tr}(\boldsymbol{A}) \qquad\qquad \mathsf{Var}[\mathsf{H}_\ell(\boldsymbol{A})] = \frac{2}{\ell}\|\boldsymbol{A}\|_F^2$$

# Hutchinson's Estimator

⊙ If $\mathbf{x} \sim \mathcal{N}(0, \boldsymbol{I})$, then

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

⊙ Hutchinson's Estimator: $\mathsf{H}_\ell(\boldsymbol{A}) := \frac{1}{\ell} \sum_{i=1}^{\ell} \mathbf{x}_i^\mathsf{T} \boldsymbol{A} \mathbf{x}_i$

$$\mathbb{E}[\mathsf{H}_\ell(\boldsymbol{A})] = \text{tr}(\boldsymbol{A}) \qquad\qquad \text{Var}[\mathsf{H}_\ell(\boldsymbol{A})] = \frac{2}{\ell}\|\boldsymbol{A}\|_F^2$$

## Proof: $\mathsf{H}_\ell(\boldsymbol{A})$ needs $\ell = O(\frac{1}{\varepsilon^2})$ for PSD $\boldsymbol{A}$

⊙ For PSD $\boldsymbol{A}$, we have $\|\boldsymbol{A}\|_F \leq \text{tr}(\boldsymbol{A})$, so that

# Hutchinson's Estimator

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

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

- Hutchinson's Estimator: $H_\ell(\mathbf{A}) := \frac{1}{\ell} \sum_{i=1}^{\ell} \mathbf{x}_i^\intercal \mathbf{A} \mathbf{x}_i$

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

## Proof: $H_\ell(\mathbf{A})$ needs $\ell = O(\frac{1}{\varepsilon^2})$ for PSD $\mathbf{A}$

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

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

For what $\boldsymbol{A}$ is this analysis tight?

$$\begin{aligned}
|\mathsf{H}_\ell(\boldsymbol{A}) - \mathsf{tr}(\boldsymbol{A})| &\leq O(\tfrac{1}{\sqrt{\ell}})\|\boldsymbol{A}\|_F \\
&\leq O(\tfrac{1}{\sqrt{\ell}})\,\mathsf{tr}(\boldsymbol{A}) \\
&= \varepsilon\,\mathsf{tr}(\boldsymbol{A})
\end{aligned}$$

For what $\boldsymbol{A}$ is this analysis tight?

$$\begin{aligned}
|H_\ell(\boldsymbol{A}) - \text{tr}(\boldsymbol{A})| &\approx O(\tfrac{1}{\sqrt{\ell}})\|\boldsymbol{A}\|_F \\
&\leq O(\tfrac{1}{\sqrt{\ell}})\,\text{tr}(\boldsymbol{A}) \\
&= \varepsilon\,\text{tr}(\boldsymbol{A})
\end{aligned}$$

For what $\boldsymbol{A}$ is this analysis tight?

$$\begin{aligned} |H_\ell(\boldsymbol{A}) - \text{tr}(\boldsymbol{A})| &\approx O(\tfrac{1}{\sqrt{\ell}}) \|\boldsymbol{A}\|_F \\ &\leq O(\tfrac{1}{\sqrt{\ell}}) \, \text{tr}(\boldsymbol{A}) \\ &= \varepsilon \, \text{tr}(\boldsymbol{A}) \end{aligned}$$

⊙ When is the bound $\|\boldsymbol{A}\|_F \leq \text{tr}(A)$ tight?

## Hutchinson's Estimator

For what $\boldsymbol{A}$ is this analysis tight?

$$|H_\ell(\boldsymbol{A}) - \text{tr}(\boldsymbol{A})| \approx O(\tfrac{1}{\sqrt{\ell}}) \|\boldsymbol{A}\|_F$$
$$\leq O(\tfrac{1}{\sqrt{\ell}}) \, \text{tr}(\boldsymbol{A})$$
$$= \varepsilon \, \text{tr}(\boldsymbol{A})$$

⊙ When is the bound $\|\boldsymbol{A}\|_F \leq \text{tr}(A)$ tight?

⊙ Let $\mathbf{v} = \begin{bmatrix} \lambda_1 & \ldots & \lambda_n \end{bmatrix}$ be the eigenvalues of PSD $\boldsymbol{A}$

## Hutchinson's Estimator

For what $A$ is this analysis tight?

$$|H_\ell(A) - tr(A)| \approx O(\tfrac{1}{\sqrt{\ell}}) \|A\|_F$$
$$\leq O(\tfrac{1}{\sqrt{\ell}}) tr(A)$$
$$= \varepsilon\, tr(A)$$

⊙ When is the bound $\|A\|_F \leq tr(A)$ tight?

⊙ Let $v = \begin{bmatrix} \lambda_1 & \dots & \lambda_n \end{bmatrix}$ be the eigenvalues of PSD $A$

⊙ When is the bound $\|v\|_2 \leq \|v\|_1$ tight?

## Hutchinson's Estimator

For what $\boldsymbol{A}$ is this analysis tight?

$$|\mathrm{H}_\ell(\boldsymbol{A}) - \mathrm{tr}(\boldsymbol{A})| \approx O(\tfrac{1}{\sqrt{\ell}}) \|\boldsymbol{A}\|_F$$
$$\leq O(\tfrac{1}{\sqrt{\ell}}) \, \mathrm{tr}(\boldsymbol{A})$$
$$= \varepsilon \, \mathrm{tr}(\boldsymbol{A})$$

⦿ When is the bound $\|\boldsymbol{A}\|_F \leq \mathrm{tr}(A)$ tight?

⦿ Let $\mathbf{v} = \begin{bmatrix} \lambda_1 & \ldots & \lambda_n \end{bmatrix}$ be the eigenvalues of PSD $\boldsymbol{A}$

⦿ 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's Estimator

For what $\boldsymbol{A}$ is this analysis tight?

$$|\mathsf{H}_\ell(\boldsymbol{A}) - \mathrm{tr}(\boldsymbol{A})| \approx O(\tfrac{1}{\sqrt{\ell}})\|\boldsymbol{A}\|_F$$
$$\leq O(\tfrac{1}{\sqrt{\ell}})\,\mathrm{tr}(\boldsymbol{A})$$
$$= \varepsilon\,\mathrm{tr}(\boldsymbol{A})$$

◎ When is the bound $\|\boldsymbol{A}\|_F \leq \mathrm{tr}(A)$ tight?

◎ Let $\mathbf{v} = \begin{bmatrix} \lambda_1 & \dots & \lambda_n \end{bmatrix}$ be the eigenvalues of PSD $\boldsymbol{A}$

◎ When is the bound $\|\mathbf{v}\|_2 \leq \|\mathbf{v}\|_1$ tight?

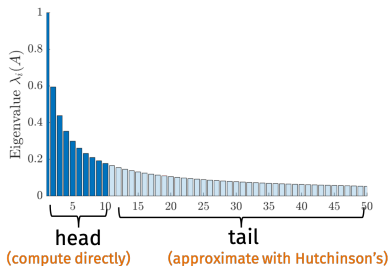    ○ Property of norms: $\|\mathbf{v}\|_2 \approx \|\mathbf{v}\|_1$ only if $\mathbf{v}$ is nearly sparse

◎ Hutchinson only requires $O(\tfrac{1}{\varepsilon^2})$ queries if $\boldsymbol{A}$ has a few large eigenvalues

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

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

1. Find a good rank-$k$ approximation $\tilde{\boldsymbol{A}}_k$
2. Notice that $\text{tr}(\boldsymbol{A}) = \text{tr}(\tilde{\boldsymbol{A}}_k) + \text{tr}(\boldsymbol{A} - \tilde{\boldsymbol{A}}_k)$
3. Compute $\text{tr}(\tilde{\boldsymbol{A}}_k)$ exactly
4. Return Hutch++$(\boldsymbol{A}) = \text{tr}(\tilde{\boldsymbol{A}}_k) + \text{H}_\ell(\boldsymbol{A} - \tilde{\boldsymbol{A}}_k)$

# Helping Hutchinson's Estimator



head (compute directly)

tail (approximate with Hutchinson's)
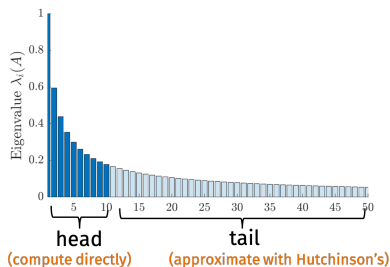
Idea: Explicitly estimate the top few eigenvalues of $\boldsymbol{A}$. Use Hutchinson's for the rest.

1. Find a good rank-$k$ approximation $\tilde{\boldsymbol{A}}_k$
2. Notice that $\operatorname{tr}(\boldsymbol{A}) = \operatorname{tr}(\tilde{\boldsymbol{A}}_k) + \operatorname{tr}(\boldsymbol{A} - \tilde{\boldsymbol{A}}_k)$
3. Compute $\operatorname{tr}(\tilde{\boldsymbol{A}}_k)$ exactly
4. Return $\mathsf{Hutch}\text{++}(\boldsymbol{A}) = \operatorname{tr}(\tilde{\boldsymbol{A}}_k) + \mathsf{H}_\ell(\boldsymbol{A} - \tilde{\boldsymbol{A}}_k)$

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

(Whiteboard)

# Finding a Good Low-Rank Approximation

Let $A_k$ be the best rank-$k$ approximation of $A$.

## Lemma [Sar06, Woo14]

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

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

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

## Finding a Good Low-Rank Approximation

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

### Lemma [Sar06, Woo14]

Let $\boldsymbol{S} \in \mathbb{R}^{d \times k}$ have i.i.d. uniform $\pm 1$ entries, $\boldsymbol{Q} = \text{orth}(\boldsymbol{AS})$, and $\tilde{\boldsymbol{A}}_k = \boldsymbol{AQQ}^\mathsf{T}$. Then, with probability $1 - \delta$,

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

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

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

$$\text{tr}(\tilde{\boldsymbol{A}}_k) = \text{tr}(\boldsymbol{AQQ}^\mathsf{T}) = \text{tr}(\boldsymbol{Q}^\mathsf{T}(\boldsymbol{AQ}))$$

**Hutch++ Algorithm:**
- ◉ Input: Number of matrix-vector queries $m$, matrix $\boldsymbol{A}$
1. Sample $\boldsymbol{S} \in \mathbb{R}^{d \times \frac{m}{3}}$ and $\boldsymbol{G} \in \mathbb{R}^{d \times \frac{m}{3}}$ with i.i.d. $\mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$ entries
2. Compute $\boldsymbol{Q} = \mathrm{qr}(\boldsymbol{AS})$
3. Return $\mathrm{tr}(\boldsymbol{Q}^T \boldsymbol{A} \boldsymbol{Q}) + \frac{3}{m} \mathrm{tr}(\boldsymbol{G}^T(\boldsymbol{I} - \boldsymbol{Q}\boldsymbol{Q}^\mathsf{T})\boldsymbol{A}(\boldsymbol{I} - \boldsymbol{Q}\boldsymbol{Q}^\mathsf{T})\boldsymbol{G})$

# Hutch++

**Hutch++ Algorithm:**

- ⊙ Input: Number of matrix-vector queries $m$, matrix $\boldsymbol{A}$
1. Sample $\boldsymbol{S} \in \mathbb{R}^{d \times \frac{m}{3}}$ and $\boldsymbol{G} \in \mathbb{R}^{d \times \frac{m}{3}}$ with i.i.d. $\mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$ entries
2. Compute $\boldsymbol{Q} = \text{qr}(\boldsymbol{A}\boldsymbol{S})$
3. Return $\text{tr}(\boldsymbol{Q}^T \boldsymbol{A} \boldsymbol{Q}) + \frac{3}{m} \text{tr}(\boldsymbol{G}^T (\boldsymbol{I} - \boldsymbol{Q}\boldsymbol{Q}^\mathsf{T}) \boldsymbol{A} (\boldsymbol{I} - \boldsymbol{Q}\boldsymbol{Q}^\mathsf{T}) \boldsymbol{G})$

This algorithm is adaptive:

$$\mathbf{x}_{k+1} \Longrightarrow \text{ORACLE} \Longrightarrow \boldsymbol{A}\mathbf{x}_k$$
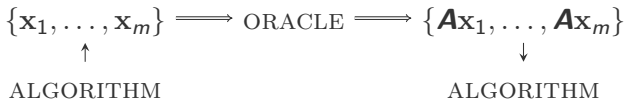$$\uparrow \qquad\qquad \text{ALGORITHM} \longleftarrow$$

**Hutch++ Algorithm:**
- ⊙ Input: Number of matrix-vector queries $m$, matrix $\boldsymbol{A}$
1. Sample $\boldsymbol{S} \in \mathbb{R}^{d \times \frac{m}{3}}$ and $\boldsymbol{G} \in \mathbb{R}^{d \times \frac{m}{3}}$ with i.i.d. $\mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$ entries
2. Compute $\boldsymbol{Q} = \mathrm{qr}(\boldsymbol{AS})$
3. Return $\mathrm{tr}(\boldsymbol{Q}^T \boldsymbol{A} \boldsymbol{Q}) + \frac{3}{m} \mathrm{tr}(\boldsymbol{G}^T(\boldsymbol{I} - \boldsymbol{Q}\boldsymbol{Q}^\mathsf{T})\boldsymbol{A}(\boldsymbol{I} - \boldsymbol{Q}\boldsymbol{Q}^\mathsf{T})\boldsymbol{G})$

This algorithm is adaptive:

$$\mathbf{x}_{k+1} \Longrightarrow \mathrm{ORACLE} \Longrightarrow \boldsymbol{A}\mathbf{x}_k$$
$$\uparrow \qquad\qquad \mathrm{ALGORITHM} \longleftarrow$$

There is a non-adaptive variant of Hutch++:

$$\{\mathbf{x}_1, \ldots, \mathbf{x}_m\} \Longrightarrow \mathrm{ORACLE} \Longrightarrow \{\boldsymbol{A}\mathbf{x}_1, \ldots, \boldsymbol{A}\mathbf{x}_m\}$$
$$\uparrow \qquad\qquad\qquad\qquad\qquad \downarrow$$
$$\mathrm{ALGORITHM} \qquad\qquad\qquad\qquad \mathrm{ALGORITHM}$$

## Experiments

When $\|\boldsymbol{A}\|_F \approx \mathrm{tr}(\boldsymbol{A})$, Hutch++ is much faster than H$_\ell$:



(a) $\|\boldsymbol{A}\|_F = 0.63\,\mathrm{tr}(\boldsymbol{A})$      (b) $\|\boldsymbol{A}\|_F = 0.02\,\mathrm{tr}(\boldsymbol{A})$

```
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
```

# Trace Estimation Lower Bounds

$$\mathbf{x} \quad \overset{input}{\Longrightarrow} \quad \text{ORACLE} \quad \overset{output}{\Longrightarrow} \quad \boldsymbol{A}\mathbf{x}$$

View oracle as a limit on information about $\boldsymbol{A}$:

1. Suppose $\boldsymbol{A} \sim \mathcal{D}$ is a random matrix
2. Then $\text{tr}(\boldsymbol{A})$ is a random variable with variance
3. If an algorithm computes few queries, it has little information about $\text{tr}(\boldsymbol{A})$
4. Then the algorithm cannot predict $\text{tr}(\boldsymbol{A})$ well

## Super Rough Intuition

$$\mathbf{x} \quad \overset{input}{\Longrightarrow} \quad \text{ORACLE} \quad \overset{output}{\Longrightarrow} \quad \boldsymbol{A}\mathbf{x}$$

View oracle as a limit on information about $\boldsymbol{A}$:

1. Suppose $\boldsymbol{A} \sim \mathcal{D}$ is a random matrix
2. Then $\text{tr}(\boldsymbol{A})$ is a random variable with variance
3. If an algorithm computes few queries, it has little information about $\text{tr}(\boldsymbol{A})$
4. Then the algorithm cannot predict $\text{tr}(\boldsymbol{A})$ well

# Removing the Algorithm's Agency

- ◎ **Problem:** The user can pick many different query vectors $\mathbf{x}$.
- ◎ If the user had no freedom, we could use **statistics** to make lower bounds.

## Removing the Algorithm's Agency

- ◉ **Problem:** The user can pick many different query vectors $\mathbf{x}$.
- ◉ If the user had no freedom, we could use **statistics** to make lower bounds.

Two Observations:

1. WLOG, the user submits orthonormal query vectors

## Removing the Algorithm's Agency

- ⊙ **Problem:** The user can pick many different query vectors $\mathbf{x}$.
- ⊙ If the user had no freedom, we could use **statistics** to make lower bounds.

Two Observations:

1. WLOG, the user submits orthonormal query vectors
2. Let $\boldsymbol{G}$ be a $\mathcal{N}(0, 1)$ Gaussian matrix
   Let $\boldsymbol{Q}$ be an orthogonal matrix
   Then $\boldsymbol{GQ}$ is a $\mathcal{N}(0, 1)$ Gaussian matrix
   - ○ (informal) If $\boldsymbol{A}$ uses Gaussians, the responses from the oracle are independent of the queries submitted.

## Removing the Algorithm's Agency

- ⊚ **Problem:** The user can pick many different query vectors $\mathbf{x}$.
- ⊚ If the user had no freedom, we could use **statistics** to make lower bounds.

Two Observations:

1. WLOG, the user submits orthonormal query vectors
2. Let $\boldsymbol{G}$ be a $\mathcal{N}(0,1)$ Gaussian matrix
   Let $\boldsymbol{Q}$ be an orthogonal matrix
   Then $\boldsymbol{GQ}$ is a $\mathcal{N}(0,1)$ Gaussian matrix
   - ○ (informal) If $\boldsymbol{A}$ uses Gaussians, the responses from the oracle are independent of the queries submitted.

- ⊚ (informal) WLOG, the user observes the first $k$ columns of $\boldsymbol{A}$.

# Wigner/Wishart Anti-Concentration Method

## Theorem (Wishart Case)

- ⊙ Let $\boldsymbol{G} \in \mathbb{R}^{d \times d}$ be a $\mathcal{N}(0, 1)$ Gaussian Matrix.

- ⊙ Let $\boldsymbol{A} = \boldsymbol{G}^\intercal \boldsymbol{G}$ be a *Wishart Matrix*.

- ⊙ An algorithm sends query vectors $\mathbf{x}_1, \ldots, \mathbf{x}_k$, gets responses $\mathbf{w}_1, \ldots, \mathbf{w}_k$

# Wigner/Wishart Anti-Concentration Method

## Theorem (Wishart Case)

⊙ Let $\boldsymbol{G} \in \mathbb{R}^{d \times d}$ be a $\mathcal{N}(0, 1)$ Gaussian Matrix.

⊙ Let $\boldsymbol{A} = \boldsymbol{G}^{\mathsf{T}} \boldsymbol{G}$ be a *Wishart Matrix*.

⊙ An algorithm sends query vectors $\mathbf{x}_1, \ldots, \mathbf{x}_k$, gets responses $\mathbf{w}_1, \ldots, \mathbf{w}_k$

⊙ Then there exists orthogonal matrix $\boldsymbol{V}$ such that

$$\boldsymbol{V} \boldsymbol{A} \boldsymbol{V}^{\mathsf{T}} = \boldsymbol{\Delta} + \begin{bmatrix} 0 & 0 \\ 0 & \tilde{\boldsymbol{A}} \end{bmatrix}$$

where $\tilde{\boldsymbol{A}} \in \mathbb{R}^{(d-k) \times (d-k)}$ is distributed as $\tilde{A} = \tilde{\boldsymbol{G}}^{\mathsf{T}} \tilde{\boldsymbol{G}}$, conditioned on all observations $\mathbf{x}_1, \ldots, \mathbf{x}_k, \mathbf{w}_1, \ldots, \mathbf{w}_k$

⊙ $\boldsymbol{\Delta}$ is known exactly

# Wigner/Wishart Anti-Concentration Method

## Theorem (Wishart Case)

⊙ Let $\boldsymbol{G} \in \mathbb{R}^{d \times d}$ be a $\mathcal{N}(0, 1)$ Gaussian Matrix.

⊙ Let $\boldsymbol{A} = \boldsymbol{G}^\mathsf{T} \boldsymbol{G}$ be a *Wishart Matrix*.

⊙ An algorithm sends query vectors $\mathbf{x}_1, \ldots, \mathbf{x}_k$, gets responses $\mathbf{w}_1, \ldots, \mathbf{w}_k$

⊙ Then there exists orthogonal matrix $\boldsymbol{V}$ such that

$$\boldsymbol{V A V}^\mathsf{T} = \boldsymbol{\Delta} + \begin{bmatrix} 0 & 0 \\ 0 & \tilde{\boldsymbol{A}} \end{bmatrix}$$

where $\tilde{\boldsymbol{A}} \in \mathbb{R}^{(d-k) \times (d-k)}$ is distributed as $\tilde{A} = \tilde{\boldsymbol{G}}^\mathsf{T} \tilde{\boldsymbol{G}}$, conditioned on all observations $\mathbf{x}_1, \ldots, \mathbf{x}_k, \mathbf{w}_1, \ldots, \mathbf{w}_k$

⊙ $\boldsymbol{\Delta}$ is known exactly

⊙ Analogous holds for Wigner Matrices: $\boldsymbol{A} = \frac{1}{2}(\boldsymbol{G} + \boldsymbol{G}^\mathsf{T})$

Consider any adaptive algorithm after $k$ steps:

1. $\operatorname{tr}(\boldsymbol{A}) = \operatorname{tr}(\boldsymbol{V}\boldsymbol{A}\boldsymbol{V}^{\mathsf{T}}) = \operatorname{tr}(\boldsymbol{\Delta}) + \operatorname{tr}(\tilde{\boldsymbol{A}})$

## Wigner/Wishart Anti-Concentration Method

Consider any adaptive algorithm after $k$ steps:

1. $\text{tr}(\boldsymbol{A}) = \text{tr}(\boldsymbol{V}\boldsymbol{A}\boldsymbol{V}^\mathsf{T}) = \text{tr}(\boldsymbol{\Delta}) + \text{tr}(\tilde{\boldsymbol{A}})$
2. Let $t$ estimate $\text{tr}(\boldsymbol{A})$. Define $\tilde{t} := t - \text{tr}(\boldsymbol{\Delta})$.

## Wigner/Wishart Anti-Concentration Method

Consider any adaptive algorithm after $k$ steps:

1. $\mathrm{tr}(\boldsymbol{A}) = \mathrm{tr}(\boldsymbol{V}\boldsymbol{A}\boldsymbol{V}^\mathsf{T}) = \mathrm{tr}(\boldsymbol{\Delta}) + \mathrm{tr}(\tilde{\boldsymbol{A}})$
2. Let $t$ estimate $\mathrm{tr}(\boldsymbol{A})$. Define $\tilde{t} := t - \mathrm{tr}(\boldsymbol{\Delta})$.
3. Note $\mathrm{tr}(\boldsymbol{A}) = \|\boldsymbol{G}\|_F^2 \sim \chi_{d^2}^2$ and $\mathrm{tr}(\tilde{\boldsymbol{A}}) \sim \chi_{(d-k)^2}^2$

Consider any adaptive algorithm after $k$ steps:

1. $\operatorname{tr}(\boldsymbol{A}) = \operatorname{tr}(\boldsymbol{V}\boldsymbol{A}\boldsymbol{V}^{\mathsf{T}}) = \operatorname{tr}(\boldsymbol{\Delta}) + \operatorname{tr}(\tilde{\boldsymbol{A}})$

2. Let $t$ estimate $\operatorname{tr}(\boldsymbol{A})$. Define $\tilde{t} := t - \operatorname{tr}(\boldsymbol{\Delta})$.

3. Note $\operatorname{tr}(\boldsymbol{A}) = \|\boldsymbol{G}\|_F^2 \sim \chi_{d^2}^2$ and $\operatorname{tr}(\tilde{\boldsymbol{A}}) \sim \chi_{(d-k)^2}^2$

   ○ $|t - \operatorname{tr}(\boldsymbol{A})| = |\tilde{t} - \operatorname{tr}(\tilde{\boldsymbol{A}})| \geq \Omega(d - k)$

Consider any adaptive algorithm after $k$ steps:

1. $\text{tr}(\boldsymbol{A}) = \text{tr}(\boldsymbol{V}\boldsymbol{A}\boldsymbol{V}^{\mathsf{T}}) = \text{tr}(\boldsymbol{\Delta}) + \text{tr}(\tilde{\boldsymbol{A}})$

2. Let $t$ estimate $\text{tr}(\boldsymbol{A})$. Define $\tilde{t} := t - \text{tr}(\boldsymbol{\Delta})$.

3. Note $\text{tr}(\boldsymbol{A}) = \|\boldsymbol{G}\|_F^2 \sim \chi_{d^2}^2$ and $\text{tr}(\tilde{\boldsymbol{A}}) \sim \chi_{(d-k)^2}^2$

   ○ $|t - \text{tr}(\boldsymbol{A})| = |\tilde{t} - \text{tr}(\tilde{\boldsymbol{A}})| \geq \Omega(d - k)$

   ○ $\text{tr}(\boldsymbol{A}) \leq O(d^2)$

Consider any adaptive algorithm after $k$ steps:

1. $\text{tr}(\boldsymbol{A}) = \text{tr}(\boldsymbol{V}\boldsymbol{A}\boldsymbol{V}^\mathsf{T}) = \text{tr}(\boldsymbol{\Delta}) + \text{tr}(\tilde{\boldsymbol{A}})$

2. Let $t$ estimate $\text{tr}(\boldsymbol{A})$. Define $\tilde{t} := t - \text{tr}(\boldsymbol{\Delta})$.

3. Note $\text{tr}(\boldsymbol{A}) = \|\boldsymbol{G}\|_F^2 \sim \chi_{d^2}^2$ and $\text{tr}(\tilde{\boldsymbol{A}}) \sim \chi_{(d-k)^2}^2$
   - $|t - \text{tr}(\boldsymbol{A})| = |\tilde{t} - \text{tr}(\tilde{\boldsymbol{A}})| \geq \Omega(d - k)$
   - $\text{tr}(\boldsymbol{A}) \leq O(d^2)$

4. Enforce $|t - \text{tr}(\boldsymbol{A})| \leq \varepsilon \, \text{tr}(\boldsymbol{A})$
$$(d - k) \leq \varepsilon \cdot Cd^2$$

Consider any adaptive algorithm after $k$ steps:

1. $\mathrm{tr}(\boldsymbol{A}) = \mathrm{tr}(\boldsymbol{V}\boldsymbol{A}\boldsymbol{V}^\mathsf{T}) = \mathrm{tr}(\boldsymbol{\Delta}) + \mathrm{tr}(\tilde{\boldsymbol{A}})$

2. Let $t$ estimate $\mathrm{tr}(\boldsymbol{A})$. Define $\tilde{t} := t - \mathrm{tr}(\boldsymbol{\Delta})$.

3. Note $\mathrm{tr}(\boldsymbol{A}) = \|\boldsymbol{G}\|_F^2 \sim \chi_{d^2}^2$ and $\mathrm{tr}(\tilde{\boldsymbol{A}}) \sim \chi_{(d-k)^2}^2$
   - $|t - \mathrm{tr}(\boldsymbol{A})| = |\tilde{t} - \mathrm{tr}(\tilde{\boldsymbol{A}})| \geq \Omega(d-k)$
   - $\mathrm{tr}(\boldsymbol{A}) \leq O(d^2)$

4. Enforce $|t - \mathrm{tr}(\boldsymbol{A})| \leq \varepsilon\,\mathrm{tr}(\boldsymbol{A})$
   $$(d-k) \leq \varepsilon \cdot Cd^2$$

5. Set $d = \frac{1}{2C\varepsilon}$ and simplify: $k \geq \frac{1}{4C\varepsilon}$

<u>Non-Adaptive Proof Framework</u>

Design distributions $\mathcal{P}_0$ and $\mathcal{P}_1$, for large enough $n$:

| $\mathcal{P}_0$ | $\boldsymbol{A} = \boldsymbol{G}^T\boldsymbol{G}$ for $\boldsymbol{G} \in \mathbb{R}^{(\frac{1}{\varepsilon}) \times d}$ Gaussian |
| --- | --- |
| $\mathcal{P}_1$ | $\boldsymbol{A} = \boldsymbol{G}^T\boldsymbol{G}$ for $\boldsymbol{G} \in \mathbb{R}^{(\frac{1}{\varepsilon}+1)\times d}$ Gaussian |

# Statistical Hypothesis Testing

Non-Adaptive Proof Framework

Design distributions $\mathcal{P}_0$ and $\mathcal{P}_1$, for large enough $n$:

| $\mathcal{P}_0$ | $A = G^T G$ for $G \in \mathbb{R}^{(\frac{1}{\varepsilon}) \times d}$ Gaussian |
|---|---|
| $\mathcal{P}_1$ | $A = G^T G$ for $G \in \mathbb{R}^{(\frac{1}{\varepsilon}+1) \times d}$ Gaussian |

1. A trace estimator can distinguish $\mathcal{P}_0$ from $\mathcal{P}_1$
   - If $A_0 \sim \mathcal{P}_0$ and $A_1 \sim \mathcal{P}_1$
   - With high probability, $\operatorname{tr}(A_0) \leq (1 - 2\varepsilon)\operatorname{tr}(A_1)$

Non-Adaptive Proof Framework

Design distributions $\mathcal{P}_0$ and $\mathcal{P}_1$, for large enough $n$:

| $\mathcal{P}_0$ | $\boldsymbol{A} = \boldsymbol{G}^T\boldsymbol{G}$ for $\boldsymbol{G} \in \mathbb{R}^{(\frac{1}{\varepsilon}) \times d}$ Gaussian |
|---|---|
| $\mathcal{P}_1$ | $\boldsymbol{A} = \boldsymbol{G}^T\boldsymbol{G}$ for $\boldsymbol{G} \in \mathbb{R}^{(\frac{1}{\varepsilon}+1) \times d}$ Gaussian |

1. A trace estimator can distinguish $\mathcal{P}_0$ from $\mathcal{P}_1$
   - If $\boldsymbol{A}_0 \sim \mathcal{P}_0$ and $\boldsymbol{A}_1 \sim \mathcal{P}_1$
   - With high probability, $\mathrm{tr}(\boldsymbol{A}_0) \leq (1 - 2\varepsilon)\,\mathrm{tr}(\boldsymbol{A}_1)$
2. No algorithm can distinguish $\mathcal{P}_0$ from $\mathcal{P}_1$ with $\Omega(\frac{1}{\varepsilon})$ queries
   - Nature samples $i \sim \{0, 1\}$, and $\boldsymbol{A} \sim \mathcal{P}_i$
   - User access $\boldsymbol{A}$ through the oracle

Non-Adaptive Proof Framework

Design distributions $\mathcal{P}_0$ and $\mathcal{P}_1$, for large enough $n$:

| $\mathcal{P}_0$ | $\boldsymbol{A} = \boldsymbol{G}^T\boldsymbol{G}$ for $\boldsymbol{G} \in \mathbb{R}^{(\frac{1}{\varepsilon}) \times d}$ Gaussian |
|---|---|
| $\mathcal{P}_1$ | $\boldsymbol{A} = \boldsymbol{G}^T\boldsymbol{G}$ for $\boldsymbol{G} \in \mathbb{R}^{(\frac{1}{\varepsilon}+1)\times d}$ Gaussian |

1. A trace estimator can distinguish $\mathcal{P}_0$ from $\mathcal{P}_1$
   ○ If $\boldsymbol{A}_0 \sim \mathcal{P}_0$ and $\boldsymbol{A}_1 \sim \mathcal{P}_1$
   ○ With high probability, $\text{tr}(\boldsymbol{A}_0) \leq (1 - 2\varepsilon)\,\text{tr}(\boldsymbol{A}_1)$
2. No algorithm can distinguish $\mathcal{P}_0$ from $\mathcal{P}_1$ with $\Omega(\frac{1}{\varepsilon})$ queries
   ○ Nature samples $i \sim \{0, 1\}$, and $\boldsymbol{A} \sim \mathcal{P}_i$
   ○ User access $\boldsymbol{A}$ through the oracle
   ○ WLOG User picks standard basis vectors

Non-Adaptive Proof Framework

Design distributions $\mathcal{P}_0$ and $\mathcal{P}_1$, for large enough $n$:

| $\mathcal{P}_0$ | $\boldsymbol{A} = \boldsymbol{G}^T\boldsymbol{G}$ for $\boldsymbol{G} \in \mathbb{R}^{(\frac{1}{\varepsilon}) \times d}$ Gaussian |
|---|---|
| $\mathcal{P}_1$ | $\boldsymbol{A} = \boldsymbol{G}^T\boldsymbol{G}$ for $\boldsymbol{G} \in \mathbb{R}^{(\frac{1}{\varepsilon}+1)\times d}$ Gaussian |

1. A trace estimator can distinguish $\mathcal{P}_0$ from $\mathcal{P}_1$
   - If $\boldsymbol{A}_0 \sim \mathcal{P}_0$ and $\boldsymbol{A}_1 \sim \mathcal{P}_1$
   - With high probability, $\text{tr}(\boldsymbol{A}_0) \leq (1 - 2\varepsilon)\,\text{tr}(\boldsymbol{A}_1)$
2. No algorithm can distinguish $\mathcal{P}_0$ from $\mathcal{P}_1$ with $\Omega(\frac{1}{\varepsilon})$ queries
   - Nature samples $i \sim \{0, 1\}$, and $\boldsymbol{A} \sim \mathcal{P}_i$
   - User access $\boldsymbol{A}$ through the oracle
   - WLOG User picks standard basis vectors
   - Bound Total Variation between first $k$ columns of $\boldsymbol{A}_0$ and $\boldsymbol{A}_1$

1. Introduced Hutchinson's Estimator for PSD $A$
2. Improved it: Hutch++ uses $O(\frac{1}{\varepsilon})$
3. Two lower bounds: Adaptive & Non-Adaptive require $\Omega(\frac{1}{\varepsilon})$
4. Trace Estimation requires $\Theta(\frac{1}{\varepsilon})$ queries

## Open Questions

- ◉ When is adaptivity helpful?

- ◉ What about inexact oracles? We often approximate $f(\boldsymbol{A})\mathbf{x}$ with iterative methods. How accurate do these computations need to be?

- ◉ Extend to include row/column sampling? This would encapsulate e.g. SGD/SCD.

- ◉ Memory-limited lower bounds? This is a realistic model for iterative methods.

THANK YOU

📄 Haim Avron and Sivan Toledo.
Randomized algorithms for estimating the trace of an implicit symmetric positive semi-definite matrix.
Journal of the ACM, 58(2), 2011.

📄 Tamas Sarlos.
Improved approximation algorithms for large matrices via random projections.
In Proceedings of the 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS), pages 143–152, 2006.

📄 David P. Woodruff.
Sketching as a tool for numerical linear algebra.
Foundations and Trends in Theoretical Computer Science, 10(1–2):1–157, 2014.

Karl Wimmer, Yi Wu, and Peng Zhang.
Optimal query complexity for estimating the trace of a matrix.
In Proceedings of the 41st International Colloquium on
Automata, Languages and Programming (ICALP), pages
1051–1062, 2014.