

# Trace Estimation



**Raphael Meyer**, Caltech

UMich EECS 598-004 Randomized Numerical Linear Algebra in Machine Learning

## Motivation 1: Road Network Connectivity



NYC is adding a new bus line. Where should it go?

**Idea:** Maximize *connectivity* of transit network using existing bus stops

**Given:** Adjacency matrix  $\mathbf{B} \in \mathbb{R}^{n \times n}$  in memory

**Algorithm:** For each possible new bus route (edge), build new  $\mathbf{B}' \in \mathbb{R}^{n \times n}$ , and compute the connectivity of  $\mathbf{B}'$ .

**Return:** Edge that maximized connectivity

**Bottleneck:** Computing connectivity

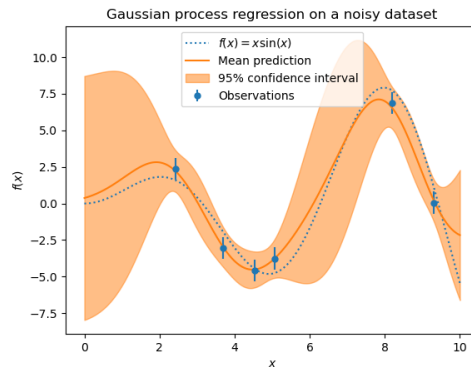
## Computational Bottleneck

- ⊙ We have matrix  $\mathbf{B} \in \mathbb{R}^{n \times n}$  in memory
- ⊙ We want to compute *connectivity* of  $\mathbf{B}$ :
  - Estrada Index:  $\text{tr}(e^{\mathbf{B}})$
  - Num of Triangles:  $\text{tr}(\frac{1}{6}\mathbf{B}^3)$
- ⊙ The **trace** is the sum of the diagonal of a matrix.
- ⊙ Computing  $\mathbf{B}^3$  takes  $O(n^3)$  time. **slow**
- ⊙ Computing  $\mathbf{B}^3 x = \mathbf{B}(\mathbf{B}(\mathbf{B}x))$  takes  $O(n^2)$  time. **fast**

*Can we approximate  $\text{tr}(\mathbf{B}^3)$  by computing few  
 $\mathbf{B}^3 x_1, \dots, \mathbf{B}^3 x_m$ ?*

- ⊙ **Yes, we can!**

## Motivation 2: Gaussian Process Optimization



We have time series data, and want to interpolate it with kernel  $k_\theta$

**Idea:** Maximize log-likelihood kernel matrix  $K_\theta$  over hyperparameter  $\theta$

**Given:** Data points  $(x_1, y_1), \dots, (x_n, y_n)$ , kernel function  $k_\theta(x, x')$

**Algorithm:** Compute Gradient Descent of Log-Likelihood to optimize  $\theta$

**Return:** Optimal  $\theta$

**Bottleneck:** Computing gradients

$$\nabla_\theta \mathcal{L}(\theta) = \frac{1}{2} y^T K^{-1} \frac{dK}{d\theta} K^{-1} y - \frac{1}{2} \text{tr} \left( K^{-1} \frac{dK}{d\theta} \right)$$

Caltech

## General Picture: Trace Estimation

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

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

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

Num. Triangles	Estrada Index	Log-Determinant
$\text{tr}(\frac{1}{6}\mathbf{B}^3)$	$\text{tr}(e^{\mathbf{B}})$	$\text{tr}(\ln(\mathbf{B}))$

- If  $\mathbf{A} = f(\mathbf{B})$ , then we can often compute  $\mathbf{A}x$  quickly  
Compute  $\mathbf{B}^q x$  in  $\tilde{O}(n^2 \sqrt{q})$  time  
Lanczos-FA: compute  $\mathbf{A}x$  in  $\tilde{O}(n^2 \sqrt{\kappa(\mathbf{B})})$  time for many  $f$
- Goal: Estimate  $\text{tr}(\mathbf{A})$  by computing  $\mathbf{A}x_1, \dots, \mathbf{A}x_m$

# Matrix-Vector Oracle Model

Think of the Matrix-Vector Product as a Computational Primitive

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

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

- ⊙ e.g. Randomized SVD, Johnson-Lindenstrauss, Power Method, Lanczos-FA
- ⊙ Some existing lower bounds; very active area of research!

**Trace Estimation:** Estimate  $\text{tr}(A)$  with as few Matrix-Vector products  $Ax_1, \dots, Ax_k$  as possible.

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

# Outline

- ⊙ Part 1: The Girard-Hutchinson Trace Estimator
  - ⊙ 2 Lines of MATLAB Code
- ⊙ Part 2: The Hutch++ Trace Estimation
  - ⊙ 5 Lines of MATLAB Code
- ⊙ Part 3: Extensions of Hutch++
  - ⊙ 3 Follow-up papers, at a very high level
- ⊙ Part 4: Lower Bounds
  - ⊙ No algorithm can get better big-Oh than Hutch++!

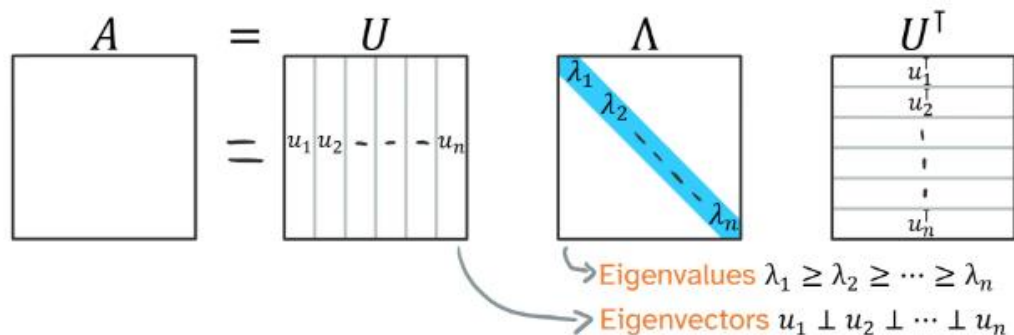
**Part 1: Girard-Hutchinson**

**Part 2: Hutch++**



# Linear Algebra Notation

- Let  $A \in \mathbb{R}^{n \times n}$  be symmetric. Then,  $A = U\Lambda U^T$



- $A$  is **Positive Semi-Definite (PSD)** if all  $\lambda_i \geq 0$

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

- Frobenius Norm  $\|A\|_F^2 = \sum_{i=1}^n A_{ii}^2 = \sum_{i=1}^n \lambda_i^2$

- For PSD matrices,  $\|A\|_F \leq tr(A)$

## Girard-Hutchinson Trace Estimator

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

$$\mathbb{E}[x^T \mathbf{A} x] = \text{tr}(\mathbf{A})$$

$$\text{Var}[x^T \mathbf{A} x] = 2\|\mathbf{A}\|_F^2$$

- ⊙ Girard-Hutchinson Estimator:  $H_\ell(\mathbf{A}) := \frac{1}{\ell} \sum_{i=1}^{\ell} x_i^T \mathbf{A} x_i$

$$\mathbb{E}[H_\ell(\mathbf{A})] = \text{tr}(\mathbf{A})$$

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

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

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

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

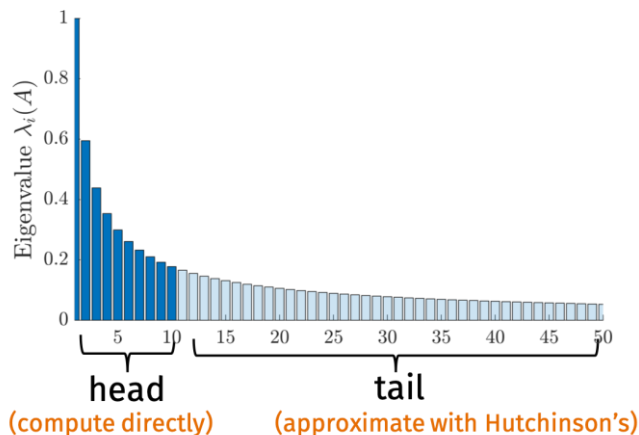
## Girard-Hutchinson Trace Estimator

- ⊙ When is this analysis tight?

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

- ⊙ When is the bound  $\|\mathbf{A}\|_F \leq \text{tr}(\mathbf{A})$  tight?
- ⊙ Let  $\mathbf{v} = [\lambda_1 \dots \lambda_n]$  be eigenvalues of  $\mathbf{A}$
- ⊙ When is  $\|\mathbf{v}\|_2 \leq \|\mathbf{v}\|_1$  tight?
  - ⊙ Property of norms:  $\|\mathbf{v}\|_2 \approx \|\mathbf{v}\|_1$  only if  $\mathbf{v}$  is nearly sparse
- ⊙ The estimator only needs  $O(\frac{1}{\varepsilon^2})$  when  $\mathbf{A}$  has few large eigenvalues!

# Helping Hutchinson's Estimator



⊙ Idea: Explicitly estimate top few eigenvals of  $\mathbf{A}$ . Use Hutchinson for rest.

1. Find a good rank- $k$  approximation  $\widetilde{\mathbf{A}}_k$
2. Notice that  $\text{tr}(\mathbf{A}) = \text{tr}(\mathbf{A}_k) + \text{tr}(\mathbf{A} - \mathbf{A}_k)$
3. Compute  $\text{tr}(\mathbf{A}_k)$  exactly
4. Return  $\text{Hutch}++(\mathbf{A}) := \text{tr}(\widetilde{\mathbf{A}}_k) + H_\ell(\mathbf{A} - \widetilde{\mathbf{A}}_k)$

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

## Finding a good Low-Rank Approximation

⊙ Recall Randomized SVD:

1. Let  $\mathbf{S} \in \mathbb{R}^{n \times (2k+1)}$  have iid  $\mathcal{N}(0,1)$  entries
2. Compute  $\mathbf{Q} = \text{orth}(\mathbf{A}\mathbf{S})$
3. Then  $\widetilde{\mathbf{A}}_k = \mathbf{A}\mathbf{Q}\mathbf{Q}^T$  is pretty good:

$$\mathbb{E} \|\mathbf{A} - \widetilde{\mathbf{A}}_k\|_F^2 \leq 2\|\mathbf{A} - \mathbf{A}_k\|_F^2$$

⊙ Notice  $\text{tr}(\widetilde{\mathbf{A}}_k)$  can be computed quickly:

$$\text{tr}(\widetilde{\mathbf{A}}_k) = \text{tr}(\mathbf{A}\mathbf{Q}\mathbf{Q}^T) = \text{tr}(\mathbf{Q}^T\mathbf{A}\mathbf{Q})$$

# Hutch++

## ⊙ Final Algorithm

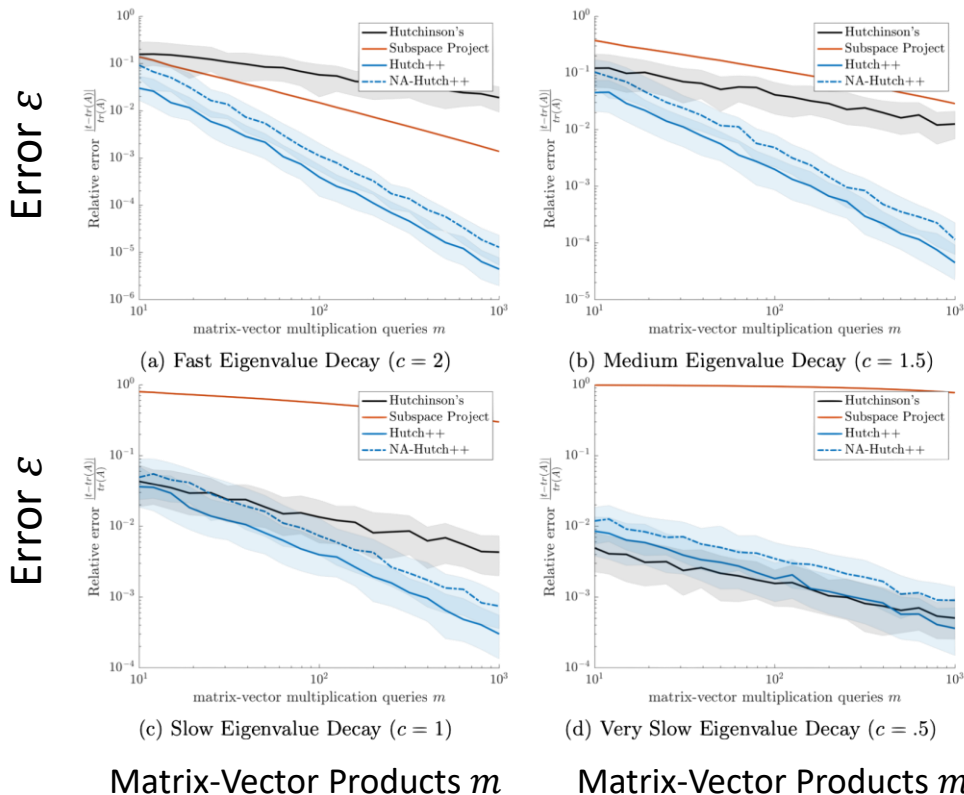
⊙ Input: Number of matrix-vector queries  $m$ , matrix  $A$

1. Sample  $S \in \mathbb{R}^{n \times \frac{m}{3}}$  and  $G \in \mathbb{R}^{n \times \frac{m}{3}}$  with iid  $\mathcal{N}(0,1)$  entries
2. Compute  $Q = qr(AS)$
3. Return  $tr(Q^T A Q) + \frac{1}{m/3} tr(G^T (I - Q Q^T) A (I - Q Q^T) G)$

```
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
```

# Experiments – Hutch++ vs Girard-Hutchinson

Hutch++ is fastest when  $A$  has fast eigenvalue decay. Here,  $\lambda_i = i^{-c}$



## Checkpoint – Hutch++

*Can we approximate  $\text{tr}(\mathbf{A})$  by computing few  $\mathbf{A}x_1, \dots, \mathbf{A}x_m$ ?*

- ⊙ **Yes, we can!**
- ⊙ Girard-Hutchinson  $H_\ell(\mathbf{A})$  uses  $O\left(\frac{1}{\varepsilon^2}\right)$  matrix-vector products
- ⊙ We can estimate  $\text{tr}(\mathbf{A})$  with  $O\left(\frac{1}{\varepsilon}\right)$  matrix-vector products
- ⊙ Four step algorithm:
  1. Find Low-Rank Approximation  $\widetilde{\mathbf{A}}_k$
  2. Compute  $\text{tr}(\widetilde{\mathbf{A}}_k)$
  3. Compute  $H_\ell(\mathbf{A} - \widetilde{\mathbf{A}}_k)$
  4. Return  $\text{tr}(\widetilde{\mathbf{A}}_k) + H_\ell(\mathbf{A} - \widetilde{\mathbf{A}}_k)$



# Part 3: Extensions of Hutch++

## Practical Considerations

- ⊙ What if my matrix is not PSD?

This is fine. Same proof shows  $|\text{Hutch++}(\mathbf{A}) - \text{tr}(\mathbf{A})| \leq \varepsilon \|\mathbf{A}\|_1$

- ⊙ Can we estimate the variance of Hutch++?
- ⊙ Can we maximize parallelism without wasting matrix-vector products?
- ⊙ Can we just be more efficient?
- ⊙ Can we take advantage of the fact that  $\mathbf{A} = f(\mathbf{B})$ ?

## A Posteriori Variance Estimation

- ⊙ What is the confidence interval / actual variance of my Hutch++ sample?
- ⊙ Intuitive answer: Sample variance of the Girard-Hutchinson Samples.
- ⊙ Better answer: use Bootstrapping to estimate the true variance below

$$\text{Lemma: } \text{Var}[\text{Hutch++}(\mathbf{A})] = \mathbb{E}[\text{Var}[H_\ell(\mathbf{A} - \widetilde{\mathbf{A}}_k) | \mathbf{Q}]]$$

*Proof.* By law of total variance,

$$\text{Var}[\text{Hutch++}(\mathbf{A})] = \mathbb{E}[\text{Var}[\text{Hutch++}(\mathbf{A}) | \mathbf{Q}]] + \text{Var}[\mathbb{E}[\text{Hutch++}(\mathbf{A}) | \mathbf{Q}]]$$

But observe that

$$\mathbb{E}[\text{Hutch++}(\mathbf{A}) | \mathbf{Q}] = \text{tr}(\widetilde{\mathbf{A}}_k) + \text{tr}(\mathbf{A} - \widetilde{\mathbf{A}}_k) = \text{tr}(\mathbf{A})$$

So

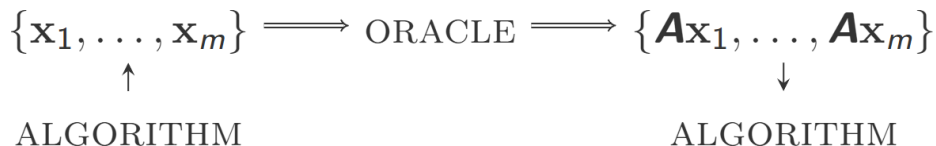
$$\text{Var}[\mathbb{E}[\text{Hutch++}(\mathbf{A}) | \mathbf{Q}]] = 0$$

## Maximizing Parallelism

- ⊙ To maximize parallelism, we want none of our matrix-vector products to depend on previous matrix-vector products.
- ⊙ In Hutch++, we must compute  $\mathbf{A}\mathbf{S}$  to compute  $\mathbf{Q}$  to compute  $\widetilde{\mathbf{A}}_k = \mathbf{A}\mathbf{Q}\mathbf{Q}^T$
- ⊙ Formally, this is an **adaptive** algorithm:



- ⊙ Non-obvious question: is trace estimation possible **non-adaptively**?



# Nyström++

- ⊙ Idea: Use the **Nyström approximation** to PSD  $A$ :
  - ⊙  $\widetilde{A}_k = A\langle S \rangle := AS(S^T AS)^+ S^T A$
  - ⊙ Unlike before, we construct  $\widetilde{A}_k$  non-adaptively!
- ⊙ Updated Algorithm:
  1. Use  $\frac{m}{2}$  matrix-vector products to build  $\widetilde{A}_k$
  2. Use  $\frac{m}{2}$  matrix-vector products to compute  $x_i^T A x_i$
  3. Return  $tr(\widetilde{A}_k) + \frac{1}{m/2} \sum_{i=1}^{m/2} (x_i^T A x_i - x_i^T \widetilde{A}_k x_i)$
- ⊙ Non-adaptive and basically always as fast/faster than Hutch++!
- ⊙ Analysis is more involved

## XTrace & Exchangeability

- ⊙ Motivation: Exchangeability Principle in Statistics
  - ⊙ Informally: any algorithm that uses iid samples, is unbiased, and has optimally small variance must be **exchangeable**.
  - ⊙ If we reorder our samples, the algorithm's output doesn't change.
- ⊙ The Girard-Hutchinson Estimator is exchangeable
  - ⊙  $H_\ell(\mathbf{A}) = \sum_{i=1}^{\ell} x_i^T \mathbf{A} x_i$
- ⊙ Hutch++ and Nyström++ are NOT exchangeable!
  - ⊙ Why?
- ⊙ XTrace: leave-one-out version of Hutch++ that is exchangeable  
XNysTrace: Nyström version of Xtrace
- ⊙ Even faster than Nyström++ in most cases!

## Krylov-Aware Trace Estimation

- ⊙ Remember we often have  $\mathbf{A} = f(\mathbf{B})$  where  $\mathbf{B}$  is in memory
- ⊙ Use Lanczos-FA Algorithm to approximately compute  $x \mapsto f(\mathbf{B})x$
- ⊙ Observation: Hutch++ style algorithms throw away pre-computed information if matrix-vector products with  $\mathbf{A}$  are computed via Lanczos-FA
- ⊙ “Krylov-Aware” Trace Estimator
  - Takes  $\mathbf{B}, f, m$  as inputs
  - Avoids throwing away information
- ⊙ Fastest algorithm I know for estimating  $f(\mathbf{B})$
- ⊙ Maybe incompatible with XTrace style algorithms?

## Checkpoint – Hutch++ Extensions

*Can we approximate  $\text{tr}(A)$  by computing few  $Ax_1, \dots, Ax_m$ ?*

- ⊙ **Yes, we can! Non-adaptively and super efficiently!**
- ⊙ Practitioner Advice:
  - ⊙ XNysTrace for PSD Matrices
  - ⊙ XTrace for non-PSD Matrices
  - ⊙ Krylov-Aware Trace Estimation for  $A = f(B)$
- ⊙ Natural theoretical question: Is better than  $O(\frac{1}{\varepsilon})$  possible?
  - ⊙ **No!**  $\Omega(\frac{1}{\varepsilon})$  is optimal!



# Part 4: Lower Bounds

## Lower Bound Super Rough Intuition

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

- ⊙ View the Matrix-Vector Oracle as a **limit on information** about  $A$ :
  1. Suppose  $A \sim \mathcal{D}$  is a random matrix
  2. Then  $\text{tr}(A)$  is a random variable with variance
  3. Any algorithm that computes few oracle queries has little information about  $\text{tr}(A)$
  4. Then, due to variance, the algorithm cannot predict  $\text{tr}(A)$  well
- ⊙ Step 3 is hard to prove rigorously.

## Removing the Algorithm's Agency

- ⊙ **Want to show:** Any algorithm that computes few oracle queries has little information about  $\text{tr}(A)$
- ⊙ **Problem:** Algorithm can pick many different query vectors  $x$
- ⊙ If instead algorithm had no freedom, we could use classical statistics to make lower bounds.

### Two Observations:

1. Without loss of generality, query vectors are orthonormal. (Why?)
2. Let  $G \in \mathbb{R}^{n \times n}$  be iid  $\mathcal{N}(0,1)$  matrix  
Let  $Q \in \mathbb{R}^{n \times k}$  have orthonormal columns  
Then  $GQ$  is a  $\mathcal{N}(0,1)$  matrix, even conditioned on knowing  $Q$ .  
(informal) If  $A$  is Gaussian, then  $\{Ax_1, \dots, Ax_m\}$  is independent of  $\{x_1, \dots, x_m\}$

# Wishart Anti-Concentration Method

## Hidden Wishart Theorem:

- ⊙ Let  $\mathbf{G} \in \mathbb{R}^{n \times n}$  be a  $\mathcal{N}(0,1)$  matrix
- ⊙ Let  $\mathbf{A} = \mathbf{G}^T \mathbf{G}$  be a **Wishart Matrix**
- ⊙ Suppose *Algorithm* computes  $\mathbf{A}x_1, \dots, \mathbf{A}x_m$ , possibly adaptively.
- ⊙ Then, there exists orthogonal  $\mathbf{V} \in \mathbb{R}^{n \times n}$  such that

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

where  $\tilde{\mathbf{A}} \in \mathbb{R}^{(n-m) \times (n-m)}$  is distributed as  $\tilde{\mathbf{G}}^T \tilde{\mathbf{G}}$  where  $\tilde{\mathbf{G}}$  is  $\mathcal{N}(0,1)$ , conditioned on all observations  $\{x_1, \mathbf{A}x_1, \dots, x_m, \mathbf{A}x_m\}$

- ⊙  $\mathbf{\Delta}$  is known deterministically by the algorithm

We can exactly separate the information we do/don't know about  $\mathbf{A}$ !

## Wishart Anti-Concentration Method

⊙ Consider any (possibly adaptive) algorithm after  $m$  queries. Then,

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

2. Let  $t$  be *Algorithm's* estimate of  $tr(\mathbf{A})$ . Define  $\tilde{t} := t - tr(\mathbf{\Delta})$

3. Note  $tr(\mathbf{A}) = \|\mathbf{G}\|_F^2 \sim \chi_{n^2}^2$  and  $tr(\tilde{\mathbf{A}}) \sim \chi_{(n-m)^2}^2$

- $|t - tr(\mathbf{A})| = |\tilde{t} - tr(\tilde{\mathbf{A}})| \geq \Omega(n - m)$
- $tr(\mathbf{A}) \leq O(n^2)$

4. Enforce  $|t - tr(\mathbf{A})| \leq \varepsilon tr(\mathbf{A})$   
 $(n - m) \leq \varepsilon C n^2$

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

## Checkpoint – Lower Bounds

*Can we approximate  $\text{tr}(A)$  by computing few  $Ax_1, \dots, Ax_m$ ?*

- ⊙ **Yes, but only if we take  $m = \Omega\left(\frac{1}{\varepsilon}\right)$  queries!**
- ⊙ Proven via Hidden Wishart, but other methods exist:
  - ⊙ Communication Complexity
  - ⊙ Statistical Hypothesis Testing
  - ⊙ Block Krylov Reduction
- ⊙ Hidden second meaning: Wishart matrices are unstructured covariance matrices.
- ⊙ If we want to estimate the trace of a covariance matrix (e.g. a Kernel matrix or Hessian), then do not expect that better than Hutch++ is possible! Otherwise, faster might exist!

**The End.**

**Thanks!**

**Questions?**