
Finding Structure with Randomness



Joel A. Tropp

Applied & Computational Mathematics
California Institute of Technology
`jtropp@acm.caltech.edu`

Joint with P.-G. Martinsson and N. Halko
Applied Mathematics, Univ. Colorado at Boulder

Top 10 Scientific Algorithms

list (here, the list is in chronological order; however, the articles appear in no particular order):

- Metropolis Algorithm for Monte Carlo
- Simplex Method for Linear Programming
- Krylov Subspace Iteration Methods
- The Decompositional Approach to Matrix Computations
- The Fortran Optimizing Compiler
- QR Algorithm for Computing Eigenvalues
- Quicksort Algorithm for Sorting
- Fast Fourier Transform
- Integer Relation Detection
- Fast Multipole Method

With each of these algorithms or approaches, there is a person or group receiving credit for inventing or

we
enc
vol
of v
of c
way
eve
rela
pro
are
high
J
ing
woi
plai
whi
not

Source: Dongarra and Sullivan, *Comput. Sci. Eng.*, 2000.

The Decompositional Approach

“The underlying principle of the decompositional approach to matrix computation is that it is not the business of the matrix algorithmicists to solve particular problems but to construct computational platforms from which a variety of problems can be solved.”

- 🐼 A decomposition solves not one but many problems
- 🐼 Often expensive to compute but can be reused
- 🐼 Shows that apparently different algorithms produce the same object
- 🐼 Facilitates rounding-error analysis
- 🐼 Can be updated efficiently to reflect new information
- 🐼 Has led to highly effective black-box software

Source: Stewart 2000.

Low-Rank Matrix Approximation

$$\begin{matrix} \mathbf{A} & \approx & \mathbf{B} & \mathbf{C}, \\ m \times n & & m \times k & k \times n. \end{matrix}$$

Benefits:

- ☞ Exposes structure of the matrix
- ☞ Allows efficient storage
- ☞ Facilitates multiplication with vectors or other matrices

Applications:

- ☞ Principal component analysis
- ☞ Low-dimensional embedding of data
- ☞ Approximating continuum operators with exponentially decaying spectra
- ☞ Model reduction for PDEs with rapidly oscillating coefficients

Approximation of Massive Data

- 🐼 **Problem:** Major cost for numerical algorithms is data transfer
- 🐼 Cost scales, roughly, with number of passes **not** amount of arithmetic
- 🐼 Random access to data is expensive, so classical algorithms may fail
- 🐼 **Assume** a matrix–matrix product with data matrix takes one pass
- 🐼 Matrix multiplication is efficient in many architectures:
 - 🐼 Graphics processing units
 - 🐼 Multi-core processors
 - 🐼 Parallel computers
 - 🐼 Distributed systems

Model Problem

Given:

- An $m \times n$ matrix A with $m \geq n$
- Target rank k
- Oversampling parameter p

Construct an $n \times (k + p)$ matrix Q with orthonormal columns s.t.

$$\|A - QQ^*A\| \approx \min_{\text{rank}(B) \leq k} \|A - B\|,$$

- QQ^* is the orthogonal projector onto the range of Q
- The basis Q can be used to construct matrix decompositions

From Basis to Decomposition

Problem: Given the basis Q , where do we get a factorization?

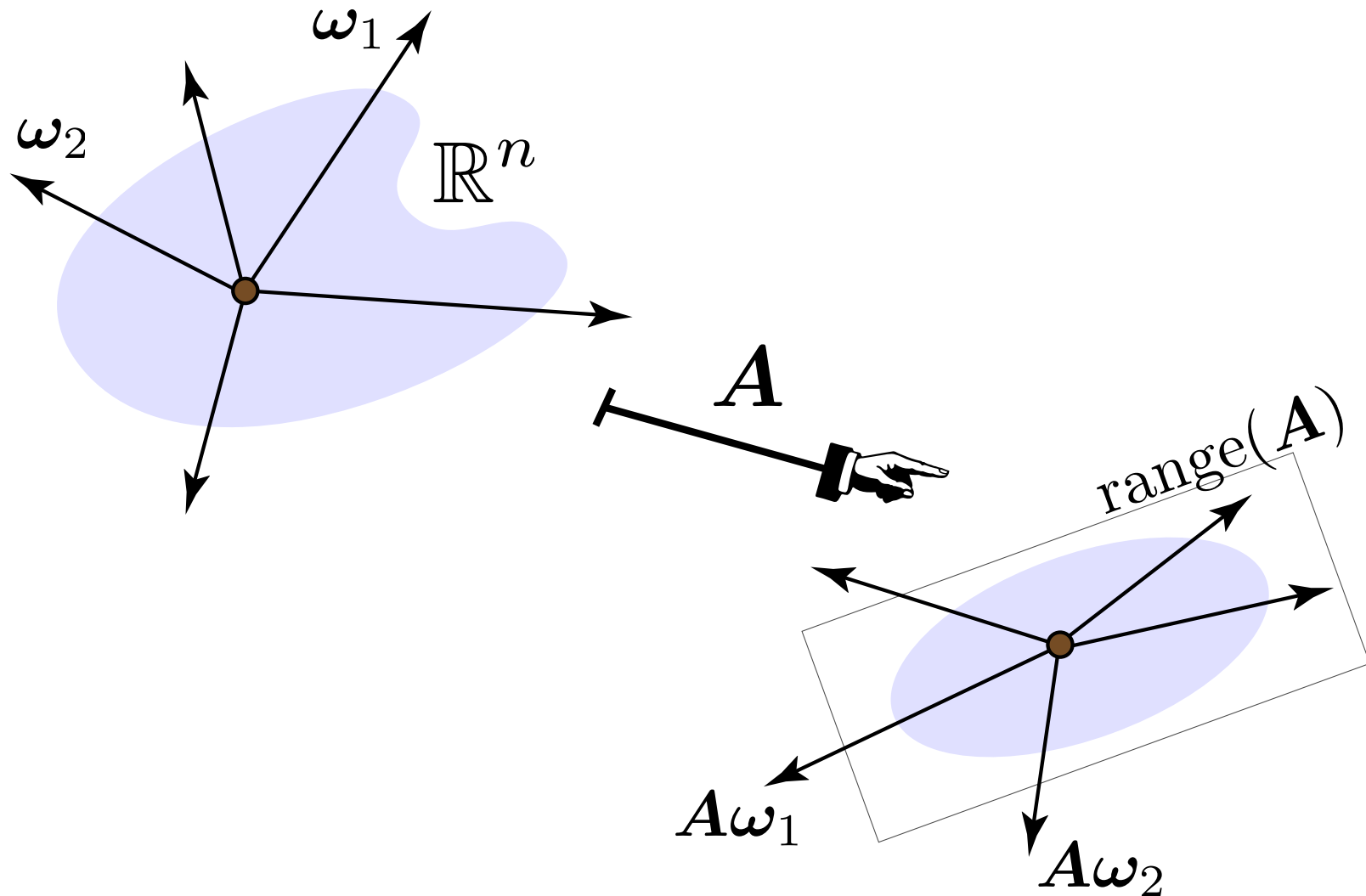
Example: Singular value decomposition

Assume A is $m \times n$ and Q is $n \times k$ where $A \approx QQ^*A$.

1. Form $k \times n$ matrix $B = Q^*A$ with one pass
2. Factor $B = U\Sigma V^*$ at cost $O(k^2n)$
3. Conclude $A \approx (QU)\Sigma V^*$.

Total Cost: One pass + one multiply $(m \times n \times k) + O(k^2n)$ flops

Random Sampling: Intuition



Proto-Algorithm for Model Problem

🐼 Converting this intuition into a computational procedure...

Input: An $m \times n$ matrix A , a target rank k , an oversampling parameter p

Output: An $m \times (k + p)$ matrix Q with orthonormal columns

1. Draw an $n \times (k + p)$ random matrix Ω .
 2. Form the matrix product $Y = A\Omega$.
 3. Construct an orthonormal basis Q for the range of Y .
-

Major Players: Deshpande, Drineas, Frieze, Kannan, Mahoney, Martinsson, Papadimitriou, Rokhlin, Sarlos, Tygert, Vempala (1998–2009)

Implementation Issues

Q: How much oversampling?

A: Remarkably, $p = 5$ or $p = 10$ is usually adequate!

Q: What random matrix?

A: For this application, standard Gaussian works nicely.

Q: How do we do the matrix–matrix multiply?

A: Exploit the computational architecture.

Q: How do we compute the orthonormal basis?

A: Carefully... Double Gram–Schmidt or Householder reflectors.

Q: How do we pick k ?

A: Can be done adaptively using a randomized error estimator.

Total Costs for Approximate k -SVD

Proto-Algorithm:

1 pass + 2 multiplies ($m \times n \times k$) + $k^2(m + n)$ flops

Classical Sparse Methods:

k passes + k multiplies ($m \times n \times 1$) + $k^2(m + n)$ flops

Classical Dense Methods:

k passes (or more) + mnk flops

Proto-Algorithm + Power Scheme

Problem: The singular values of the data matrix often decay slowly

Remedy: Apply the proto-algorithm to $(\mathbf{A}\mathbf{A}^*)^q \mathbf{A}$ for small q

Input: An $m \times n$ matrix \mathbf{A} , a target rank k , an oversampling parameter p

Output: An $m \times (k + p)$ matrix \mathbf{Q} with orthonormal columns

1. Draw an $n \times (k + p)$ random matrix $\mathbf{\Omega}$.
2. Form the matrix product $\mathbf{Y} = (\mathbf{A}\mathbf{A}^*)^q \mathbf{A}\mathbf{\Omega}$ by sequential multiplication.
3. Construct an orthonormal basis \mathbf{Q} for the range of \mathbf{Y} .

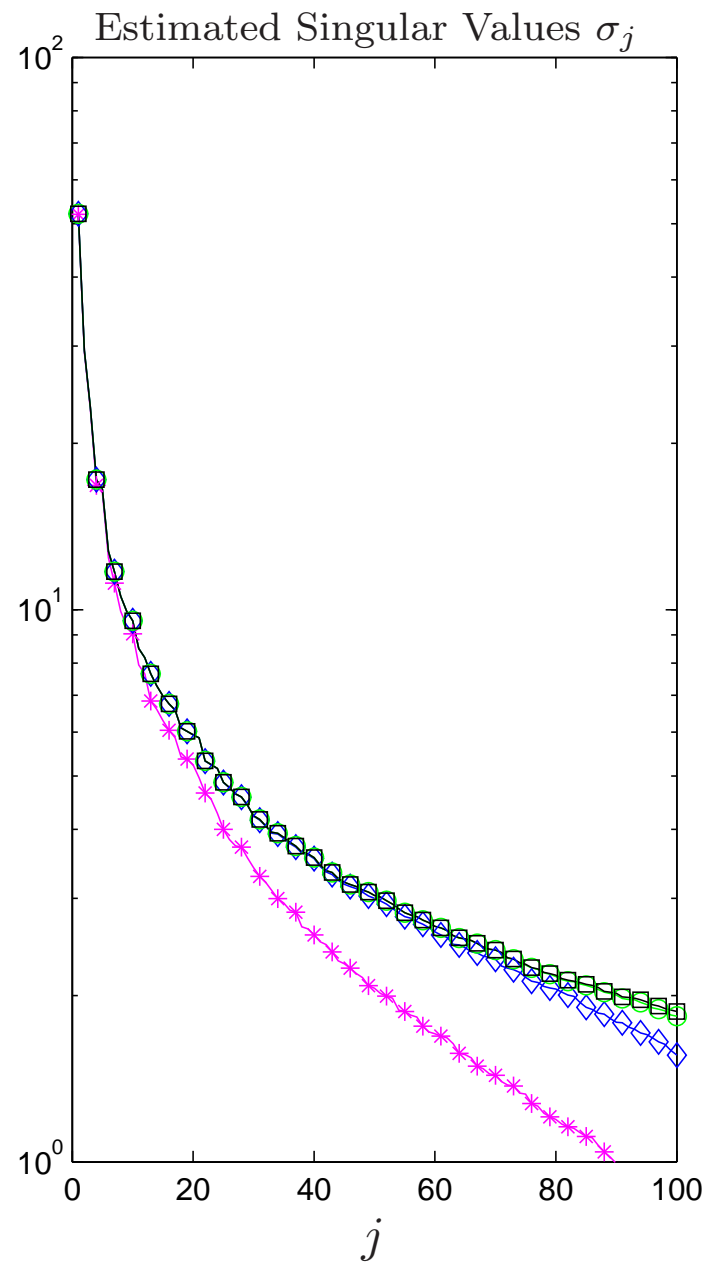
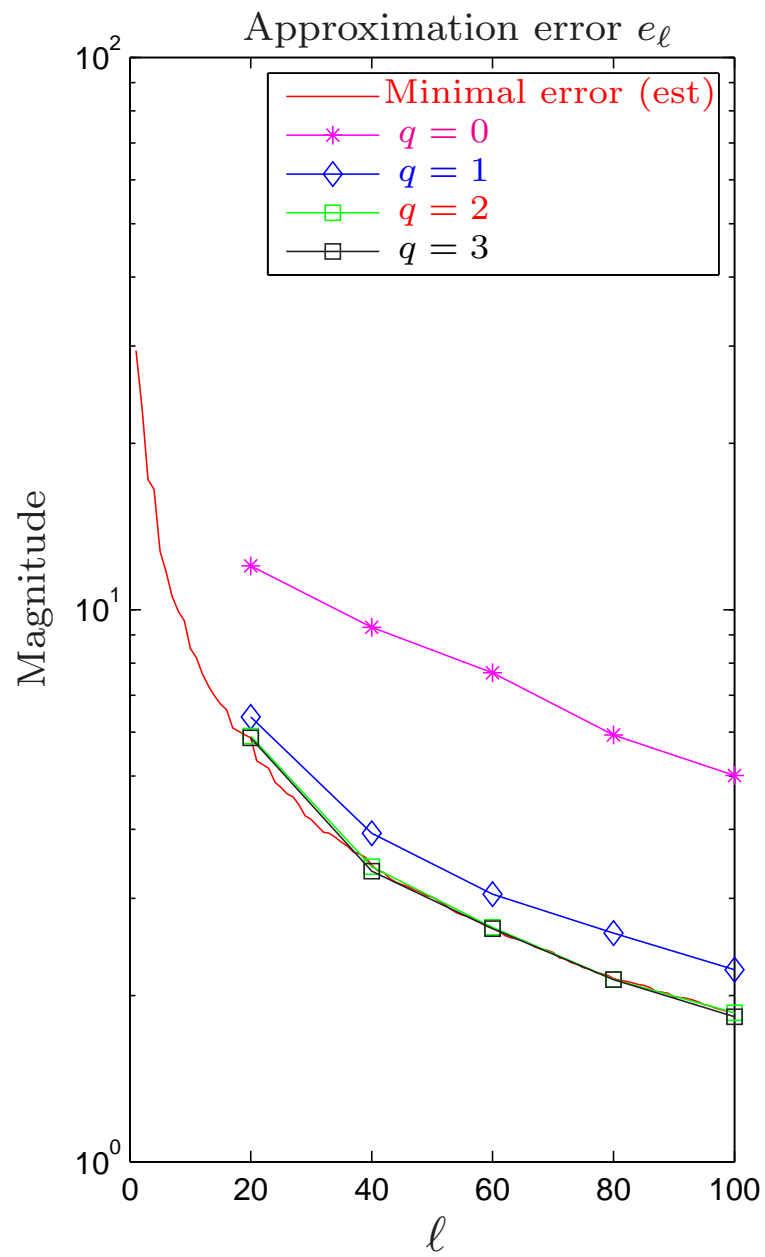
Total Cost: q passes + q multiplies $(m \times n \times k) + O(kmn)$ flops.

Eigenfaces

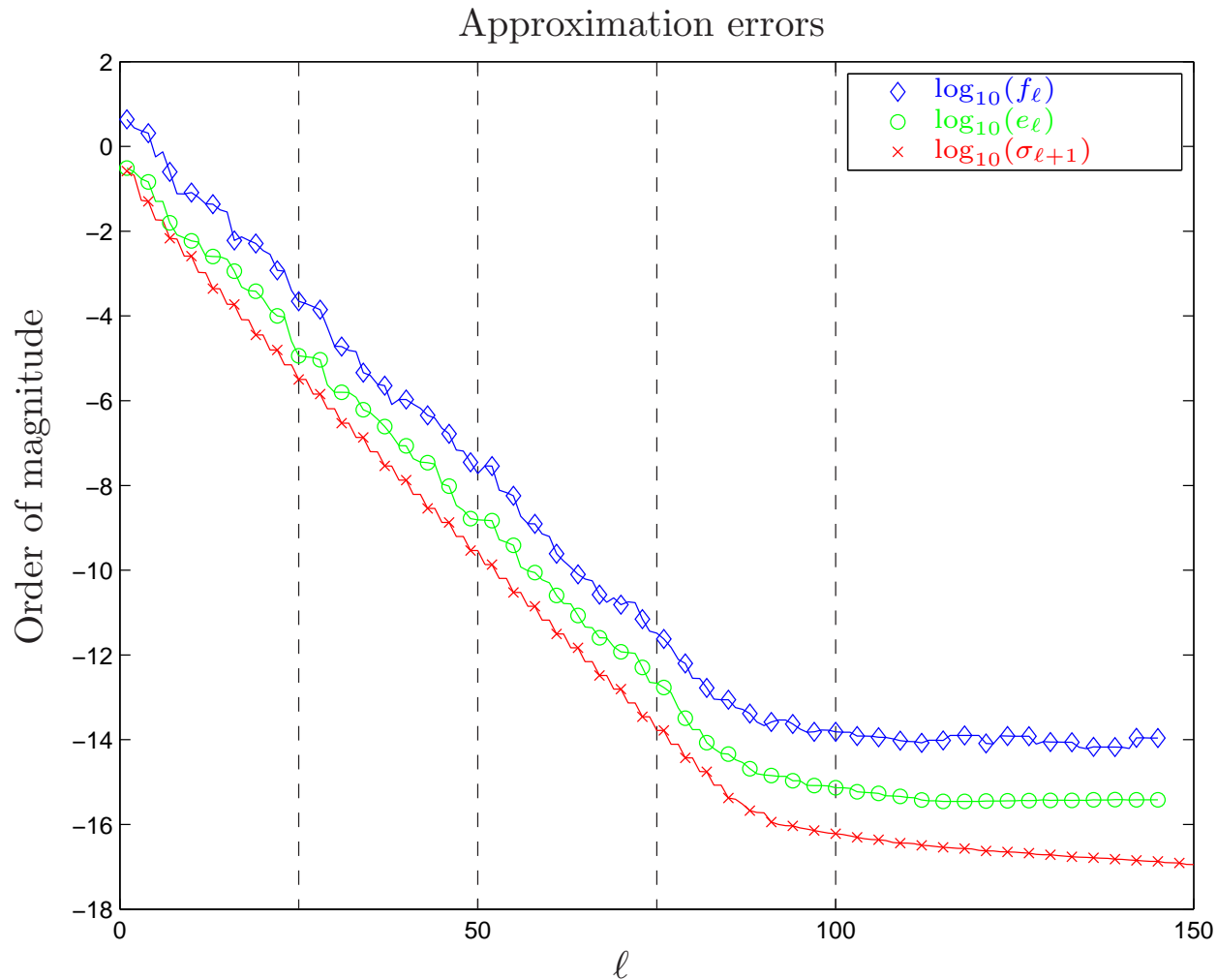
- Database consists of 7,254 photographs with 98,304 pixels each
- Form $98,304 \times 7,254$ data matrix \tilde{A}
- **Total storage:** 5.4 Gigabytes (uncompressed)
- Center each column and scale to unit norm to obtain A
- The dominant left singular vectors are called **eigenfaces**
- Attempt to compute first 100 eigenfaces using power scheme



Image: Scholarpedia article “Eigenfaces,” 12 October 2009



Approximating a Helmholtz Integral Operator



Error Bound for Proto-Algorithm

Theorem 1. [HMT 2009] Assume

- the matrix \mathbf{A} is $m \times n$ with $m \geq n$;
- the optimal error $\sigma_{k+1} = \min_{\text{rank}(\mathbf{B}) \leq k} \|\mathbf{A} - \mathbf{B}\|$;
- the test matrix $\mathbf{\Omega}$ is $n \times (k + p)$ standard Gaussian.

Then the basis \mathbf{Q} computed by the proto-algorithm satisfies

$$\mathbb{E} \|\mathbf{A} - \mathbf{Q}\mathbf{Q}^* \mathbf{A}\| \leq \left[1 + \frac{4\sqrt{k+p}}{p-1} \cdot \sqrt{n} \right] \sigma_{k+1}.$$

The probability of a substantially larger error is negligible.

Error Bound for Power Scheme

Theorem 2. [HMT 2009] Assume

- the matrix \mathbf{A} is $m \times n$ with $m \geq n$;
- the optimal error $\sigma_{k+1} = \min_{\text{rank}(\mathbf{B}) \leq k} \|\mathbf{A} - \mathbf{B}\|$;
- the test matrix $\mathbf{\Omega}$ is $n \times (k + p)$ standard Gaussian.

Then the basis \mathbf{Q} computed by the proto-algorithm satisfies

$$\mathbb{E} \|\mathbf{A} - \mathbf{Q}\mathbf{Q}^* \mathbf{A}\| \leq \left[1 + \frac{4\sqrt{k+p}}{p-1} \cdot \sqrt{n} \right]^{1/q} \sigma_{k+1}.$$

The probability of a substantially larger error is negligible.

- The power scheme drives the extra factor to one exponentially fast!

Inner Workings I

Assume

• \mathbf{A} is $m \times n$ with SVD

$$\mathbf{A} = \mathbf{U} \begin{bmatrix} & k & n-k \\ \boldsymbol{\Sigma}_1 & & \\ & & \boldsymbol{\Sigma}_2 \end{bmatrix} \begin{bmatrix} \mathbf{V}_1^* \\ \mathbf{V}_2^* \end{bmatrix} \begin{matrix} k \\ n-k \end{matrix}$$

• Let $\boldsymbol{\Omega}$ be a test matrix, decomposed as

$$\boldsymbol{\Omega}_1 = \mathbf{V}_1^* \boldsymbol{\Omega} \quad \text{and} \quad \boldsymbol{\Omega}_2 = \mathbf{V}_2^* \boldsymbol{\Omega}.$$

• Construct the sample matrix $\mathbf{Y} = \mathbf{A}\boldsymbol{\Omega}$.

Theorem 3. [BMD09, HMT09] *When $\boldsymbol{\Omega}_1$ has full row rank,*

$$\|(\mathbf{I} - \mathbf{P}_Y)\mathbf{A}\|^2 \leq \|\boldsymbol{\Sigma}_2\|^2 + \|\boldsymbol{\Sigma}_2\boldsymbol{\Omega}_2\boldsymbol{\Omega}_1^\dagger\|^2.$$

Inner Workings II

- When Ω is Gaussian, Ω_1 and Ω_2 are independent.
- Taking the expectation w.r.t. Ω_2 first...

$$\mathbb{E}_2 \left\| \Sigma_2 \Omega_2 \Omega_1^\dagger \right\| \leq \|\Sigma_2\| \left\| \Omega_1^\dagger \right\|_{\text{F}} + \|\Sigma_2\|_{\text{F}} \left\| \Omega_1^\dagger \right\|.$$

- The expectations of the norms w.r.t. Ω_1 satisfy

$$\mathbb{E} \left\| \Omega_1^\dagger \right\|_{\text{F}} \leq \sqrt{\frac{k}{p-1}} \quad \text{and} \quad \mathbb{E} \left\| \Omega_1^\dagger \right\| \leq \frac{e\sqrt{k+p}}{p}.$$

- Conclude

$$\mathbb{E} \left\| (\mathbf{I} - \mathbf{P}_Y) \mathbf{A} \right\| \leq \left[1 + \sqrt{\frac{k}{p-1}} \right] \sigma_{k+1} + \frac{e\sqrt{k+p}}{p} \left(\sum_{j=k+1}^{\infty} \sigma_j^2 \right)^{1/2}$$

Result for Structured Random Matrices

Theorem 4. [HMT09] *Suppose that Ω is an $n \times \ell$ SRFT matrix where*

$$\ell \gtrsim (k + \log n) \log k.$$

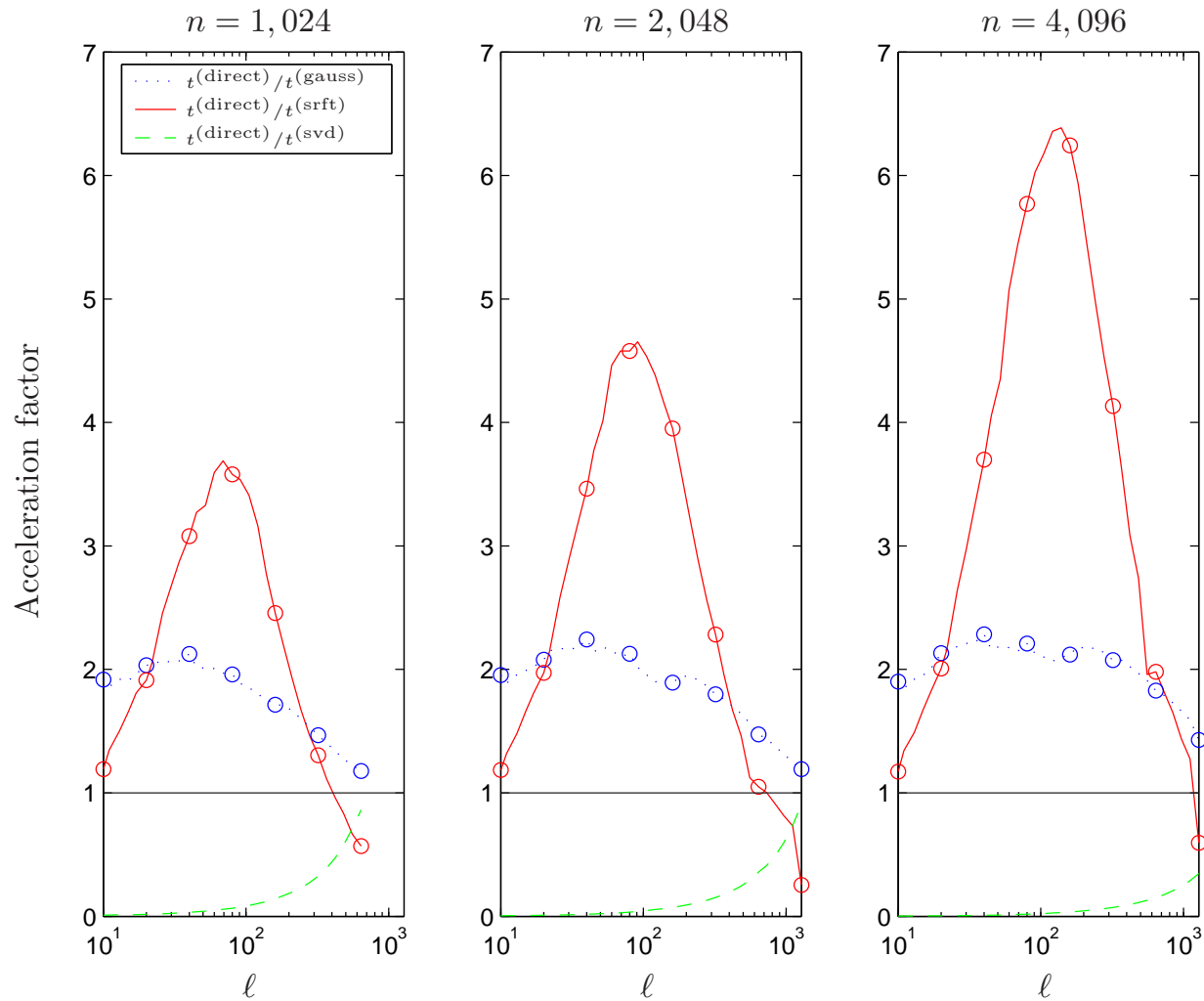
Then

$$\|(\mathbf{I} - \mathbf{P}_Y)\mathbf{A}\| \leq \sqrt{1 + \frac{Cn}{\ell}} \cdot \sigma_{k+1},$$

except with probability k^{-c} .

- Follows from same approach
- Uses Rudelson's lemma to show that random rows from a randomized Fourier transform form a well-conditioned set

Faster SVD with Structured Randomness



To learn more...

E-mail:

✉ jtropp@acm.caltech.edu

Web: <http://www.acm.caltech.edu/~jtropp>

Papers:

✉ HMT, “Finding Structure with Randomness: Stochastic Algorithms for Computing Approximate Matrix Decompositions,” submitted 2009