# Applied Random Matrix Theory

## Joel A. Tropp

Steele Family Professor of
Applied & Computational Mathematics

**Computing + Mathematical Sciences**

**California Institute of Technology**

`jtropp@cms.caltech.edu`

1

# SIAM Journal on Mathematics of Data Science (SIMODS)

https://simods.siam.org/

# What is a *Random* Matrix?

**Definition.** A `random matrix` is a matrix whose entries are random variables, not necessarily independent.

### A random matrix in captivity:

$$\begin{bmatrix} 0.0000 & -1.3077 & -1.3499 & 0.2050 & 0.0000 \\ 1.8339 & 0.0000 & -1.3077 & 0.0000 & 0.2050 \\ -2.2588 & 1.8339 & 0.0000 & -1.3077 & -1.3499 \\ 2.7694 & 0.0000 & 1.8339 & 0.0000 & -1.3077 \\ 0.0000 & 2.7694 & -2.2588 & 1.8339 & 0.0000 \end{bmatrix}$$

### What do we want to understand?

- ✏ Eigenvalues
- ✏ Eigenvectors
- ✏ Singular values
- ✏ Singular vectors
- ✏ Operator norms
- ✏ …

# Random Matrices in Statistics



**John Wishart**



3. *Multi-variate Distribution. Use of Quadratic co-ordinates.*

A comparison of equation (8) with the corresponding results (1) and (2) for uni-variate and bi-variate sampling, respectively, indicates the form the general result may be expected to take. In fact, we have for the simultaneous distribution in random samples of the $n$ variances (squared standard deviations) and the $\frac{n\,(n-1)}{2}$ product moment coefficients the following expression:
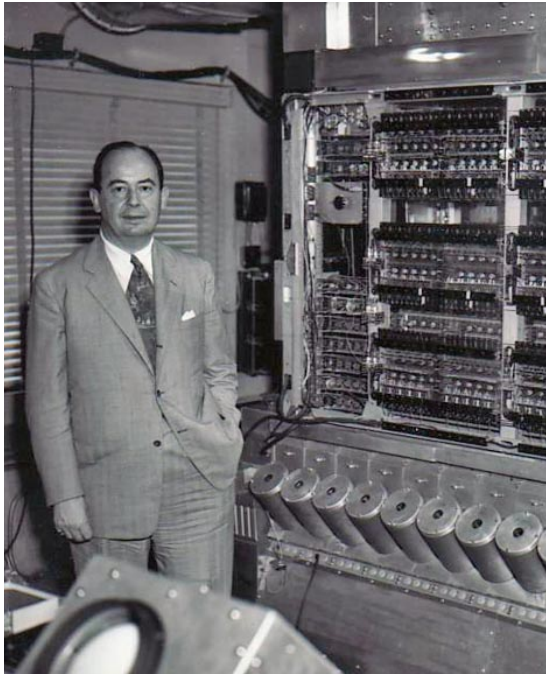
$$dp = \frac{\begin{vmatrix} A_{11} & A_{12} \ldots A_{1n} \\ A_{21} & A_{22} \ldots A_{2n} \\ \vdots & \vdots \quad \vdots \\ A_{n1} & A_{n2} \ldots A_{nn} \end{vmatrix}^{\frac{N-1}{2}}}{(\sqrt{\pi})^{\frac{1}{2}n\,(n-1)}\,\Gamma\left(\frac{N-1}{2}\right)\Gamma\left(\frac{N-2}{2}\right)\ldots\Gamma\left(\frac{N-n}{2}\right)}$$

$$\times\, e^{-A_{11}a_{11}-A_{22}a_{22}-\ldots-A_{nn}a_{nn}-2A_{12}a_{12}-2A_{13}a_{13}-\ldots-2A_{n-1\,n}a_{n-1\,n}}$$

$$\times\, \begin{vmatrix} a_{11} & a_{12} \ldots a_{1n} \\ a_{21} & a_{22} \ldots a_{2n} \\ \vdots & \vdots \quad \vdots \\ a_{n1} & a_{n2} \ldots a_{nn} \end{vmatrix}^{\frac{N-n-2}{2}} \quad da_{11}da_{12}\ldots da_{nn} \quad \ldots\ldots\ldots\ldots(9),$$

where $a_{pq}=s_p s_q r_{pq}$, and $A_{pq}=\frac{N}{2\sigma_p\sigma_q}\cdot\frac{\Delta_{pq}}{\Delta}$, $\Delta$ being the determinant

$$|\rho_{pq}|,\ p,\ q=1,\ 2,\ 3,\ \ldots\ n,$$

and $\Delta_{pq}$ the minor of $\rho_{pq}$ in $\Delta$.

❧ Sample covariance matrix for the multivariate normal distribution

# Random Matrices in Numerical Linear Algebra



**John von Neumann**

now combining (8.6) and (8.7) we obtain our desired result:

$$\text{Prob}\,(\lambda > 2\sigma^2 rn) < \frac{(rn)^{n-1/2}e^{-rn}\pi^{1/2}e^n \cdot 2^{n-2}}{\pi n^{n-1}(r-1)n}$$

(8.8)

$$= \left(\frac{2r}{e^{r-1}}\right)^n \times \frac{1}{4(r-1)(r\pi n)^{1/2}}\,.$$

We sum up in the following theorem:

(8.9) The probability that the upper bound $|A|$ of the matrix $A$ of (8.1) exceeds $2.72\sigma n^{1/2}$ is less than $.027 \times 2^{-n}n^{-1/2}$, that is, with probability greater than 99% the upper bound of $A$ is less than $2.72\sigma n^{1/2}$ for $n = 2, 3, \cdots$.

This follows at once by taking $r = 3.70$.

❧ Model for floating-point errors in LU decomposition

# Random Matrices in Nuclear Physics

**Random sign symmetric matrix**

The matrices to be considered are $2N + 1$ dimensional real symmetric matrices; $N$ is a very large number. The diagonal elements of these matrices are zero, the non diagonal elements $v_{ik} = v_{ki} = \pm v$ have all the same absolute value but random signs. There are $\mathfrak{N} = 2^{N(2N+1)}$ such matrices. We shall calculate, after an introductory remark, the averages of $(H'')_{00}$ and hence the strength function $S'(x) = \sigma(x)$. This has, in the present case, a second interpretation: it also gives the density of the characteristic values of these matrices. This will be shown first.
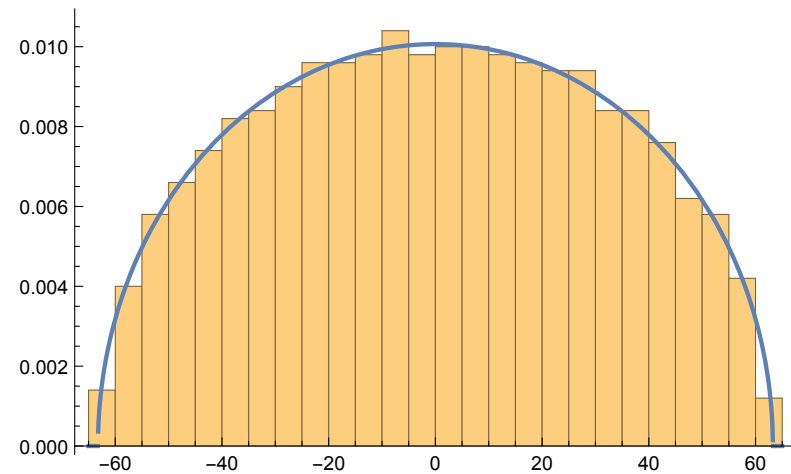
**Eugene Wigner**

❧ Model for the Hamiltonian of a heavy atom in a slow nuclear reaction

**Sources:** Wigner, *Ann. Math.* 1955. Photo from Nobel Foundation.

# Classical RMT

$$\begin{bmatrix} 0 & + & - & + & + & - & + \\ & 0 & + & - & - & - & + \\ & & 0 & + & - & + & + \\ & & & 0 & - & - & - \\ & & & & 0 & + & - \\ & * & & & & 0 & + \\ & & & & & & 0 \end{bmatrix}$$

Wigner $(d = 7)$

Distribution of eigenvalues $(d = 10^3)$

- Highly symmetric models
- Very precise results
- Strong resonances with other fields of mathematics

# Contemporary Applications of RMT

- Numerical linear algebra
- Numerical analysis
- Uncertainty quantification
- High-dimensional statistics
- Econometrics
- Approximation theory
- Sampling theory
- *Machine learning*

- Learning theory
- Mathematical signal processing
- Optimization
- Computer graphics and vision
- Quantum information theory
- Theory of algorithms
- Combinatorics
- …

**Sources:** (Drawn at random, nonuniformly) Halko et al. 2011; March & Biros 2014; Constantine & Gleich 2015; Koltchinskii 2011; Chen & Christensen 2013; Cohen et al. 2013; Bass & Groechenig 2013; Djolonga et al. 2013; Lopez-Paz et al. 2014; Fornasier et al. 2012; Morvant et al. 2012; Chen et al. 2014; Cheung et al. 2012; Chen et al. 2014; Holevo 2012; Harvey & Olver 2014; Cohen et al. 2014; Oliveira 2014.
Per Google Scholar, over $33,900$ papers with key "Random Matrix Theory."

# Contemporary RMT

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \end{bmatrix}$$

$\downarrow$     (sample random columns)     $\downarrow$

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 & 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 & 1 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 & 1 & 0 & 0 & 0 \end{bmatrix}$$

- ❧ Wide range of examples, many data-driven
- ❧ Results may sacrifice precision for applicability
- ❧ Theory is still developing

# Thesis Statement

# Modern applications demand new random matrix models and new analytical tools

# Matrix Concentration

❧ **Goal:** For a random matrix $Z$, find probabilistic bounds for

$$\| Z - \mathbb{E} Z \|$$

❧ An upper bound on this quantity ensures that

 ❧ Singular values of $Z$ and $\mathbb{E} Z$ are close
 ❧ Singular vectors of $Z$ and $\mathbb{E} Z$ are close (for isolated singular values)
 ❧ Linear functionals of $Z$ and $\mathbb{E} Z$ are close
 ❧ Spectral norm of $Z$ is controlled: $\| Z \| = \| \mathbb{E} Z \| \pm \| Z - \mathbb{E} Z \|$

$\| \cdot \|$ = spectral norm = largest singular value = $\ell_2$ operator norm

# The Independent Sum Model

$$Z = \sum_k S_k$$

## with $S_k$ independent

**Useful observation:** $\mathbb{E}\, Z = \sum_k \mathbb{E}\, S_k$

**Exercise:** Express the sample covariance matrix in this model

**Exercise:** Express column sampling (with replacement) from a fixed matrix

# The Bernstein Inequality

**Fact 1** (Bernstein 1920s)**. Suppose**

- ✒ *$S_1, S_2, S_3, \ldots$ are independent real random variables*
- ✒ *Each one is centered: $\mathbb{E}\, S_k = 0$*
- ✒ *Each one is bounded: $|S_k| \leq L$*

**Then**, *for $t > 0$,*

$$
\mathbb{P}\left\{\left|\textstyle\sum_k S_k\right| \geq t\right\} \quad \leq \quad 2 \cdot \begin{cases} \mathrm{e}^{-\mathrm{c}t^2/v}, & t \leq v/L \\ \mathrm{e}^{-\mathrm{c}t/L}, & t \geq v/L \end{cases} \qquad (\mathrm{c} = 3/8)
$$

*where the variance proxy is*

$$
v = \mathrm{Var}\left(\textstyle\sum_k S_k\right) = \textstyle\sum_k \mathbb{E}\, S_k^2
$$

# The Matrix Bernstein Inequality I

**Theorem 2** (T 2011). **Suppose**

- ☙ $S_1, S_2, S_3, \ldots$ *are independent random matrices with dimension* $d_1 \times d_2$
- ☙ *Each one is centered:* $\mathbb{E}\, S_k = 0$
- ☙ *Each one is bounded:* $\|S_k\| \leq L$

**Then**, *for* $t > 0$,

$$\mathbb{P}\left\{\left\|\sum_k S_k\right\| \geq t\right\} \quad \leq \quad (d_1 + d_2)\cdot \begin{cases} \mathrm{e}^{-ct^2/v}, & t \leq v/L \\ \mathrm{e}^{-ct/L}, & t \geq v/L \end{cases} \qquad (c = 3/8)$$

*where the matrix variance proxy is*

$$v = \max\left\{\left\|\sum_k \mathbb{E}(S_k S_k^*)\right\|, \ \left\|\sum_k \mathbb{E}(S_k^* S_k)\right\|\right\}$$

**Sources:** Tomczak–Jaegermann 1973; Lust-Piquard 1986; Pisier 1998; Rudelson 1999; Ahlswede & Winter 2002; Junge & Xu 2003, 2008; Rudelson & Vershynin 2005; Gross 2011; Recht 2011; Oliveira 2011; Tropp 2011–2015.

# The Matrix Bernstein Inequality II

**Theorem 3** (T 2011). **Suppose**

- 👉 $S_1, S_2, S_3, \ldots$ *are independent random matrices with dimension* $d_1 \times d_2$
- 👉 *Each one is centered:* $\mathbb{E}\, S_k = 0$
- 👉 *Each one is bounded:* $\|S_k\| \le L$

**Then**

$$\mathbb{E}\left\| \sum_k S_k \right\| \quad \le \quad \sqrt{2v \cdot \log(d_1 + d_2)} \;+\; \frac{1}{3}L \cdot \log(d_1 + d_2)$$

*where the matrix variance proxy is*

$$v = \max\left\{ \left\| \sum_k \mathbb{E}(S_k S_k^*) \right\|, \; \left\| \sum_k \mathbb{E}(S_k^* S_k) \right\| \right\}$$

**Sources:** Tomczak–Jaegermann 1973; Lust-Piquard 1986; Pisier 1998; Rudelson 1999; Ahlswede & Winter 2002; Junge & Xu 2003, 2008; Rudelson & Vershynin 2005; Gross 2011; Recht 2011; Oliveira 2011; Tropp 2011–2015.

# Example: Matrix Sparsification

$$A = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 2 & 4 & 6 & 8 \\ 3 & 6 & 9 & 12 \\ 4 & 8 & 12 & 16 \end{bmatrix} \longrightarrow \hat{A} = \begin{bmatrix} & 2 & & \\ & 4 & & 8 \\ 3 & 6 & 9 & 12 \\ & & 12 & 16 \end{bmatrix}$$

- ✎ **Goal:** Find a sparse matrix $\hat{A}$ for which $\|A - \hat{A}\|$ is small
- ✎ **Approach:** Non-uniform randomized sampling

**Sources:** Achlioptas & McSherry 2001, 2007; Arora et al. 2006; d'Asprémont 2008; Gittens & Tropp 2009; Nguyen et al. 2009; Drineas & Zouzias 2011; Achlioptas et al. 2013; Kundu & Drineas 2014; Tropp 2015.

# Sparsification: Sampling Model

🕮 Let $A$ be a fixed $d_1 \times d_2$ matrix

🕮 Construct a probability mass $\{p_{ij}\}$ on the matrix indices

🕮 Define a 1-sparse random matrix $S$ where

$$S = \frac{a_{ij}}{p_{ij}} \mathbf{E}_{ij} \quad \text{with probability } p_{ij}$$

🕮 The random matrix $S$ is an unbiased estimator for $A$

$$\mathbb{E}\, S = \sum\nolimits_{ij} \frac{a_{ij}}{p_{ij}} \mathbf{E}_{ij} \cdot p_{ij} = \sum\nolimits_{ij} a_{ij} \mathbf{E}_{ij} = A$$

🕮 To reduce the variance, average $r$ independent copies of $S$

$$\hat{A}_r = \frac{1}{r} \sum\nolimits_{k=1}^{r} S_k \quad \text{where } S_k \sim S$$

🕮 By construction, $\hat{A}_r$ has at most $r$ nonzero entries and approximates $A$

# Sparsification: Analysis

- **Recall:** $S = (a_{ij}/p_{ij})\mathbf{E}_{ij}$ with probability $p_{ij}$

- Bound for spectral norm:

$$\|S - \mathbb{E}\,S\| \le 2 \cdot \max_{ij} \frac{|a_{ij}|}{p_{ij}}$$

- Bound for variance:

$$\left\| \mathbb{E}(S - \mathbb{E}\,S)(S - \mathbb{E}\,S)^* \right\| \le \left\| \mathbb{E}\,SS^* \right\| = \left\| \sum_i \left( \sum_j \frac{|a_{ij}|^2}{p_{ij}} \right) \mathbf{E}_{ii} \right\| = \max_i \sum_j \frac{|a_{ij}|^2}{p_{ij}}$$

$$\left\| \mathbb{E}(S - \mathbb{E}\,S)^*(S - \mathbb{E}\,S) \right\| \le \left\| \mathbb{E}\,S^*S \right\| = \left\| \sum_j \left( \sum_i \frac{|a_{ij}|^2}{p_{ij}} \right) \mathbf{E}_{jj} \right\| = \max_j \sum_i \frac{|a_{ij}|^2}{p_{ij}}$$

- Construct probability mass $p_{ij} \propto |a_{ij}| + |a_{ij}|^2$ to control all terms

# Sparsification: Result

**Proposition 4** (Kundu & Drineas 2014; T 2015). **Suppose**

$$r \geq \varepsilon^{-2} \cdot \mathrm{srank}(\boldsymbol{A}) \cdot \max\{d_1, d_2\} \log(d_1 + d_2) \qquad (0 < \varepsilon \leq 1)$$

**Then** *the relative error in the $r$-sparse approximation $\hat{\boldsymbol{A}}_r$ satisfies*

$$\frac{\mathbb{E} \|\boldsymbol{A} - \hat{\boldsymbol{A}}_r\|}{\|\boldsymbol{A}\|} \leq 4\varepsilon$$

*The stable rank*

$$\mathrm{srank}(\boldsymbol{A}) := \frac{\|\boldsymbol{A}\|_{\mathrm{F}}^2}{\|\boldsymbol{A}\|^2} \leq \mathrm{rank}(\boldsymbol{A})$$

❧  The proof is an immediate consequence of matrix Bernstein

# Application: Fast Laplacian Solvers

**Theorem 5** (Kyng & Sachdeva 2016). **Suppose**

- *G is a weighted, undirected graph with $n$ vertices and $m$ edges*
- *$L$ is the combinatorial Laplacian of the graph $G$*

**Then**, *with high probability, the* SparseCholesky *algorithm produces*

- *A lower-triangular matrix $C$ with $\tilde{O}(m)$ nonzero entries that satisfies*

$$\frac{1}{2}L \preccurlyeq CC^* \preccurlyeq \frac{3}{2}L$$

- *The running time is $\tilde{O}(m)$*

*In particular, we can solve $Lx = b$ to machine precision in time $\tilde{O}(m)$*
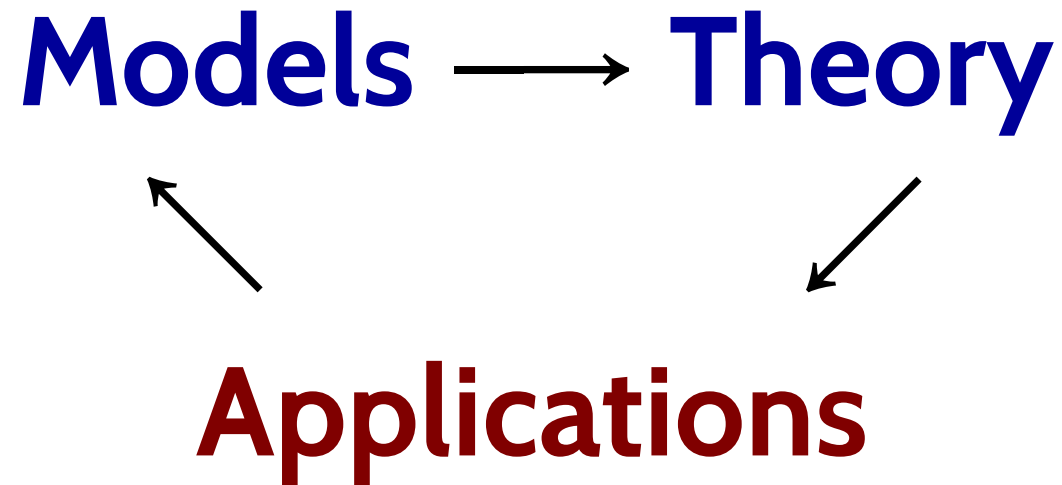
# SPARSECHOLESKY for a Graph Laplacian

$$L = \begin{bmatrix} a & \boldsymbol{u}^* \\ \boldsymbol{u} & \boldsymbol{L}_2 \end{bmatrix}_{n \times n} \quad \rightarrow \quad \boldsymbol{L}_2 - a^{-1} \begin{bmatrix} & \\ & \boldsymbol{u}\boldsymbol{u}^* \\ & \end{bmatrix}_{(n-1) \times (n-1)} \qquad \text{Subtract rank-1}$$

$$\rightarrow \quad \boldsymbol{L}_2 - a^{-1} \begin{bmatrix} & \times & \\ \times & & \times \\ & \times & \times \end{bmatrix}_{(n-1) \times (n-1)} \qquad \text{Sparsify rank-1}$$

- ✎ Direct computation of Cholesky factorization requires $O(n^2)$ operations per step
- ✎ Randomized approximation in $\tilde{O}(m/n)$ operations per step (amortized)
- ✎ Sampling probabilities are computed using graph theory
- ✎ Analysis depends on Bernstein inequality for matrix martingales!

**Sources:** Pisier & Xu 1997; Junge & Xu 2003, 2008; Oliveira 2011; Tropp 2011; Kyng & Sachdeva 2016.

# A Virtuous Cycle

**Models** $\longrightarrow$ **Theory**

**Applications**

# Contact & Papers

**email:** `jtropp@cms.caltech.edu`

**web:** `http://users.cms.caltech.edu/~jtropp`

**Monograph:**

- ❧ *An Introduction to Matrix Concentration Inequalities*. *Found. Trends Mach. Learn.*, 2015. Preprint: `arXiv:1501.01571`

**Papers:**

- ❧ "User-friendly tail bounds for sums of random matrices." *FoCM*, 2011
- ❧ "User-friendly tail bounds for matrix martingales." Caltech ACM Report 2011-01
- ❧ "Freedman's inequality for matrix martingales." *ECP*, 2011
- ❧ "From the joint convexity of relative entropy to a concavity theorem of Lieb." *PAMS*, 2012
- ❧ "Improved analysis of the subsampled randomized Hadamard transform." *AADA*, 2011
- ❧ "The masked sample covariance estimator" with R. Chen & A. Gittens. *I&I*, 2012
- ❧ "Tail bounds for all eigenvalues of a sum of random matrices" with A. Gittens. Caltech ACM Report 2014-02
- ❧ "Matrix concentration inequalities via the method of exchangeable pairs" with L. Mackey et al. *Ann. Probab.*, 2014
- ❧ "Subadditivity of matrix $\varphi$-entropy and concentration of random matrices" with R. Chen. *EJP*, 2014
- ❧ "Efron–Stein inequalities for random matrices" with D. Paulin & L. Mackey. *Ann. Probab.*, 2016
- ❧ "Second-order matrix concentration inequalities." *ACHA*, 2016
- ❧ "The expected norm of a sum of independent random matrices: An elementary approach," *HDP 7*, 2016