

# COMPSCI 614: Randomized Algorithms with Applications to Data Science

---

Prof. Cameron Musco

University of Massachusetts Amherst. Spring 2024.

Lecture 12

- The midterm is the Thursday after break in class.
- I will hold a review session Monday from 3-4:30pm and Tuesday in class.
- There is no real quiz this week, but see Weekly Quizzes section on Moodle for a single question quiz where you can mark if you attended Sally Dong's job talk for extra credit.
- Practice midterm in moodle.
  - ↳ Solutions
  - ↳ study guide

## Last Time:

- Finish up fast low-rank approximation via approximate matrix multiplication.
- Start on stochastic trace estimation and motivation for matrix-vector query algorithms.

# Summary

## Last Time:

- Finish up fast low-rank approximation via approximate matrix multiplication.
- Start on stochastic trace estimation and motivation for matrix-vector query algorithms.

## Today:

- Finish stochastic trace estimation.
- Hutchinson's estimator and full analysis.

• ~~K-means ++~~

# Matrix Trace

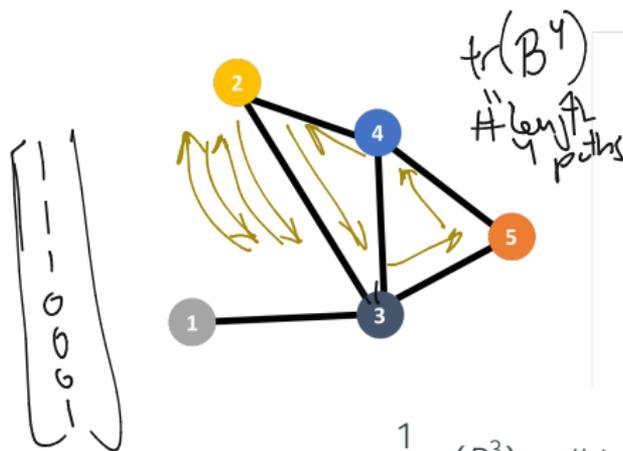
The trace of a matrix  $A \in \mathbb{R}^{n \times n}$  is the sum of its diagonal entries.

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

When  $A$  is diagonalizable (e.g., when it is symmetric) with eigenvalues  $\lambda_1, \dots, \lambda_n$ ,  $\text{tr}(A) = \sum_{i=1}^n \lambda_i$ .

**Main question:** How many matrix-vector multiplications “queries”  $Ax_1, \dots, Ax_m$  are required to approximate  $\text{tr}(A)$ ?

# Motivating Example



$\text{tr}(B^4)$   
# length 4 paths

$B^3$				
<u>0</u>	1	4	②	1
1	<u>2</u>	6	5	2
4	⑥	<u>4</u>	⑥	6
②	5	6	<u>4</u>	5
1	2	6	5	<u>2</u>

sparse matrix approx

$$\frac{1}{6} \text{tr}(B^3) = \# \text{ triangles.}$$

$B(B(Bx))$        $\text{tr}(B^3) = \#$  length 3 paths starting at vertex  $v$  and ending at vertex  $v$  (not counting  $v$ )

- Explicitly forming  $B^3$  and computing  $\text{tr}(B^3)$  takes  $O(n^3)$  time.
- Can multiply  $B^3$  by a vector in  $3 \cdot |E| = O(n^2)$  operations.
- So a trace estimation algorithm using  $m$  queries, yields an  $O(m \cdot |E|)$  time approximate triangle counting algorithm.

## Other Examples

**Example 2:** Hessian/Jacobian matrix-vector products.

*$n \times n$  Jacobian /  $n \times n$  Hessian*

- For vector  $x$ ,  $\nabla f(y)x$  and  $\nabla^2 f(y)x$  can often be computed efficiently using finite difference methods or explicit differentiation (e.g., via backpropagation).
- Do not need to fully form  $\nabla f(y)$  or  $\nabla^2 f(y)$ .
- Many applications of estimating the traces of these matrices, e.g., in analyzing neural network convergence, in optimization of score-based methods, etc.
- $\text{tr}(\nabla^2 f(y)x)$ : Laplacian
- $\text{tr}(\nabla f(y)x)$ : Divergence

## Other Examples

**Example 3:**  $A$  is a function of another (explicit) matrix  $B$ ,  $A = f(B)$  that can be applied efficiently via an iterative method.

## Other Examples

$$\Sigma^{1/2} \cdot \Sigma^{1/2} = \Sigma$$

**Example 3:**  $A$  is a function of another (explicit) matrix  $B$ ,  $A = f(B)$  that can be applied efficiently via an iterative method.

$$x \sim N(0, \Sigma^{n \times n})$$

$\Sigma^{1/2} g$  - iid Gaussian

$$\Sigma^{1/2} E[g g^T] \Sigma^{1/2}$$

- Repeated multiplication to apply  $A = B^3$ .
- Conjugate gradient, MINRES, or any linear system solver:

$$A = B^{-1}.$$

- Lanczos method, polynomial/rational approximation:

$$A = \exp(B), A = \sqrt{B}, A = \log(B), \text{ etc.}$$

- These methods run in  $n^2 \cdot C$  time, where  $C$  depends on properties of  $B$ . Typically  $C \ll n$  so  $n^2 \cdot C \ll n^3$ .

## Matrix Function Examples

- Log-likelihood computation in Bayesian optimization, experimental design.  $\text{tr}(\log(B)) = \log \det(B) = \sum \log(\lambda_i)$
- Estrada index, a measure of protein folding degree and more generally, network connectivity.  $\text{tr}(\exp(B))$ .
- Trace inverse, which is important in uncertainty quantification and many other scientific computing applications.  $\text{tr}(B^{-1})$
- Information about the matrix eigenvalue spectrum, since  $\text{tr}(f(B)) = \sum_{i=1}^n f(\lambda_i)$ , where  $\lambda_i$  is  $B$ 's  $i^{\text{th}}$  eigenvalue.
- E.g., counting the number of eigenvalues in an interval, spectral density estimation, matrix norms
- See e.g., [Ubaru, and Saad 2017].

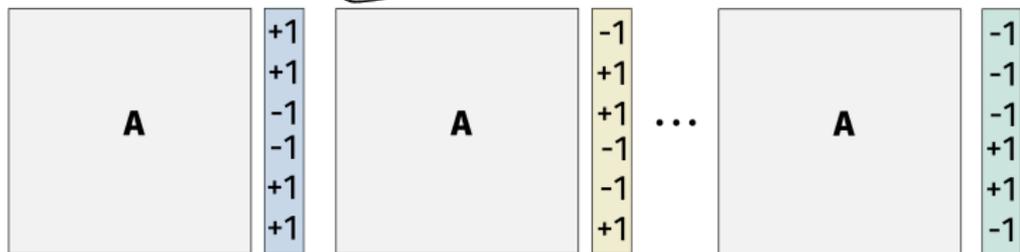


# Hutchinson's Method

Hutchinson 1991, Girard 1987:

$$E[-1 \ -1 \ 1] \left\{ \begin{matrix} A \\ \vdots \\ \vdots \end{matrix} \right\}$$

- Draw  $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$  i.i.d. with random  $\{+1, -1\}$  entries.
- Return  $\bar{\mathbf{T}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^T \mathbf{A} \mathbf{x}_i$  as an approximation to  $\text{tr}(\mathbf{A})$ .



- One of the earliest examples I know of a randomized algorithm for linear algebraic computation.

# Hutchinson's Method Error Bound

## Theorem

Let  $\bar{\mathbf{T}}$  be the trace estimate returned by Hutchinson's method. If  $m = O\left(\frac{1}{\delta\epsilon^2}\right)$ , then with probability  $\geq 1 - \delta$ ,

$$|\bar{\mathbf{T}} - \text{tr}(A)| \leq \epsilon \|A\|_F$$

# Hutchinson's Method Error Bound

## Theorem

Let  $\bar{T}$  be the trace estimate returned by Hutchinson's method.

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

$$|\bar{T} - \text{tr}(A)| \leq \epsilon \|A\|_F$$

*non-negative eigenvalues*

If  $A$  is symmetric positive semidefinite (PSD) then

$$\|A\|_F = \sqrt{\sum_{i=1}^n \lambda_i^2} \leq \sum_{i=1}^n \lambda_i = \text{tr}(A).$$

So for PSD  $A$ :  $(1 - \epsilon) \text{tr}(A) \leq \bar{T} \leq (1 + \epsilon) \text{tr}(A)$ .

## Theorem

Let  $\bar{T}$  be the trace estimate returned by Hutchinson's method. If  $m = O\left(\frac{1}{\delta\epsilon^2}\right)$ , then with probability  $\geq 1 - \delta$ ,

$$|\bar{T} - \text{tr}(A)| \leq \epsilon \|A\|_F$$

1. Show that  $\mathbb{E}[\bar{T}] = \text{tr}(A)$ .
2. Bound  $\text{Var}[\bar{T}]$ .
3. Apply Chebyshev's inequality.

# Proof Approach

## Theorem

Let  $\bar{T}$  be the trace estimate returned by Hutchinson's method. If  $m = O\left(\frac{1}{\delta\epsilon^2}\right)$ , then with probability  $\geq 1 - \delta$ ,

$$\underline{\underline{|\bar{T} - \text{tr}(A)|}} \leq \epsilon \|A\|_F$$

1. Show that  $\mathbb{E}[\bar{T}] = \text{tr}(A)$ .
2. Bound  $\text{Var}[\bar{T}]$ .
3. Apply Chebyshev's inequality.

$$\mathbb{E}[\bar{T}] = \frac{1}{m} \mathbb{E}[\text{tr}(A)] = \text{tr}(A)$$

A tighter proof that uses the **Hanson-Wright inequality**, an exponential concentration inequality for quadratic forms, can improve the  $\delta$  dependence to  $\log(1/\delta)$  – we'll cover this later in the class.

# Expectation Analysis

Hutchinson's Estimator::

- Draw  $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$  i.i.d. with random  $\{+1, -1\}$  entries.
- Return  $\bar{\mathbf{T}} = \frac{1}{m} \sum_{i=1}^m \underline{\mathbf{x}_i^T \mathbf{A} \mathbf{x}_i}$  as an approximation to  $\text{tr}(\mathbf{A})$ .

---

By linearity of expectation,  $\mathbb{E}[\bar{\mathbf{T}}] = \mathbb{E}[\mathbf{x}^T \mathbf{A} \mathbf{x}]$  for a single random  $\pm 1$  vector  $\mathbf{x}$ .  $\rightarrow \text{tr}(\mathbf{A})$

# Expectation Analysis

Hutchinson's Estimator::

- Draw  $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$  i.i.d. with random  $\{+1, -1\}$  entries.
  - Return  $\bar{\mathbf{T}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^T \mathbf{A} \mathbf{x}_i$  as an approximation to  $\text{tr}(\mathbf{A})$ .
- 

By linearity of expectation,  $\mathbb{E}[\bar{\mathbf{T}}] = \mathbb{E}[\mathbf{x}^T \mathbf{A} \mathbf{x}]$  for a single random  $\pm 1$  vector  $\mathbf{x}$ .

$$\mathbb{E}[\mathbf{x}^T \mathbf{A} \mathbf{x}] = \mathbb{E} \left[ \sum_{i=1}^n \sum_{j=1}^n \mathbf{x}_i \mathbf{x}_j A_{ij} \right] = \sum_{i=1}^n \sum_{j=1}^n A_{ij} \cdot \mathbb{E}[\mathbf{x}_i \mathbf{x}_j]$$

$\downarrow$   
0 if  $i \neq j$   
1 if  $i = j$

# Expectation Analysis

Hutchinson's Estimator::

- Draw  $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$  i.i.d. with random  $\{+1, -1\}$  entries.
  - Return  $\bar{\mathbf{T}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^T \mathbf{A} \mathbf{x}_i$  as an approximation to  $\text{tr}(\mathbf{A})$ .
- 

By linearity of expectation,  $\mathbb{E}[\bar{\mathbf{T}}] = \mathbb{E}[\mathbf{x}^T \mathbf{A} \mathbf{x}]$  for a single random  $\pm 1$  vector  $\mathbf{x}$ .

$$\mathbb{E}[\mathbf{x}^T \mathbf{A} \mathbf{x}] = \mathbb{E} \sum_{i=1}^n \sum_{j=1}^n \mathbf{x}_i \mathbf{x}_j A_{ij} = \sum_{i=1}^n \sum_{j=1}^n A_{ij} \cdot \mathbb{E}[\mathbf{x}_i \mathbf{x}_j]$$

- When  $i \neq j$ ,  $\mathbf{x}_i \mathbf{x}_j = 1$  with probability  $1/2$  and  $-1$  with probability  $1/2$ , so  $\mathbb{E}[\mathbf{x}_i \mathbf{x}_j] = 0$ . When  $i = j$ ,  $\mathbf{x}_i \mathbf{x}_j = 1$ , so  $\mathbb{E}[\mathbf{x}_i \mathbf{x}_j] = 1$ .

# Expectation Analysis

Hutchinson's Estimator::

- Draw  $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$  i.i.d. with random  $\{+1, -1\}$  entries.
  - Return  $\bar{\mathbf{T}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^T \mathbf{A} \mathbf{x}_i$  as an approximation to  $\text{tr}(\mathbf{A})$ .
- 

By linearity of expectation,  $\mathbb{E}[\bar{\mathbf{T}}] = \mathbb{E}[\mathbf{x}^T \mathbf{A} \mathbf{x}]$  for a single random  $\pm 1$  vector  $\mathbf{x}$ .

$$\mathbb{E}[\mathbf{x}^T \mathbf{A} \mathbf{x}] = \mathbb{E} \sum_{i=1}^n \sum_{j=1}^n \mathbf{x}_i \mathbf{x}_j A_{ij} = \sum_{i=1}^n \sum_{j=1}^n A_{ij} \cdot \mathbb{E}[\mathbf{x}_i \mathbf{x}_j] = \sum_{i=1}^n A_{ii}.$$

- When  $i \neq j$ ,  $\mathbf{x}_i \mathbf{x}_j = 1$  with probability  $1/2$  and  $-1$  with probability  $1/2$ , so  $\mathbb{E}[\mathbf{x}_i \mathbf{x}_j] = 0$ . When  $i = j$ ,  $\mathbf{x}_i \mathbf{x}_j = 1$ , so  $\mathbb{E}[\mathbf{x}_i \mathbf{x}_j] = 1$ .

# Expectation Analysis

Hutchinson's Estimator::  $(XX^T)_{ij} = x_i x_j = 1$  if  $i=j$  only

- Draw  $x_1, \dots, x_m \in \mathbb{R}^n$  i.i.d. with random  $\{+1, -1\}$  entries.
- Return  $\bar{T} = \frac{1}{m} \sum_{i=1}^m x_i^T A x_i$  as an approximation to  $\text{tr}(A)$ .

$$\mathbb{E} \text{tr}(X^T A X) = \mathbb{E} \text{tr}(X X^T A) = \text{tr}(\mathbb{E} X X^T A) = \text{tr}(A)$$

---

By linearity of expectation,  $\mathbb{E}[\bar{T}] = \mathbb{E}[x^T A x]$  for a single random  $\pm 1$  vector  $x$ .

$$\mathbb{E}[x^T A x] = \mathbb{E} \sum_{i=1}^n \sum_{j=1}^n x_i x_j A_{ij} = \sum_{i=1}^n \sum_{j=1}^n A_{ij} \cdot \mathbb{E}[x_i x_j] = \sum_{i=1}^n A_{ii} = \text{tr}(A)$$

- When  $i \neq j$ ,  $x_i x_j = 1$  with probability  $1/2$  and  $-1$  with probability  $1/2$ , so  $\mathbb{E}[x_i x_j] = 0$ . When  $i = j$ ,  $x_i x_j = 1$ , so  $\mathbb{E}[x_i x_j] = 1$ .
- So the estimator is correct in expectation:  $\mathbb{E}[\bar{T}] = \text{tr}(A)$ .

# Variance Bound

Hutchinson's Estimator::

- Draw  $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$  i.i.d. with random  $\{+1, -1\}$  entries.
- Return  $\bar{\mathbf{T}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^T \mathbf{A} \mathbf{x}_i$  as an approximation to  $\text{tr}(\mathbf{A})$ .

problem on  
ppt  $\rightarrow$  loads  
a lot like  
tw.

$\text{Var}[\bar{\mathbf{T}}]$

# Variance Bound

Hutchinson's Estimator::

- Draw  $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$  i.i.d. with random  $\{+1, -1\}$  entries.
  - Return  $\bar{\mathbf{T}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^T \mathbf{A} \mathbf{x}_i$  as an approximation to  $\text{tr}(\mathbf{A})$ .
- 

$$\text{Var}[\bar{\mathbf{T}}] = \frac{1}{m} \text{Var}[\mathbf{x}^T \mathbf{A} \mathbf{x}]$$

# Variance Bound

Hutchinson's Estimator::

- Draw  $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$  i.i.d. with random  $\{+1, -1\}$  entries.
- Return  $\bar{\mathbf{T}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^T \mathbf{A} \mathbf{x}_i$  as an approximation to  $\text{tr}(\mathbf{A})$ .

---

$$\text{Var}[\bar{\mathbf{T}}] = \frac{1}{m} \text{Var}[\mathbf{x}^T \mathbf{A} \mathbf{x}] = \frac{1}{m} \text{Var} \left[ \sum_{i=1}^n \sum_{j=1}^n x_i x_j A_{ij} \right]$$

# Variance Bound

Hutchinson's Estimator::

- Draw  $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$  i.i.d. with random  $\{+1, -1\}$  entries.
- Return  $\bar{\mathbf{T}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^T \mathbf{A} \mathbf{x}_i$  as an approximation to  $\text{tr}(\mathbf{A})$ .

---

$$\text{Var}[\bar{\mathbf{T}}] = \frac{1}{m} \text{Var}[\mathbf{x}^T \mathbf{A} \mathbf{x}] = \frac{1}{m} \text{Var} \left[ \sum_{i=1}^n \sum_{j=1}^n x_i x_j A_{ij} \right]$$

Can we apply linearity of variance here?

Handwritten examples of terms from the variance formula:

- $x_1 x_1 A_{11}$
- $x_1 x_2 A_{12}$
- $x_2 x_1 A_{21}$

# Variance Bound

Hutchinson's Estimator::

- Draw  $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$  i.i.d. with random  $\{+1, -1\}$  entries.
- Return  $\bar{\mathbf{T}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^T \mathbf{A} \mathbf{x}_i$  as an approximation to  $\text{tr}(\mathbf{A})$ .

---

$$\text{Var}[\bar{\mathbf{T}}] = \frac{1}{m} \text{Var}[\mathbf{x}^T \mathbf{A} \mathbf{x}] = \frac{1}{m} \text{Var} \left[ \sum_{i=1}^n \sum_{j=1}^n x_i x_j A_{ij} \right]$$

Can we apply linearity of variance here? Almost – need to remove repeated terms, and then can use pairwise independence.

# Variance Bound

Hutchinson's Estimator::

- Draw  $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$  i.i.d. with random  $\{+1, -1\}$  entries.
- Return  $\bar{\mathbf{T}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^T \mathbf{A} \mathbf{x}_i$  as an approximation to  $\text{tr}(\mathbf{A})$ .

---

$$\text{Var}[\bar{\mathbf{T}}] = \frac{1}{m} \text{Var}[\mathbf{x}^T \mathbf{A} \mathbf{x}] = \frac{1}{m} \text{Var} \left[ \sum_{i=1}^n \sum_{j=1}^n x_i x_j A_{ij} \right]$$

*n<sup>2</sup> terms*  
*"reorganized"*

Can we apply linearity of variance here? Almost - need to remove repeated terms, and then can use pairwise independence.

$$\text{Var}[\bar{\mathbf{T}}] = \frac{1}{m} \text{Var} \left[ \sum_{i=1}^n A_{ii} + \sum_{i=1}^n \sum_{j>i} x_i x_j (A_{ij} + A_{ji}) \right]$$

*x<sub>1</sub>x<sub>2</sub> (A<sub>12</sub> + A<sub>21</sub>)*

# Variance Bound

Hutchinson's Estimator::

- Draw  $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$  i.i.d. with random  $\{+1, -1\}$  entries.
- Return  $\bar{\mathbf{T}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^T A \mathbf{x}_i$  as an approximation to  $\text{tr}(A)$ .

---

$$\text{Var}[\bar{\mathbf{T}}] = \frac{1}{m} \text{Var}[\mathbf{x}^T A \mathbf{x}] = \frac{1}{m} \text{Var} \left[ \sum_{i=1}^n \sum_{j=1}^n \mathbf{x}_i \mathbf{x}_j A_{ij} \right]$$

Can we apply linearity of variance here? Almost – need to remove repeated terms, and then can use pairwise independence.

$$\begin{aligned} \text{Var}[\bar{\mathbf{T}}] &= \frac{1}{m} \text{Var} \left[ \sum_{i=1}^n A_{ii} + \sum_{i=1}^n \sum_{j>i} \mathbf{x}_i \mathbf{x}_j (A_{ij} + A_{ji}) \right] \\ &= \frac{1}{m} \sum_{i=1}^n \sum_{j>i} \underbrace{\text{Var}[\mathbf{x}_i \mathbf{x}_j]}_{\substack{1 \\ \text{if } i \neq j}} \cdot \underbrace{(A_{ij} + A_{ji})^2}_{\substack{2 \\ \text{if } i \neq j}} \end{aligned}$$

# Variance Bound

Hutchinson's Estimator::

- Draw  $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$  i.i.d. with random  $\{+1, -1\}$  entries.
- Return  $\bar{T} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^T A \mathbf{x}_i$  as an approximation to  $\text{tr}(A)$ .

$$\begin{aligned} & (A-B)^2 \\ &= A^2 + B^2 - 2AB \\ & 2AB \leq A^2 + B^2 \end{aligned}$$

$$\text{Var}[\bar{T}] = \frac{1}{m} \text{Var}[\mathbf{x}^T A \mathbf{x}] = \frac{1}{m} \text{Var} \left[ \sum_{i=1}^n \sum_{j=1}^n x_i x_j A_{ij} \right]$$

Can we apply linearity of variance here? Almost - need to remove repeated terms, and then can use pairwise independence.

$$\begin{aligned} \text{Var}[\bar{T}] &= \frac{1}{m} \text{Var} \left[ \sum_{i=1}^n A_{ii} + \sum_{i=1}^n \sum_{j>i} x_i x_j (A_{ij} + A_{ji}) \right] \\ &= \frac{1}{m} \sum_{i=1}^n \sum_{j>i} \text{Var}[x_i x_j] \cdot (A_{ij} + A_{ji})^2 \leq \frac{1}{m} \sum_{i=1}^n \sum_{j>i} 2A_{ij}^2 + 2A_{ji}^2 \\ & \quad \underbrace{A_{ij}^2 + A_{ji}^2 + 2A_{ij}A_{ji}}_{\leq 2 \cdot \frac{A_{ij}^2 + A_{ji}^2}{2}} \end{aligned}$$

"arithmetic  
geometric  
mean  
inequality"  
 $\sqrt{AB} \leq \frac{A+B}{2}$

# Variance Bound

Hutchinson's Estimator::

- Draw  $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$  i.i.d. with random  $\{+1, -1\}$  entries.
- Return  $\bar{\mathbf{T}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^T \mathbf{A} \mathbf{x}_i$  as an approximation to  $\text{tr}(\mathbf{A})$ .

$$\begin{array}{c|c} \mathbf{A}_1 & \vdots \\ \hline \vdots & \mathbf{A}_2 \end{array}$$

$$\text{Var}[\bar{\mathbf{T}}] = \frac{1}{m} \text{Var}[\mathbf{x}^T \mathbf{A} \mathbf{x}] = \frac{1}{m} \text{Var} \left[ \sum_{i=1}^n \sum_{j=1}^n \mathbf{x}_i \mathbf{x}_j \mathbf{A}_{ij} \right]$$

Can we apply linearity of variance here? Almost - need to remove repeated terms, and then can use pairwise independence.

$$\begin{aligned} \boxed{\text{Var}[\bar{\mathbf{T}}]} &= \frac{1}{m} \text{Var} \left[ \sum_{i=1}^n \mathbf{A}_{ii} + \sum_{i=1}^n \sum_{j>i} \mathbf{x}_i \mathbf{x}_j (\mathbf{A}_{ij} + \mathbf{A}_{ji}) \right] \\ &= \frac{1}{m} \sum_{i=1}^n \sum_{j>i} \text{Var}[\mathbf{x}_i \mathbf{x}_j] \cdot (\mathbf{A}_{ij} + \mathbf{A}_{ji})^2 \leq \frac{1}{m} \sum_{i=1}^n \sum_{j>i} 2\mathbf{A}_{ij}^2 + 2\mathbf{A}_{ji}^2 \\ &\leq \frac{1}{m} \cdot 2 \cdot \sum_{i=1}^n \sum_{j=1}^n \mathbf{A}_{ij}^2 \end{aligned}$$

$\frac{2 \|\mathbf{A}\|_F^2}{m}$

$\leq \frac{2 \|\mathbf{A}\|_F^2}{m}$

Hutchinson's Estimator::

- Draw  $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$  i.i.d. with random  $\{+1, -1\}$  entries.
  - Return  $\bar{\mathbf{T}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^T \mathbf{A} \mathbf{x}_i$  as an approximation to  $\text{tr}(\mathbf{A})$ .
- 

Chebyshev's inequality implies that, for  $m = \frac{2}{\delta \epsilon^2}$ :

$$\Pr \left[ \underbrace{|\bar{\mathbf{T}} - \text{tr}(\mathbf{A})| \geq \epsilon \|\mathbf{A}\|_F} \right] \leq \frac{2 \|\mathbf{A}\|_F^2 / m}{\epsilon^2 \|\mathbf{A}\|_F^2} = \delta.$$
$$\frac{2}{m \epsilon^2}$$

# Final Analysis

Hutchinson's Estimator::

- Draw  $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$  i.i.d. with random  $\{+1, -1\}$  entries.
- Return  $\bar{\mathbf{T}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^T A \mathbf{x}_i$  as an approximation to  $\text{tr}(A)$ .

Chebyshev's inequality implies that, for  $m = \frac{2}{\delta \epsilon^2}$ :

$$\mathbf{x}^T A \mathbf{x}$$

$$\Pr [|\bar{\mathbf{T}} - \text{tr}(A)| \geq \epsilon \|A\|_F] \leq \frac{2\|A\|_F^2/m}{\epsilon^2 \|A\|_F^2} = \delta.$$
$$x_1 x_2 = 1 \quad x_2 x_3 = -1 \Rightarrow x_1 x_3 = -1$$

Could we have gotten a better bound by applying Bernstein's inequality to  $\sum_{i=1}^n \sum_{j>i} \mathbf{x}_i \mathbf{x}_j (A_{ij} + A_{ji})$ ?

- upper bounds on  $A_{ij} + A_{ji}$ ? ✓
- pairwise independent ✗

$$\mathbf{x}^T A \mathbf{x}$$

# Final Analysis

Hutchinson's Estimator::

- Draw  $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$  i.i.d. with random  $\{+1, -1\}$  entries.
  - Return  $\bar{\mathbf{T}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^T \mathbf{A} \mathbf{x}_i$  as an approximation to  $\text{tr}(\mathbf{A})$ .
- 

Chebyshev's inequality implies that, for  $m = \frac{2}{\delta \epsilon^2}$ :

$$\Pr [|\bar{\mathbf{T}} - \text{tr}(\mathbf{A})| \geq \epsilon \|\mathbf{A}\|_F] \leq \frac{2\|\mathbf{A}\|_F^2/m}{\epsilon^2 \|\mathbf{A}\|_F^2} = \delta.$$

Could we have gotten a better bound by applying Bernstein's inequality to  $\sum_{i=1}^n \sum_{j>i} \mathbf{x}_i \mathbf{x}_j (A_{ij} + A_{ji})$ ?

Hanson-Wright is an exponential concentration bound that can be used in the specific case – improves bound to  $m = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ .

# Optimality of Hutchinson's Method

The  $m = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$  bound given by the Hanson-Wright inequality is tight.

- Any algorithm that only uses queries of the form  $\mathbf{x}_i^T \mathbf{A} \mathbf{x}_i$  requires  $\Omega\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$  samples to estimate  $\text{tr}(\mathbf{A})$  to error  $\pm \epsilon \text{tr}(\mathbf{A})$  for PSD  $\mathbf{A}$  [Wimmer, Wu, Zhang 2014].

↳  $\epsilon \|\mathbf{A}\|_F$

# Optimality of Hutchinson's Method

The  $m = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$  bound given by the Hanson-Wright inequality is tight.

- Any algorithm that only uses queries of the form  $\underline{x_i^T A x_i}$  requires  $\Omega\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$  samples to estimate  $\text{tr}(A)$  to error  $\pm\epsilon \text{tr}(A)$  for PSD  $A$  [Wimmer, Wu, Zhang 2014].  $\underline{A x_1 \dots A x_m}$
- We recently showed that using the full power of matrix-vector queries, one can achieve  $O\left(\frac{\log(1/\delta)}{\epsilon}\right)$  queries for PSD matrices.

(Hutch + Meyer et al.)