
Finding Structure with Randomness



Joel A. Tropp

Computing + Mathematical Sciences
California Institute of Technology
`jtropp@cms.caltech.edu`

Thanks: Alex Gittens (eBay), Michael Mahoney (Berkeley Stat),
Gunnar Martinsson (Boulder Appl. Math), Mark Tygert (Facebook)

Primary Sources for Tutorial

- 🐼 T. *User-Friendly Tools for Random Matrices: An Introduction*.
Submitted to *FnTML*, 2014.

tinyurl.com/pobvezn

- 🐼 Halko, Martinsson, and T. “Finding structure with randomness...”
SIAM Rev., 2011.

tinyurl.com/p5b5uw6

[Refs] <http://users.cms.caltech.edu/~jtropp/notes/Tro14-User-Friendly-Tools-FnTML-draft.pdf>
<http://users.cms.caltech.edu/~jtropp/papers/HMT11-Finding-Structure-SIREV.pdf>

Download the slides:

tinyurl.com/nbq2erb

[Ref] <http://users.cms.caltech.edu/~jtropp/slides/Tro14-Finding-Structure-ICML.pdf>

Matrix Decompositions & Approximations

Top 10 Scientific Algorithms

(Note: 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

[Ref] Dongarra and Sullivan 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

[Ref] Stewart 2000.

Matrix Approximations

“Matrix nearness problems arise in many areas... A common situation is where a matrix A approximates a matrix B and B is known to possess a property P ... An intuitively appealing way of improving A is to replace it by a nearest matrix X with property P .

“Conversely, in some applications it is important that A does not have a certain property P and it useful to know how close A is to having the undesirable property.”

- ☞ Approximations can purge an undesirable property (ill-conditioning)
- ☞ Can enforce a property the matrix lacks (sparsity, low rank)
- ☞ Can identify structure in a matrix
- ☞ Perform regularization, denoising, compression, ...

[Ref] Higham 1989; Dhillon & T 2006.

What's Wrong with Classical Approaches?

👉 Nothing... when the matrices are small and fit in core memory

👉 **Challenges:**

👉 Medium- to large-scale data (Megabytes+)

👉 New architectures (multi-core, distributed, data centers, ...)

👉 **Why Randomness?**

👉 It works...

👉 Randomized approximations can be very effective

👉 Leads to multiplication-rich algorithms (low communication costs; highly optimized primitives)

Hour 1: Approximation via Random Sampling

🐼 **Goal:** Find a structured approximation to a given matrix

🐼 **Approach:**

- 🐼 Construct a simple unbiased estimator of the matrix
- 🐼 Average independent copies to reduce variance

🐼 **Examples:**

- 🐼 Matrix sparsification
- 🐼 Random features

🐼 **Analysis:** Matrix Bernstein inequality

🐼 **Shortcomings:** Low precision; no optimality properties

Hour 2: Two-Stage Randomized Algorithms

🐼 **Goal:** Construct near-optimal, low-rank matrix approximations

🐼 **Approach:**

- 🐼 Use randomness to find a subspace that captures most of the action
- 🐼 Compress the matrix to this subspace, and apply classical NLA

🐼 **Randomized Range Finder:**

- 🐼 Multiply random test vectors into the target matrix and orthogonalize
- 🐼 Apply several steps of subspace iteration to improve precision

🐼 **Some Low-Rank Matrix Decompositions:**

- 🐼 Truncated singular value decomposition
- 🐼 Interpolative approximations, matrix skeleton, and CUR
- 🐼 Nyström approximations for psd matrices

Some Wisdom* from Scientific Computing

**“Who cares about the optimality of an approximation?
Who cares if I solve a specified computational problem?
My algorithm does great on the test set.”** —Nemo

- 🐛 **Optimality.** If your approximation is suboptimal, you could do better.
- 🐛 **Validation.** If your algorithm does not fit the model reliably, you cannot attribute success to *either* the model *or* the algorithm.
- 🐛 **Verification.** If your algorithm does not solve a specified problem, you cannot easily check whether it has bugs.
- 🐛 **Modularity.** To build a large system, you want each component to solve a specified problem under specified conditions.
- 🐛 **Reproducibility.** To use an approach for a different problem, you need the method to have consistent behavior.

Approximation by Random Sampling

Matrix Approximation via Sampling

- Let \mathbf{A} be a fixed matrix we want to approximate
- Represent $\mathbf{A} = \sum_i \mathbf{A}_i$ as a sum (or integral) of **simple matrices**
- Construct a **simple random matrix** \mathbf{Z} by sampling terms; e.g.,

$$\mathbf{Z} = p_i^{-1} \mathbf{A}_i \quad \text{with probability } p_i$$

- Ensures that \mathbf{Z} is an unbiased estimator: $\mathbb{E} \mathbf{Z} = \mathbf{A}$
- **Average independent copies** to reduce variance: $\hat{\mathbf{A}} = \frac{1}{r} \sum_{r=1}^r \mathbf{Z}_r$
- **Examples? Analysis?**

[Refs] Maurey 1970s; Carl 1985; Barron 1993; Rudelson 1999; Achlioptas & McSherry 2002, 2007; Drineas et al. 2006; Rudelson & Vershynin 2007; Rahimi & Recht 2007, 2008; Shalev-Shwartz & Srebro 2008; ...

The Matrix Bernstein Inequality

Theorem 1. [T 2012] Assume

- $\mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3, \dots$ are indep. random matrices with dimension $m \times n$
- $\mathbb{E} \mathbf{X}_r = \mathbf{0}$ and $\|\mathbf{X}_r\| \leq L$ for each index r
- Compute the variance measure

$$v := \max \left\{ \left\| \sum_r \mathbb{E}[\mathbf{X}_r \mathbf{X}_r^*] \right\|, \left\| \sum_r \mathbb{E}[\mathbf{X}_r^* \mathbf{X}_r] \right\| \right\}$$

Then

$$\mathbb{P} \left\{ \left\| \sum_r \mathbf{X}_r \right\| \geq t \right\} \leq d \cdot \exp \left\{ \frac{-t^2/2}{v + Lt/3} \right\}$$

where $d := m + n$.

$\|\cdot\|$ = spectral norm; $*$ = conjugate transpose

[Refs] Oliveira 2009–2011; T 2010–2014. This version from T 2014, *User-Friendly Tools*, Chap. 6.

The Matrix Bernstein Inequality

Theorem 2. [T 2014] Assume

- $\mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3, \dots$ are indep. random matrices with dimension $m \times n$
- $\mathbb{E} \mathbf{X}_r = \mathbf{0}$ and $\|\mathbf{X}_r\| \leq L$ for each index r
- Compute the variance measure

$$v := \max \left\{ \left\| \sum_r \mathbb{E}[\mathbf{X}_r \mathbf{X}_r^*] \right\|, \left\| \sum_r \mathbb{E}[\mathbf{X}_r^* \mathbf{X}_r] \right\| \right\}$$

Then

$$\mathbb{E} \left\| \sum_r \mathbf{X}_r \right\| \leq \sqrt{2v \log d} + \frac{1}{3} L \log d$$

where $d := m + n$.

$\|\cdot\|$ = spectral norm; $*$ = conjugate transpose

[Refs] Chen et al. 2012; Mackey et al. 2014. This version from T 2014, *User-Friendly Tools*, Chap. 6.

Short History of Matrix Bernstein Inequality

Operator Khintchine and Noncommutative Martingales

- **Tomczak-Jaegermann 1974.** First operator Khintchine inequality; suboptimal variance.
- **Lust-Piquard 1986.** Operator Khintchine; optimal variance; suboptimal constants.
- Lust-Piquard & Pisier 1991. Operator Khintchine for trace class.
- Pisier & Xu 1997. Initiates study of noncommutative martingales.
- **Rudelson 1999.** First use of operator Khintchine for random matrix theory.
- Buchholz 2001, 2005. Optimal constants for operator Khintchine.
- Many more works in 2000s.

Matrix Concentration Inequalities

- **Ahlsvede & Winter 2002.** Matrix Chernoff inequalities; suboptimal variance.
- Christofides & Markström 2007. Matrix Hoeffding; suboptimal variance.
- Gross 2011; Recht 2011. Matrix Bernstein; suboptimal variance.
- **Oliveira 2011; T 2012.** Matrix Bernstein; optimal variance. Independent works!
- Chen et al. 2012; Mackey et al. 2014; T 2014. Expectation form of matrix inequalities.
- Hsu et al. 2012. Intrinsic dimension bounds; suboptimal form.
- **Minsker 2012.** Intrinsic dimension bounds; optimal form.
- T 2014. Simplified proof of intrinsic dimension bounds.
- **Mackey et al. 2014.** New proofs and results via exchangeable pairs.

[Ref] See T 2014, *User-Friendly Tools* for more historical information.

Error Estimate for Matrix Sampling

Corollary 3. [Matrix Sampling Estimator] Assume

- \mathbf{A} is a fixed $m \times n$ matrix and $d := m + n$
- \mathbf{Z} is a random matrix with $\mathbb{E} \mathbf{Z} = \mathbf{A}$ and $\|\mathbf{Z}\| \leq L$
- $\mathbf{Z}_1, \dots, \mathbf{Z}_r$ are iid copies of \mathbf{Z}
- Compute the per-sample variance

$$v := \max \left\{ \|\mathbb{E}[\mathbf{Z}\mathbf{Z}^*]\|, \|\mathbb{E}[\mathbf{Z}^*\mathbf{Z}]\| \right\}$$

Then the estimator $\hat{\mathbf{A}} = \frac{1}{r} \sum_{r=1}^r \mathbf{Z}_r$ satisfies

$$\mathbb{E} \|\mathbf{A} - \hat{\mathbf{A}}\| \leq \sqrt{\frac{2v \log d}{r}} + \frac{2L \log d}{3r}$$

[Refs] This version is new. Variants appear in Rudelson 1999; Ahlswede & Winter 2002; Gross 2011; Recht 2011. See T 2014, *User-Friendly Tools*, Chap. 1.

Comments on Matrix Sampling Estimators

- ☞ Constructing simple random matrix Z requires **insight + cleverness**
- ☞ Approximation $\|A - \hat{A}\|$ in **spectral norm** controls
 - ☞ All linear functions of the approximation (marginals)
 - ☞ All singular values and singular vectors
 - ☞ **Warning:** Frobenius-norm error bounds are usually vacuous!
- ☞ Sampling estimators are typically low precision
 - ☞ **Bottleneck:** Central Limit Theorem
 - ☞ Cost of reducing error is exorbitant (ε^{-2})
 - ☞ Error is generally not comparable with best possible approximation

[Ref] Tygert 2014.

Matrix Sparsification

Matrix Sparsification

🐼 **Challenge:** Can we replace a dense matrix by a sparse proxy?

$$\begin{bmatrix} x & x & x & x & x \\ x & x & x & x & x \\ x & x & x & x & x \\ x & x & x & x & x \\ x & x & x & x & x \end{bmatrix} \rightsquigarrow \begin{bmatrix} x & & x & x & x \\ & x & x & & x \\ x & & x & x & \\ & x & & & x \\ x & & x & x & \end{bmatrix}$$

🐼 **Idea:** Achlioptas & McSherry (2002, 2007) propose random sampling

🐼 **Goals:**

- 🐼 Accelerate spectral computations via Krylov methods
- 🐼 Compression, denoising, regularization, etc.

[Refs] Achlioptas & McSherry 2002, 2007; Arora et al. 2005; Gittens & Tropp 2009; d'Aspremont 2009; Drineas & Zouzias 2011; Achlioptas et al. 2013; Drineas & Kundra 2014. See T 2014, *User-Friendly Tools*, Chap. 6.

Writing a Matrix as a Sparse Sum

$$\mathbf{A} = \sum_{ij} a_{ij} \mathbf{E}_{ij}$$

a_{ij} denotes the (i, j) entry of \mathbf{A}

\mathbf{E}_{ij} is the standard basis matrix with a one in the (i, j) position and zeroes elsewhere

Sparsification by Random Sampling

☛ Let \mathbf{A} be a fixed $n \times n$ matrix

☛ Define sampling probabilities

$$p_{ij} = \frac{1}{2} \left[\frac{|a_{ij}|^2}{\|\mathbf{A}\|_F^2} + \frac{|a_{ij}|}{\|\mathbf{A}\|_{\ell_1}} \right] \quad \text{for } i, j = 1, \dots, n$$

☛ Let $\mathbf{Z} = p_{ij}^{-1} \cdot a_{ij} \cdot \mathbf{E}_{ij}$ with probability p_{ij}

☛ Let $\mathbf{Z}_1, \dots, \mathbf{Z}_r$ be iid copies of the estimator \mathbf{Z}

☛ Thus, $\hat{\mathbf{A}} = \frac{1}{r} \sum_{r=1}^r \mathbf{Z}_r$ is an r -sparse unbiased estimator of \mathbf{A}

$\|\cdot\|_F$ = Frobenius norm; $\|\cdot\|_{\ell_1}$ = elementwise ℓ_1 norm

[Refs] Kunder & Drineas 2014. T 2014, *User-Friendly Tools*, Chap. 6.

Analysis of Sparsification

☛ Recall: $\mathbf{Z} = p_{ij}^{-1} \cdot a_{ij} \cdot \mathbf{E}_{ij}$ with probability p_{ij}

☛ Observe: $p_{ij} \geq \frac{|a_{ij}|}{2 \|\mathbf{A}\|_{\ell_1}}$ and $p_{ij} \geq \frac{|a_{ij}|^2}{2 \|\mathbf{A}\|_F^2}$

☛ Uniform bound: $\|\mathbf{Z}\| \leq \max_{ij} \frac{|a_{ij}|}{p_{ij}} \cdot \|\mathbf{E}_{ij}\| \leq 2 \|\mathbf{A}\|_{\ell_1}$

☛ Variance computation:

$$\mathbb{E}[\mathbf{Z}\mathbf{Z}^*] = \sum_{ij} \frac{|a_{ij}|^2}{p_{ij}^2} \cdot \mathbf{E}_{ij}\mathbf{E}_{ij}^* \cdot p_{ij} \preccurlyeq \sum_{ij} 2 \|\mathbf{A}\|_F^2 \cdot \mathbf{E}_{ii} = 2n \|\mathbf{A}\|_F^2 \cdot \mathbf{I}$$

$$\mathbb{E}[\mathbf{Z}^*\mathbf{Z}] = \sum_{ij} \frac{|a_{ij}|^2}{p_{ij}^2} \cdot \mathbf{E}_{ij}^*\mathbf{E}_{ij} \cdot p_{ij} \preccurlyeq \sum_{ij} 2 \|\mathbf{A}\|_F^2 \cdot \mathbf{E}_{jj} = 2n \|\mathbf{A}\|_F^2 \cdot \mathbf{I}$$

☛ Thus, $\max\{\|\mathbb{E}[\mathbf{Z}\mathbf{Z}^*]\|, \|\mathbb{E}[\mathbf{Z}^*\mathbf{Z}]\|\} \leq 2n \|\mathbf{A}\|_F^2$

[Refs] Kundra & Drineas 2014. T 2014, *User-Friendly Tools*, Chap. 6.

Error Bound for Sparsification

☞ Corollary 3 provides the error bound

$$\mathbb{E} \|\mathbf{A} - \hat{\mathbf{A}}\| \leq \sqrt{\frac{2n \|\mathbf{A}\|_{\text{F}}^2 \log(2n)}{r}} + \frac{2 \|\mathbf{A}\|_{\ell_1} \log(2n)}{3r}$$

☞ Note that $\|\mathbf{A}\|_{\ell_1} \leq n \|\mathbf{A}\|_{\text{F}}$, and place error on relative scale:

$$\frac{\mathbb{E} \|\mathbf{A} - \hat{\mathbf{A}}\|}{\|\mathbf{A}\|} \leq \frac{\|\mathbf{A}\|_{\text{F}}}{\|\mathbf{A}\|} \left[\sqrt{\frac{2n \log(2n)}{r}} + \frac{2n \log(2n)}{3r} \right]$$

☞ **[Q]** What does this mean?

[Refs] Kundra & Drineas 2014. T 2014, *User-Friendly Tools*, Chap. 6.

Detour: The Stable Rank

• The *stable rank* of a matrix is defined as

$$\text{srank}(\mathbf{A}) := \frac{\|\mathbf{A}\|_{\text{F}}^2}{\|\mathbf{A}\|^2}$$

• In general, $1 \leq \text{srank}(\mathbf{A}) \leq \text{rank}(\mathbf{A})$

• When \mathbf{A} has either n rows or n columns, $1 \leq \text{srank}(\mathbf{A}) \leq n$

• **Assume** that \mathbf{A} has n unit-norm columns, so that $\|\mathbf{A}\|_{\text{F}}^2 = n$

• When all columns of \mathbf{A} are the same, $\|\mathbf{A}\|^2 = n$ and $\text{srank}(\mathbf{A}) = 1$

• When all columns of \mathbf{A} are orthogonal, $\|\mathbf{A}\|^2 = 1$ and $\text{srank}(\mathbf{A}) = n$

[Refs] Bourgain & Tzafriri 1987; Rudelson & Vershynin 2007; T 2009; Vershynin 2011. See T 2014, *User-Friendly Tools*, Chap. 6.

Error Bound for Sparsification II

• Fix a tolerance $\varepsilon \in (0, 1)$

• Suppose the sparsity r satisfies

$$r \geq \varepsilon^{-2} \cdot 2n \log(2n) \cdot \text{srank}(\mathbf{A})$$

• Then the r -sparse matrix $\hat{\mathbf{A}}$ achieves relative error

$$\frac{\mathbb{E} \|\mathbf{A} - \hat{\mathbf{A}}\|}{\|\mathbf{A}\|} \leq 2\varepsilon$$

• If $\text{srank}(\mathbf{A}) \ll n$, then $\text{nnz}(\hat{\mathbf{A}}) \ll n^2$ suffices

[Refs] Kundu & Drineas 2014. See T 2014, *User-Friendly Tools*, Chap. 6.

Random Features

Kernel Matrices

- Consider data points $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{X} \subset \mathbb{R}^d$
- Let $k : \mathcal{X} \times \mathcal{X} \rightarrow [-1, +1]$ be a bounded similarity measure
- Form $n \times n$ kernel matrix \mathbf{K} with entries $k_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$
- Useful for regression, classification, feature selection, ...
- **Challenge:** Kernel is big ($n \times n$) and expensive to evaluate $\mathcal{O}(dn^2)$
- **Opportunity:** The kernel often has low effective dimension
- **Idea:** Construct a randomized low-rank approximation of the kernel

[Refs] Schölkopf & Smola 2002; Rahimi & Recht 2007, 2008.

Random Features

- Let \mathcal{W} be an auxiliary space
- Let $\varphi : \mathcal{X} \times \mathcal{W} \rightarrow [-1, +1]$ be a bounded *feature map*
- Let w be a random variable taking values in \mathcal{W}
- **Assume** the random feature $\varphi(x; w)$ has the *reproducing property*

$$k(x, y) = \mathbb{E}_w [\varphi(x; w) \cdot \varphi(y; w)] \quad \text{for all } x, y \in \mathcal{X}$$

- **Example:** Random Fourier features for shift-invariant kernels
- **Related:** Random Walsh features for inner-product kernels

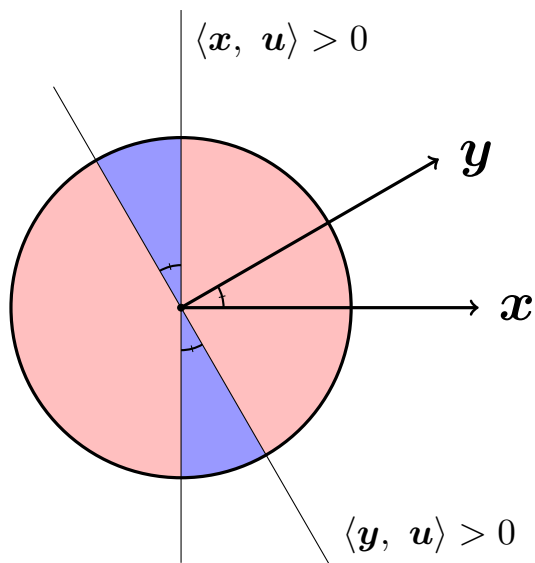
[Refs] Rahimi & Recht 2007, 2008; Maji & Berg 2009; Li et al. 2010; Vedaldi & Zisserman 2010; Vempati et al. 2010; Kar & Karnack 2012; Le et al. 2013; Pham & Pagh 2013; Hamid et al. 2014; ...

Example: Angular Similarity

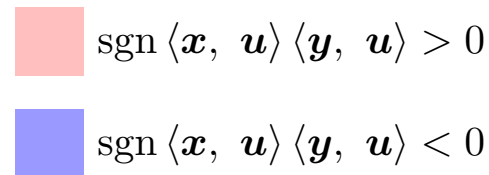
• For $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, consider the angular similarity

$$k(\mathbf{x}, \mathbf{y}) = \frac{2}{\pi} \arcsin \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|} = 1 - \frac{2\angle(\mathbf{x}, \mathbf{y})}{\pi} \in [-1, +1]$$

• Define random feature $\varphi(\mathbf{x}; \mathbf{w}) = \text{sgn} \langle \mathbf{x}, \mathbf{w} \rangle$ with $\mathbf{w} \sim \text{UNIF}(\mathbb{S}^{d-1})$



$$k(\mathbf{x}, \mathbf{y}) = \mathbb{E}[\varphi(\mathbf{x}; \mathbf{w}) \cdot \varphi(\mathbf{y}; \mathbf{w})]$$



[Refs] Early 1900s; Grothendieck 1953; Goemans & Williamson 1996.

Low-Rank Approximation of Kernel Matrix

- Draw a random vector $\mathbf{w} \in \mathcal{W}$
- Define $z_i = \varphi(\mathbf{x}_i; \mathbf{w})$ for each datum $i = 1, \dots, n$
- By the reproducing property, $k_{ij} = k(\mathbf{x}_i, \mathbf{x}_j) = \mathbb{E}[z_i z_j]$
- In matrix form: $\mathbf{K} = \mathbb{E}[\mathbf{z}\mathbf{z}^*]$ where $\mathbf{z} = [z_i] \in \mathbb{R}^n$
- That is, $\mathbf{Z} = \mathbf{z}\mathbf{z}^*$ is an unbiased rank-one estimator for \mathbf{K}
- Draw iid copies $\mathbf{Z}_1, \dots, \mathbf{Z}_r$ of the estimator \mathbf{Z}
- Thus, $\widehat{\mathbf{K}} = \frac{1}{r} \sum_{r=1}^r \mathbf{Z}_r$ is an unbiased rank- r estimator for \mathbf{K}

[Refs] Rahimi & Recht 2007, 2008; Lopez-Paz et al. ICML 2014.

Analysis of Random Features

☞ Recall: $\mathbf{Z} = \mathbf{z}\mathbf{z}^*$ where $\mathbb{E} \mathbf{Z} = \mathbf{K}$ and $\|\mathbf{z}\|_{\ell_\infty} \leq 1$

☞ Uniform bound:

$$\|\mathbf{Z}\| = \|\mathbf{z}\mathbf{z}^*\| = \|\mathbf{z}\|^2 \leq n$$

☞ Variance computation:

$$0 \preceq \mathbb{E}[\mathbf{Z}^2] = \mathbb{E}[\|\mathbf{z}\|^2 \cdot \mathbf{z}\mathbf{z}^*] \preceq n \cdot \mathbb{E}[\mathbf{z}\mathbf{z}^*] = n \cdot \mathbf{K}$$

☞ Thus, $\|\mathbb{E}[\mathbf{Z}^2]\| \leq n \|\mathbf{K}\|$

[Ref] Lopez-Paz et al. ICML 2014.

Error Bound for Random Features

🐼 Corollary 3 implies that

$$\mathbb{E} \|\mathbf{K} - \widehat{\mathbf{K}}\| \leq \sqrt{\frac{2n \|\mathbf{K}\| \log(2n)}{\mathbf{r}}} + \frac{2n \log(2n)}{3\mathbf{r}}$$

🐼 Define the *effective rank* $\rho := n / \|\mathbf{K}\|$ of the kernel to obtain

$$\frac{\mathbb{E} \|\mathbf{K} - \widehat{\mathbf{K}}\|}{\|\mathbf{K}\|} \leq \sqrt{\frac{2\rho \log(2n)}{\mathbf{r}}} + \frac{2\rho \log(2n)}{3\mathbf{r}}$$

🐼 For relative error $\mathcal{O}(\varepsilon)$, need at most \mathbf{r} random features where

$$\mathbf{r} \geq 2\varepsilon^{-2} \rho \log(2n)$$

🐼 Random features are efficient when $\rho \ll n$

[Ref] Lopez-Paz et al., ICML 2014.

Entr'acte

Optimal Low-Rank Approximation

- Let \mathbf{A} be an $n \times n$ matrix with SVD $\mathbf{A} = \sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{v}_i^*$
- Assume** normalization $\sum_{i=1}^n \sigma_i = 1$
- “Optimal” low-rank sampling: $\mathbf{Z} = \mathbf{u}_i \mathbf{v}_i^*$ with probability σ_i
- Corollary 3 gives error bound for $\hat{\mathbf{A}} = \frac{1}{r} \sum_{r=1}^r \mathbf{Z}_r$:

$$\mathbb{E} \|\mathbf{A} - \hat{\mathbf{A}}\| \leq \sqrt{\frac{2\sigma_1 \log(2n)}{r}} + \frac{2 \log(2n)}{3r}$$

- Mirsky's Theorem (1969)**. The best rank- r approximation satisfies

$$\min_{\text{rank}(\mathbf{B})=r} \|\mathbf{A} - \mathbf{B}\| = \sigma_{r+1} \leq \frac{1}{r+1}$$

- Sampling is really (really) suboptimal!**

Two-stage Low-rank Approximations

Low-Rank Approximations

🐛 **Goal:** Given a large matrix A , find low-rank factors:

$$\begin{matrix} & n & & & r & n & \\ & & & & & & \\ m & \boxed{A} & \approx & m & \boxed{B} & \boxed{C} & r \end{matrix}$$

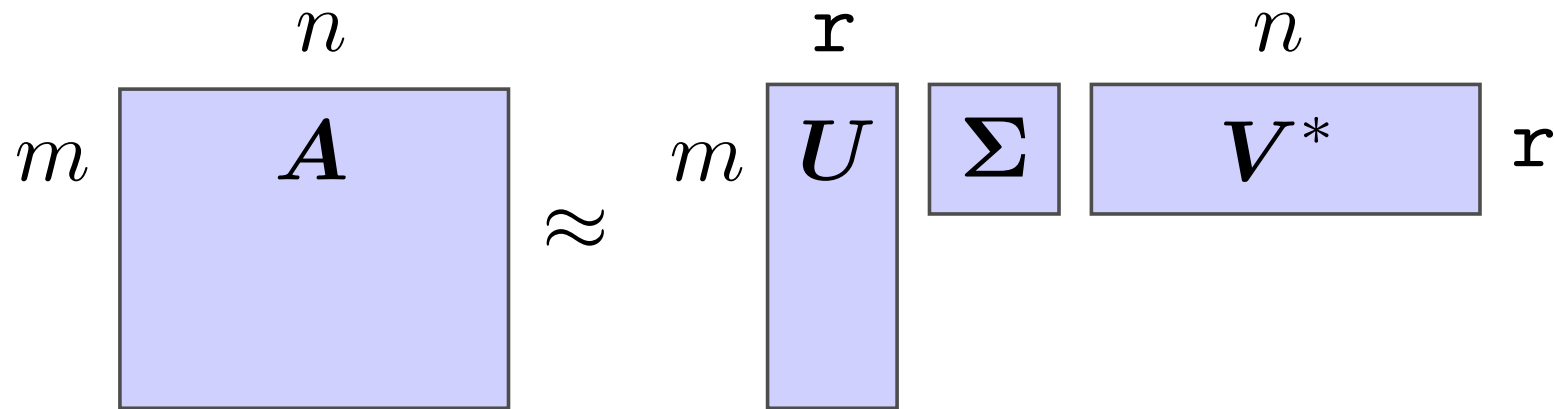
🐛 **Solving this problem gives algorithms for computing**

- 🐛 Leading singular vectors of a general matrix
- 🐛 Leading eigenvectors of a symmetric matrix
- 🐛 Spanning sets of rows or columns

🐛 We will focus on controlling backward error: $\|A - BC\| \leq \text{tol}$

Example: Truncated SVD

$A \approx U\Sigma V^*$ where U, V have orthonormal columns and Σ is diagonal:



Interpretation: r -truncated SVD = optimal rank- r approximation

Applications:

- ☞ Least-squares computations
- ☞ Principal component analysis
- ☞ Summarization and data reduction
- ☞ ...

Overview of Two-Stage Randomized SVD

Goal: Construct a truncated SVD of an input matrix A

Stage A: Finding the range

- Use a *randomized algorithm* to compute a low-dimensional basis Q that captures most of the action of A :

$$Q \text{ has orthonormal columns and } A \approx QQ^* A.$$

Stage B: Forming the decomposition

- Use the basis Q to reduce the problem size
- Apply *classical linear algebra* to reduced problem
- Obtain truncated SVD in factored form

[Ref] Halko et al. 2011, §1.

Stage A: Finding the Range

Given:

- An $m \times n$ matrix A
- Target rank $r_0 \ll \min\{m, n\}$
- Actual rank $r = r_0 + s$ where s is a small oversampling parameter

Construct an $m \times r$ matrix Q with orthonormal columns s.t.

$$\|A - QQ^*A\| \approx \min_{\text{rank}(B)=r_0} \|A - B\|,$$

- QQ^* is the orthogonal projector onto the range of Q

Approach: Use a randomized algorithm!

Total Cost: One multiply $(m \times n \times r) + \mathcal{O}(r^2n)$ flops

Stage B: Forming the SVD

Assume A is $m \times n$ and Q is $m \times r$ with ON columns and $A \approx QQ^*A$

Approach: Reduce problem size; apply classical numerical linear algebra!

1. Form $r \times n$ matrix $B = Q^*A$
2. Factor $B = U\Sigma V^*$
3. Conclude $A \approx (QU)\Sigma V^*$

$$\begin{aligned} A &\approx Q Q^* A = Q Q^* A \\ &= Q U \Sigma V^* = QU \Sigma V^* \end{aligned}$$

Total Cost: One multiply $(m \times n \times r) + \mathcal{O}(r^2n)$ flops

Total Costs for Truncated SVD

Two-Stage Randomized Algorithm:

$$2 \text{ multiplies } (m \times n \times r) + r^2(m + n) \text{ flops}$$

Classical Sparse Methods (Krylov):

$$r_0 \text{ multiplies } (m \times n \times 1) + r_0^2(m + n) \text{ flops}$$

Classical Dense Methods (RRQR + full SVD):

$$\text{Not based on multiplies} + r_0 mn \text{ flops}$$

Roadmap

1. How to use the range information to construct matrix approximations
2. How to find a subspace that captures the range of the matrix
3. Other types of sampling schemes
4. How to get high precision approximations by iteration

From Range to Decomposition

Interpolative Approximations

- *Interpolative approximation* represents matrix via a few rows (or cols)
- Takes the form $\mathbf{A} \approx \mathbf{W} \mathbf{A}_I$: where
 - I is a small set of row indices
 - $\mathbf{W}_I = \mathbf{I}$
 - All entries of \mathbf{W} have magnitude ≤ 2
- Can be constructed efficiently with rank-revealing QR
- Useful when we must **preserve meaning of rows** (or cols)

Assume \mathbf{A} is $m \times n$ and \mathbf{Q} is $m \times r$ and $\mathbf{A} \approx \mathbf{Q} \mathbf{Q}^* \mathbf{A}$

1. Form interpolative decomposition $\mathbf{Q} = \mathbf{W} \mathbf{Q}_I$: with $|I| = r$
2. Conclude $\mathbf{A} \approx \mathbf{W} \mathbf{A}_I$:

Total cost: $\mathcal{O}(r^2 m)$ flops

[Ref] Gu & Eisenstat 1996; Goreinov et al. 1997; Cheng et al. 2005. See Halko et al. 2011, §3.2.3 and §5.2.

Matrix Skeletons

- A *skeleton approximation* represents a matrix via a small submatrix
- Takes the form $A \approx W A_{IJ} X$ where
 - I is a small set of row indices; J is a small set of column indices
 - $W_{I,:} = \mathbf{I}$ and $X_{:,J} = \mathbf{I}$
 - All entries of W and X have magnitude ≤ 2
- Useful when we must preserve meaning of rows and columns

Assume A is $m \times n$ and Q is $m \times r$ and $A \approx QQ^* A$

1. Form row interpolative decomposition $Q = W Q_I$ with $|I| = r$
2. Form column interpolative decomposition $A_{I,:} = A_{IJ} X$ with $|J| = r$
3. Conclude $A \approx W A_{IJ} X$

Total cost: $\mathcal{O}(r^2(m + n))$ flops

[Ref] Gu & Eisenstat 1996; Goreinov et al. 1997; Cheng et al. 2005. See Halko et al. 2011, §3.2.3 and §5.2.

CUR Approximations

- A *CUR approximation* represents a matrix via a few rows and columns
- Takes the form $\mathbf{A} \approx \mathbf{A}_{:J} \mathbf{T} \mathbf{A}_I$: where I and J are small
- Useful when we must preserve meaning of rows and columns

Assume \mathbf{A} is $m \times n$ and \mathbf{Q} is $m \times r$ and \mathbf{P} is $n \times r$ and

$$\mathbf{A} \approx \mathbf{Q} \mathbf{Q}^* \mathbf{A} \quad \text{and} \quad \mathbf{A} \approx \mathbf{A} \mathbf{P} \mathbf{P}^*$$

1. Form row interpolative decomposition $\mathbf{Q} = \mathbf{W} \mathbf{Q}_I$: with $|I| = r$
2. Form column interpolative decomposition $\mathbf{P} = \mathbf{P}_{:J} \mathbf{X}$ with $|J| = r$
3. Compute the pseudoinverse $\mathbf{T} = (\mathbf{A}_{IJ})^\dagger$
4. Conclude $\mathbf{A} \approx \mathbf{A}_{:J} \mathbf{T} \mathbf{A}_I$:

Total cost: $\mathcal{O}(r^2(m + n))$ flops

[Ref] Goreinov et al. 1997; Drineas et al. 2009; Mahoney & Drineas 2009; Bien et al. 2010.

Nyström Approximation for PSD Matrices

🦉 *Nyström approximations* represent psd matrices (much) more accurately

Assume A is $n \times n$ psd and Q is $n \times r$ with ON columns and $A \approx QQ^*A$

Nyström approximation:

$$A \approx (AQ)(Q^*AQ)^\dagger(AQ)^*$$

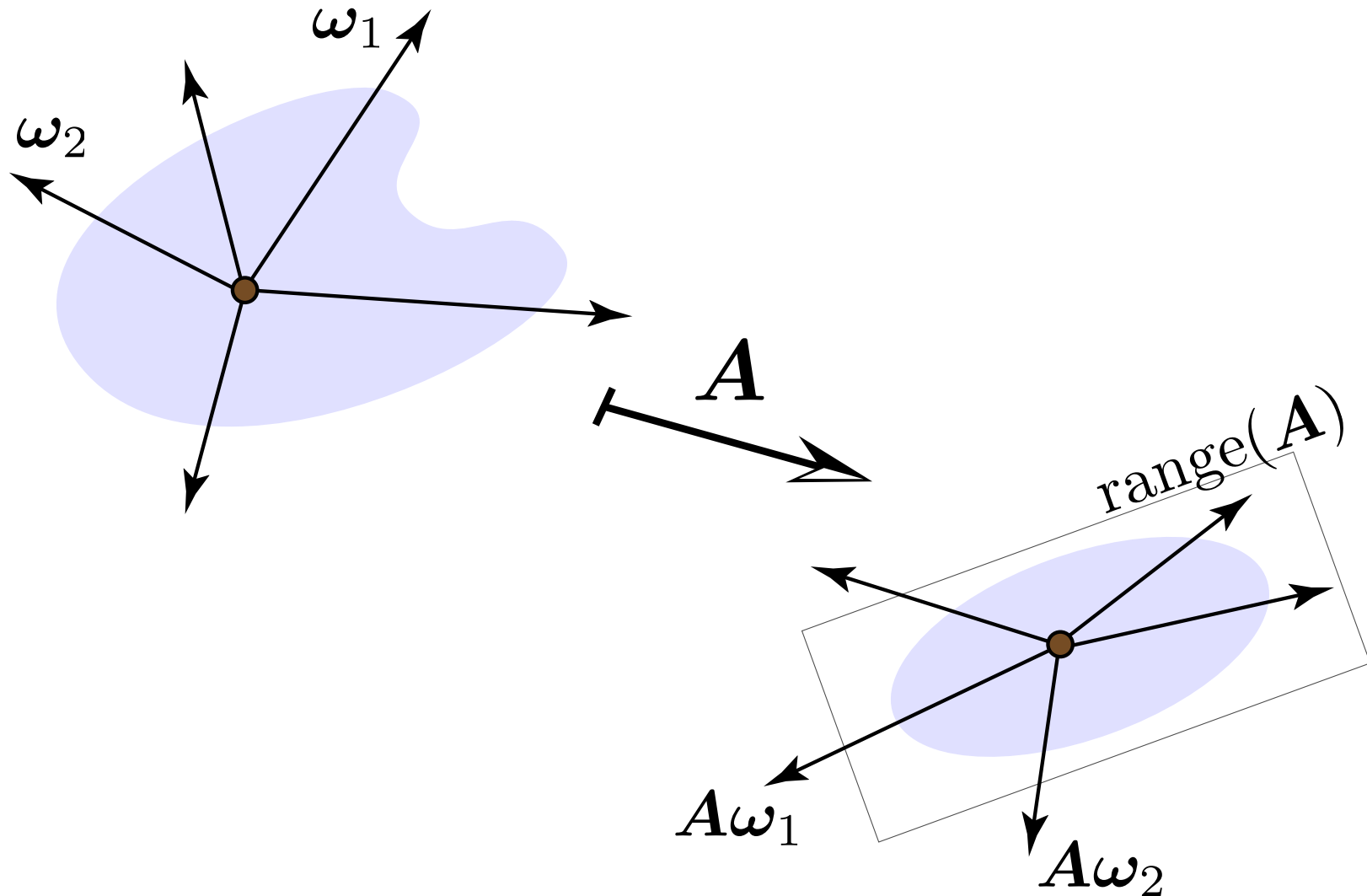
1. Form $n \times n \times r$ product $F_0 = AQ$
2. Form $r \times n \times r$ product $F = Q^*F_0$
3. Compute $r \times r$ Cholesky decomposition $F = T^*T$
4. Form $n \times r$ matrix $H = F_0T^{-1}$ by triangular solve
5. Conclude $A \approx HH^*$

Total cost: One multiply ($n \times n \times r$) + $\mathcal{O}(r^2n)$ flops

[Ref] Williams & Seeger 2002; Drineas & Mahoney 2005; Halko et al. 2011, §5.4; Gittens 2011; Gittens & Mahoney 2013.

Randomized Range Finder

Randomized Range Finder: Intuition



Prototype for Randomized Range Finder

Input: An $m \times n$ matrix A , number r of samples

Output: An $m \times r$ matrix Q with orthonormal columns

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

Total Cost: 1 multiply $(m \times n \times r) + \mathcal{O}(r^2n)$ flops

[Refs] NLA community: Stewart (1970s). GFA: Johnson & Lindenstrauss (1984) et seq. TCS: Boutsidis, Deshpande, Drineas, Frieze, Kannan, Mahoney, Papadimitriou, Sarlós, Vempala (1998–present). SciComp: Martinsson, Rokhlin, Szlam, Tygert (2004–present). See Halko et al. 2011, §2 for more history.

Implementation Issues

Q: How do we pick the number r of samples?

A: Adaptively, using a randomized error estimator.

Q: How does the number r of samples compare with the target rank r_0 ?

A: Remarkably, $r = r_0 + 5$ or $r = r_0 + 10$ is usually adequate!

Q: What random matrix Ω ?

A: For many applications, standard Gaussian works brilliantly.

Q: How do we perform 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.

[Ref] Halko et al. 2011, §4.

Simple Error Bound for Random Range Finder

Theorem 4. [HMT 2011] Assume

- The matrix \mathbf{A} is $m \times n$ with $m \geq n$
- The target rank is r_0
- The optimal error $\sigma_{r_0+1} = \min_{\text{rank}(\mathbf{B})=r_0} \|\mathbf{A} - \mathbf{B}\|$
- The test matrix $\mathbf{\Omega}$ is $n \times (r_0 + s)$ *standard Gaussian*

Then the random range finder yields an $(r_0 + s)$ -dimensional basis \mathbf{Q} with

$$\mathbb{E} \|\mathbf{A} - \mathbf{Q}\mathbf{Q}^* \mathbf{A}\| \leq \left[1 + \frac{4\sqrt{r_0 + s}}{s - 1} \cdot \sqrt{n} \right] \sigma_{r_0+1}$$

The probability of a substantially larger error is negligible.

[Ref] Halko et al. 2011, §10.2.

Error Bound for Random Range Finder

Theorem 5. [HMT 2011] Assume

- The matrix \mathbf{A} is $m \times n$
- The target rank is r_0
- The singular values of \mathbf{A} are $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots$
- The test matrix $\mathbf{\Omega}$ is $n \times (r_0 + s)$ *standard Gaussian*

Then the random range finder yields an $(r_0 + s)$ -dimensional basis \mathbf{Q} with

$$\mathbb{E} \|\mathbf{A} - \mathbf{Q}\mathbf{Q}^* \mathbf{A}\| \leq \left[1 + \sqrt{\frac{r_0}{s-1}} \right] \sigma_{r_0+1} + \frac{e\sqrt{r_0+s}}{s} \left(\sum_{j>r_0} \sigma_j^2 \right)^{1/2}$$

The probability of a substantially larger error is negligible.

This bound is essentially optimal.

[Ref] Halko et al. 2011, §10.2.

Adaptive Error Estimation

🐼 **In practice, we do not know the target rank!**

🐼 A priori error bounds are motivating—but not so useful

🐼 **Solution:** Estimate the error, and stop when we reach a target

🐼 Let $\omega_1, \dots, \omega_{10}$ be iid standard Gaussian vectors

$$\text{err_est} := \max_{1 \leq j \leq 10} \|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)\mathbf{A}\omega_j\|$$

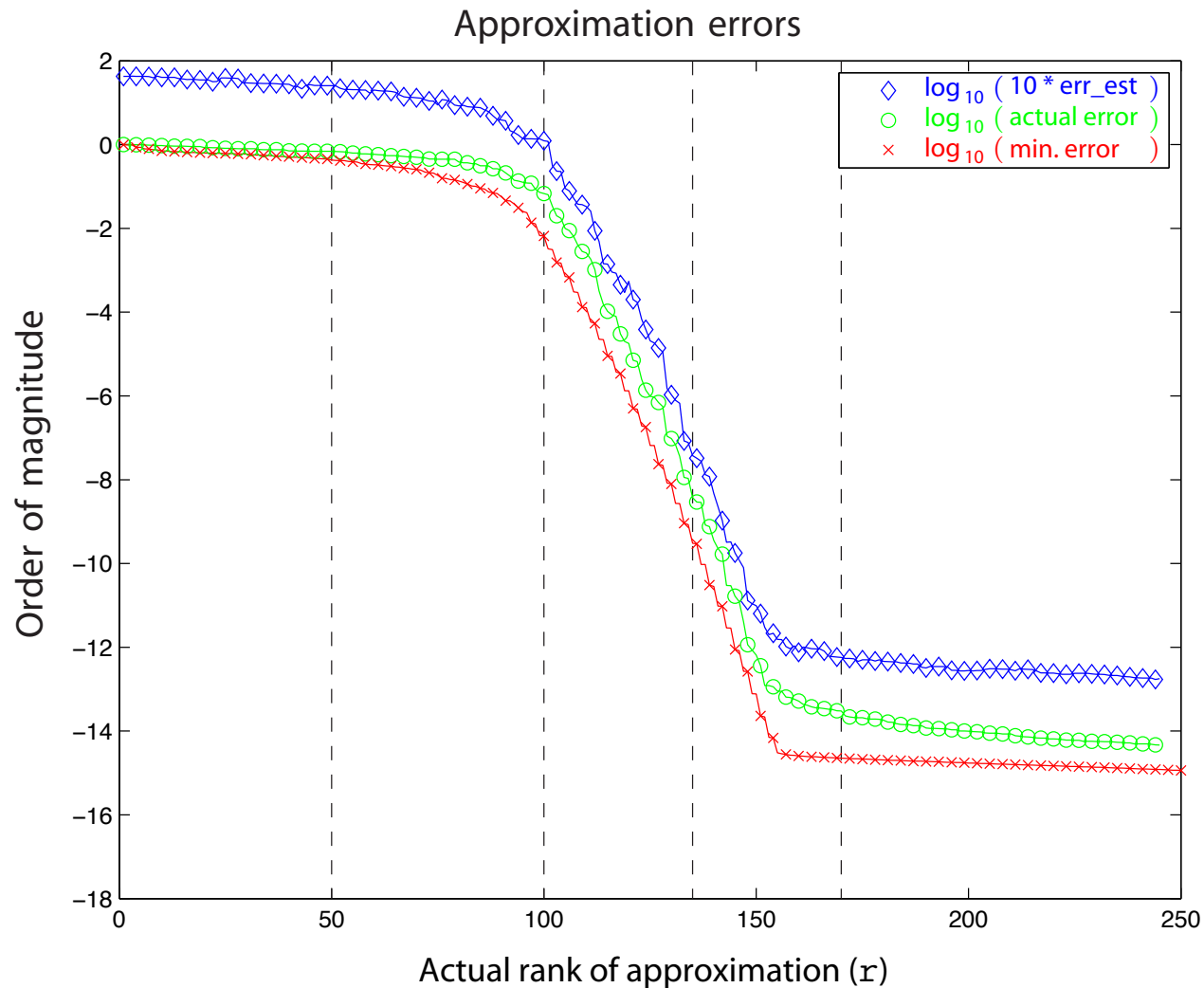
🐼 How good is this estimate?

$$\mathbb{P} \{ \|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)\mathbf{A}\| \geq 10 \cdot \text{err_est} \} \leq 10^{-10}$$

🐼 Can modify the random range finder to form `err_est` (almost) for free

[Refs] Woolfe et al. 2008. See Halko et al. 2011, §4.4.

Approximating a Helmholtz Integral Operator



Other Sampling Strategies

Other Sampling Techniques?

🐼 In randomized range finder, can change distribution of test matrix Ω

🐼 **Why?**

🐼 To exploit Ω with fast multiply

🐼 To avoid dense $Y = A\Omega$ when A is sparse

🐼 **When are these techniques worth considering?**

🐼 The matrix A *must* have a rapidly decaying spectrum

🐼 The dominant right singular vectors of the matrix A are flat

🐼 **Executive summary:**

🐼 Sampling with an SRFT is effective when A has spectral decay

🐼 **I do not recommend the other schemes, except in special cases**

Deterministic Error Bound for Range Finder

• A is $m \times n$ with SVD partitioned as

$$A = U \begin{bmatrix} \Sigma_1 & \\ & \Sigma_2 \end{bmatrix} \begin{bmatrix} V_1^* \\ V_2^* \end{bmatrix} \begin{matrix} r_0 \\ n - r_0 \end{matrix}$$

• Let Ω be any test matrix, decomposed as

$$\Omega_1 = V_1^* \Omega \quad \text{and} \quad \Omega_2 = V_2^* \Omega.$$

• Construct the sample matrix $Y = A\Omega$ and a basis Q for $\text{range}(Y)$

Theorem 6. [BMD11; HMT11] *When Ω_1 has full row rank,*

$$\|A - QQ^*A\|^2 \leq \|\Sigma_2\|^2 + \|\Sigma_2\Omega_2\Omega_1^\dagger\|^2.$$

Roughly: *Range finder works when Ω_1 has well-conditioned rows.*

[Refs] Boutsidis et al. 2011. See Halko et al. 2011, §9.

Leverage Score Sampling

- Let V_1^* be the $r_0 \times n$ matrix of top right singular vectors of A
- The squared column norms ℓ_1, \dots, ℓ_n of V_1^* are called *leverage scores*
- We can obtain a sampling distribution by normalizing: $p_i = \ell_i / r_0$
- Let the number r of samples satisfy $r \sim r_0 \log r_0$
- Let $\omega = e_i$ with probability p_i
- Test matrix Ω has r columns, each an independent copy of ω
- Then $\Omega_1 = V_1^* \Omega$ is likely to have well-conditioned rows
- Leverage score sampling respects columns
- **But...** Cost of approximating leverage scores is $\mathcal{O}(mn \log m)$
- **Better method:** Use SRFT (p. 63) + interpolative approximation

$e_i = i$ th standard basis vector

[Ref] See Mahoney 2011 for details about leverage score sampling.

Simple Column Sampling

- **Assume** that the columns of V_1^* have comparable norms
- Then we can take a sampling distribution $p_i = 1/n$ for each i
- Let the number r of samples satisfy $r \sim r_0 \log r_0$
- Let $\omega = e_i$ with probability p_i
- Test matrix Ω has r columns, each an independent copy of ω
- Then $\Omega_1 = V_1^* \Omega$ is likely to have well-conditioned rows
- Simple random sampling respects columns and is basically free
- Works when leverage scores are approximately uniform
- **But...** the approach fails when leverage scores are nonuniform

[Refs] See Chen & Demanet 2009; Gittens 2009; Gittens & Mahoney 2013 for analysis of simple random sampling.

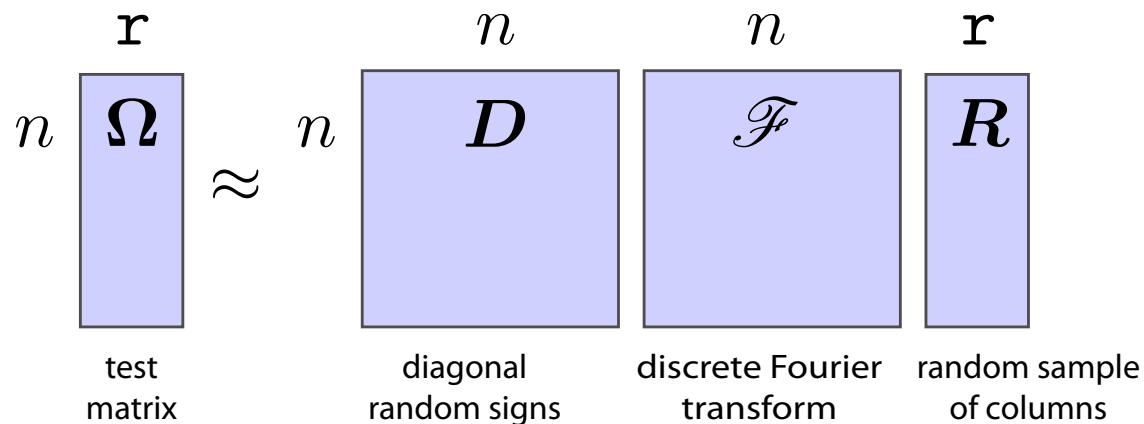
Sparse Projections

- 🐼 **Assume** that the matrix A is very sparse
- 🐼 May be valuable to try to preserve this sparsity
- 🐼 Consider test matrix Ω has $\log^{\mathcal{O}(1)}(n)$ random ± 1 entries per column
- 🐼 The number r of columns in Ω has order $r_0 \log^{\mathcal{O}(1)}(n)$
- 🐼 Then $\Omega_1 = V_1^* \Omega$ is likely to have well-conditioned rows
- 🐼 Sparse projections respect sparsity and are relatively cheap
- 🐼 **But...** this approach tends not to work well for sparse matrices
- 🐼 **Ironically,** may be more valuable for dense matrices

[Refs] Clarkson & Woodruff 2012; Nelson & Nguyen 2012; Meng & Mahoney 2013; Bourgain & Nelson 2013.

Structured Random Matrices

- Choose a *structured* random matrix Ω to accelerate multiply
- Example:** Subsampled randomized Fourier transform (SRFT)



- Intuition:** The transform $D\mathcal{F}$ uniformizes leverage scores
- Cost of forming $Y = A\Omega$ by FFT is $\mathcal{O}(mn \log n)$
- In practice, if A has spectral decay, SRFT works as well as Gaussian

[Refs] Ailon & Chazelle 2006, 2009; Rokhlin & Tygert 2008; Liberty 2009; T 2011. See Halko et al. 2011, §4.6.

An Error Bound for the SRFT

Theorem 7. [Boutsidis & Gittens 2013] Assume

- The matrix \mathbf{A} is $m \times n$
- The target rank is r_0
- The singular values of \mathbf{A} are $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots$
- The test matrix $\mathbf{\Omega}$ is an $n \times r$ SRFT with

$$r \gtrsim (r_0 + \log n) \log r_0$$

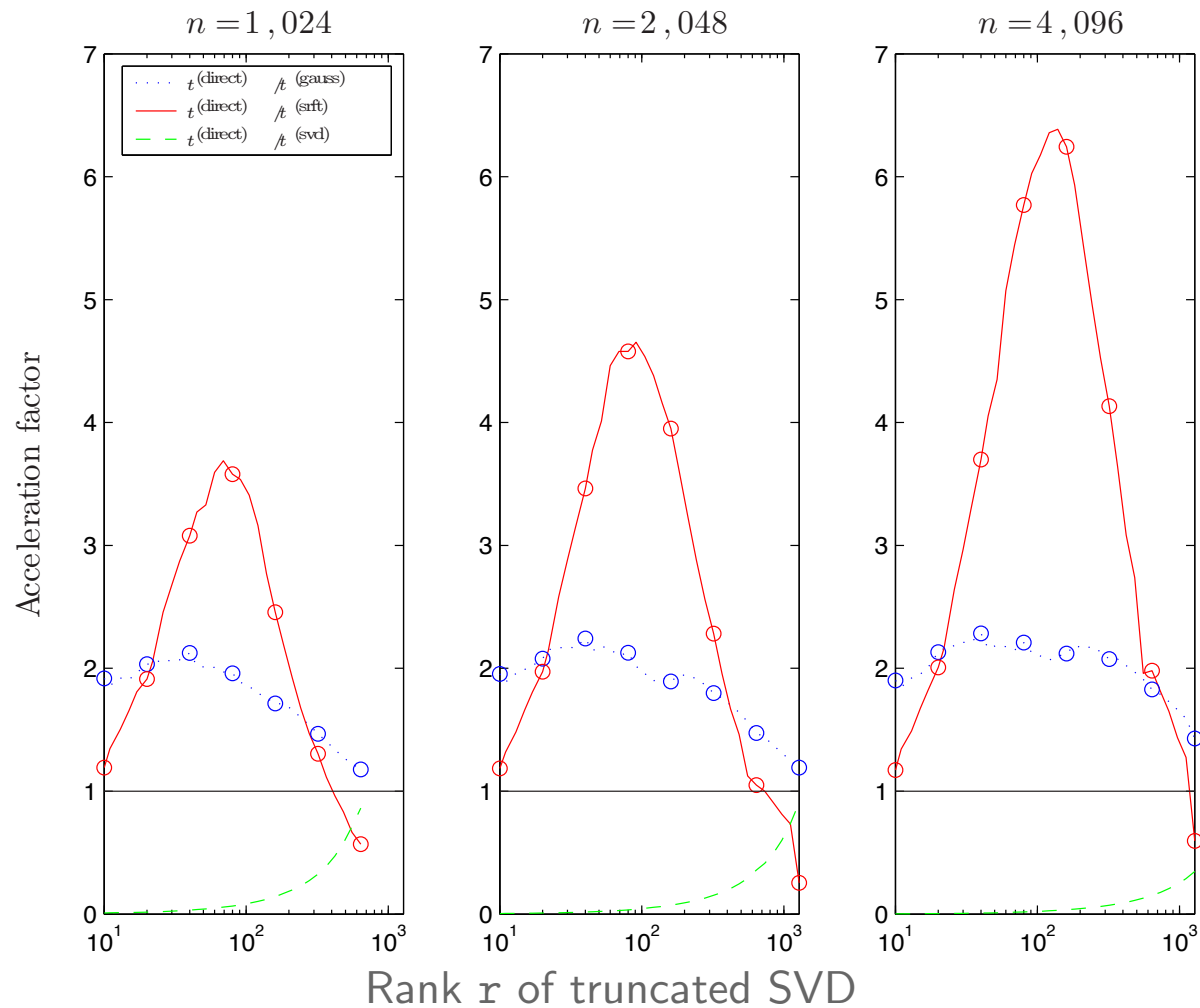
Then

$$\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^*\mathbf{A}\| \lesssim \left[1 + \frac{\log n}{\sqrt{r}}\right] \sigma_{r_0+1} + \sqrt{\frac{\log n}{r}} \cdot \left(\sum_{j>r_0} \sigma_j^2\right)^{1/2}$$

except with probability $\mathcal{O}(r_0^{-1})$

[Refs] T 2011; Halko et al. 2011; Boutsidis & Gittens 2013.

SRFT: Faster SVD when Spectrum Decays



[Ref] Halko et al. 2011, §7.4.

Randomized Power Scheme

The Role of Spectral Decay

- The random range finder works when the spectrum of \mathbf{A} decays quickly (Helmholtz)
- **Problem:** This behavior is not common in data analysis problems
- **Examples:**
 - Matrix is contaminated with noise: $\mathbf{A} = \mathbf{A}_0 + \mathbf{N}$ with $\|\mathbf{N}\| \propto \|\mathbf{A}_0\|$
 - Matrix spectrum decays slowly: $\sigma_j \sim j^{-\alpha}$ for $\alpha < \frac{1}{2}$
- In these settings, the random range finder **will give unreliable results**
- **Remedy:** Use subspace iteration

Randomized Range Finder + Power Scheme

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

Remedy: Apply the randomized range finder to $(\mathbf{A}\mathbf{A}^*)^q\mathbf{A}$ for small q

Input: An $m \times n$ matrix \mathbf{A} , number r of samples

Output: An $m \times r$ matrix \mathbf{Q} with orthonormal columns

1. Draw an $n \times r$ random matrix $\mathbf{\Omega}$.
2. Carefully form the matrix product $\mathbf{Y} = (\mathbf{A}\mathbf{A}^*)^q\mathbf{A}\mathbf{\Omega}$.
3. Construct an orthonormal basis \mathbf{Q} for the range of \mathbf{Y} .

Total Cost: $2q + 1$ multiplies $(m \times n \times r) + \mathcal{O}(qr^2n)$

[Refs] Stewart 1970s; Roweis 1997; Gu 2007; Rokhlin et al. 2008. See Halko et al. 2011, §4.5.

Implementation Issues

Q: How do we pick the number r of samples?

A: Adaptively, using a randomized error estimator.

Q: How do we pick the number q of power iterations?

A: Remarkably, $q = 2$ or $q = 3$ is usually adequate!

Q: What random matrix Ω ?

A: Standard Gaussian. Other sampling distributions have limited value.

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

A: Alternately multiply by A and A^* .

Q: How do we compute the orthonormal basis?

A: Carefully... Perform QR factorization at each step.

[Ref] Halko et al. 2011, §4.

Simple Error Bound for Power Scheme

Theorem 8. [HMT 2011] Assume

- The matrix \mathbf{A} is $m \times n$ with $m \geq n$
- The target rank is r_0
- The optimal error $\sigma_{r_0+1} = \min_{\text{rank}(\mathbf{B}) \leq r_0} \|\mathbf{A} - \mathbf{B}\|$
- The test matrix $\mathbf{\Omega}$ is $n \times (r_0 + s)$ *standard Gaussian*

Then the basis \mathbf{Q} computed by the power scheme satisfies

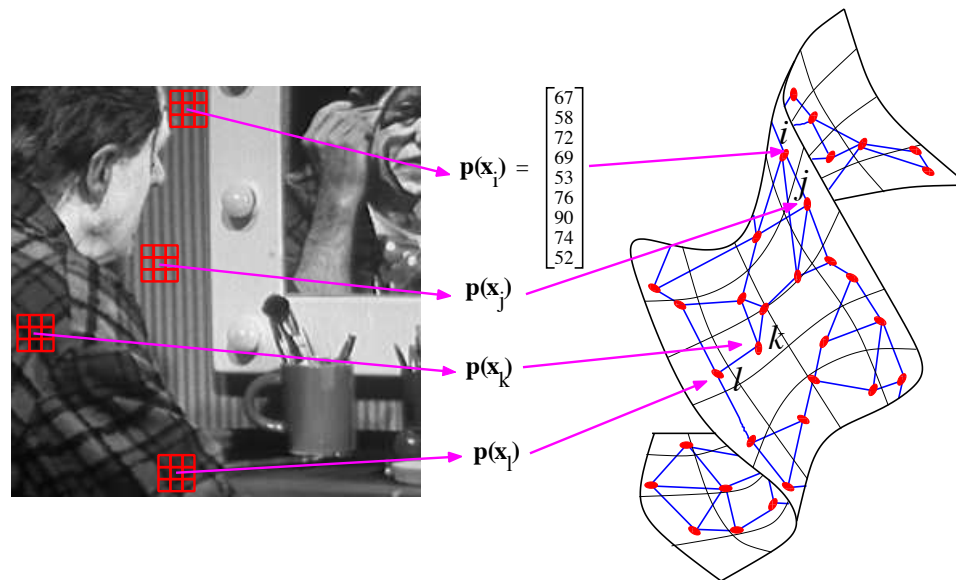
$$\mathbb{E} \|\mathbf{A} - \mathbf{Q}\mathbf{Q}^* \mathbf{A}\| \leq \left[1 + \frac{4\sqrt{r_0 + s}}{s - 1} \cdot \sqrt{n} \right]^{1/(2q+1)} \sigma_{r_0+1}$$

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

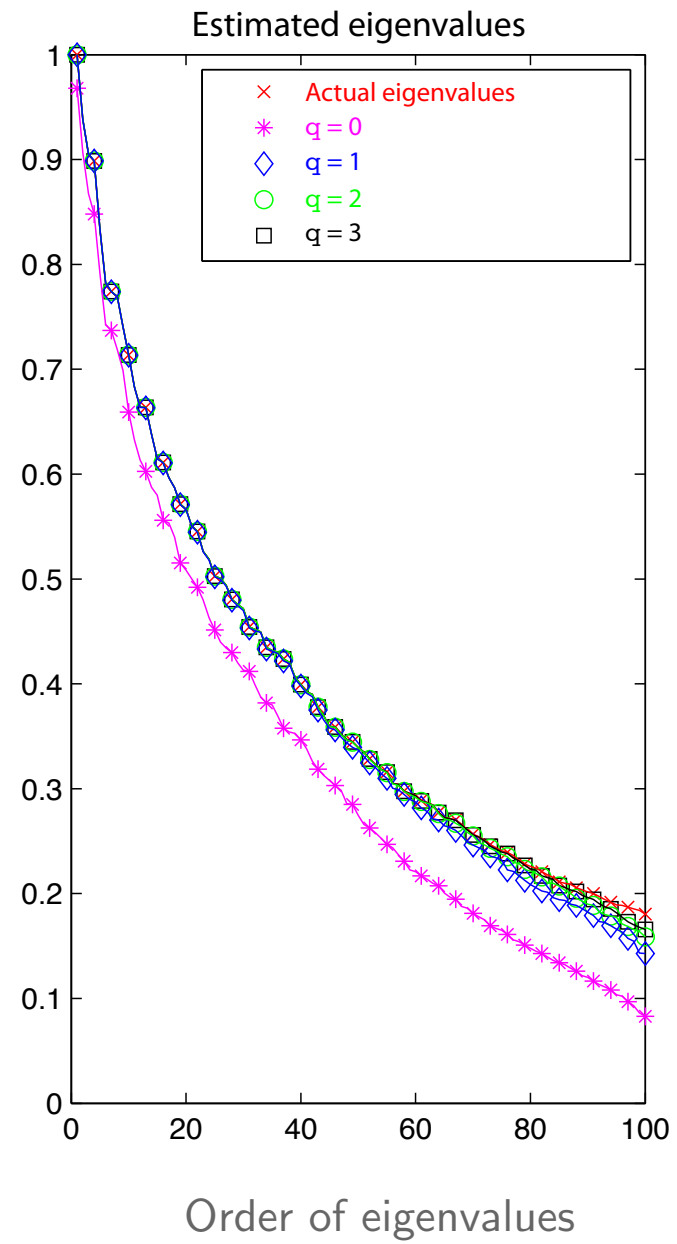
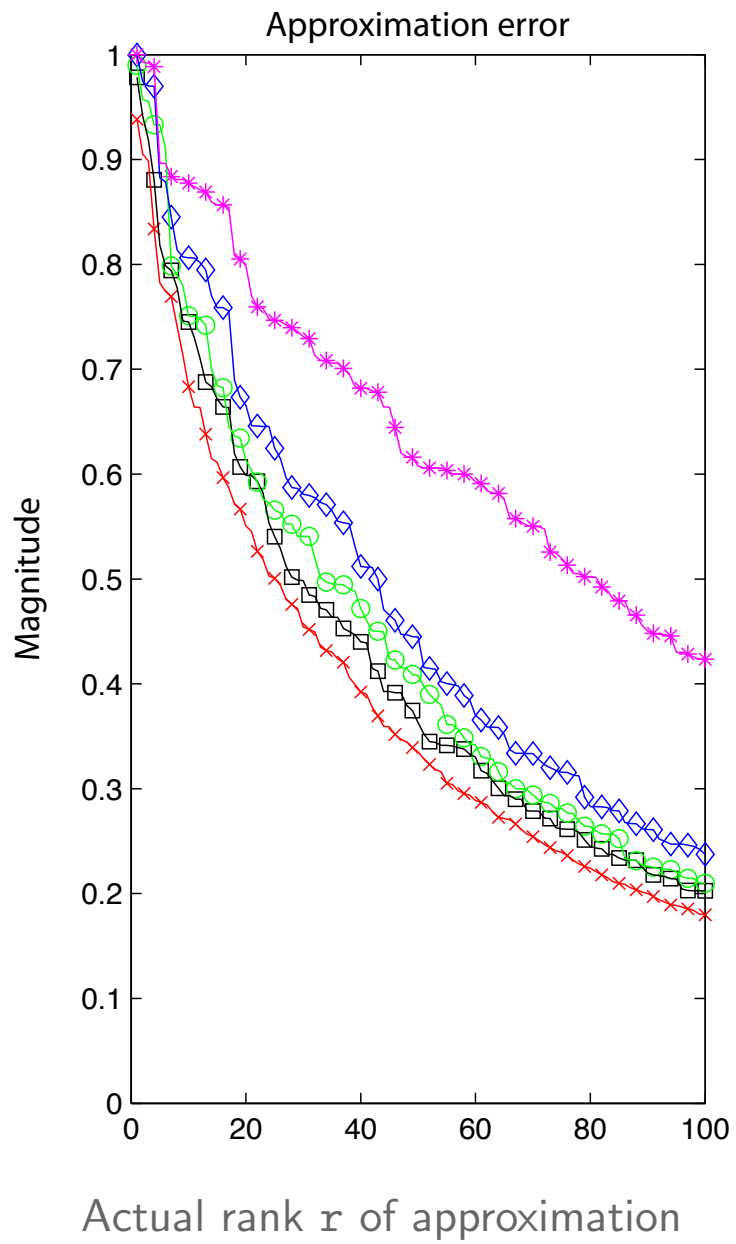
[Ref] Halko et al. 2011, §10.

A Graph Laplacian

- 🐼 Graph Laplacian on a point cloud of 9×9 image patches
- 🐼 Form $9,025 \times 9,025$ symmetric data matrix A
- 🐼 **Total storage:** 311 Megabytes (uncompressed)
- 🐼 The eigenvectors give information about graph structure
- 🐼 Attempt to compute first 100 eigenvalues/vectors using power scheme



[Ref] Courtesy of Gunnar Martinsson & François Meyer.

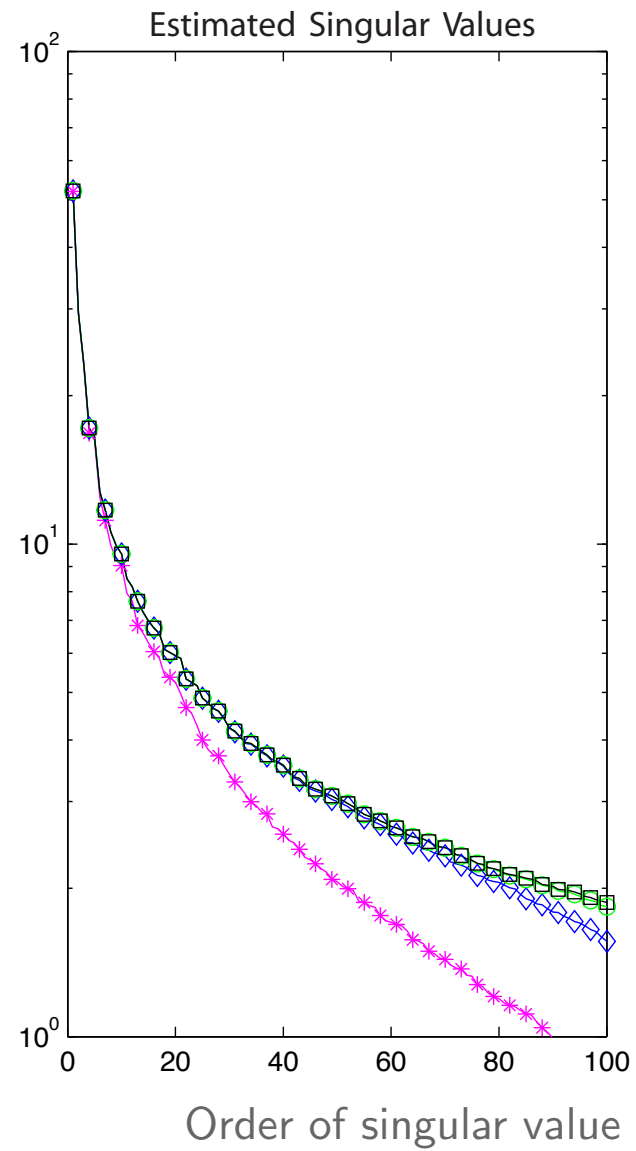
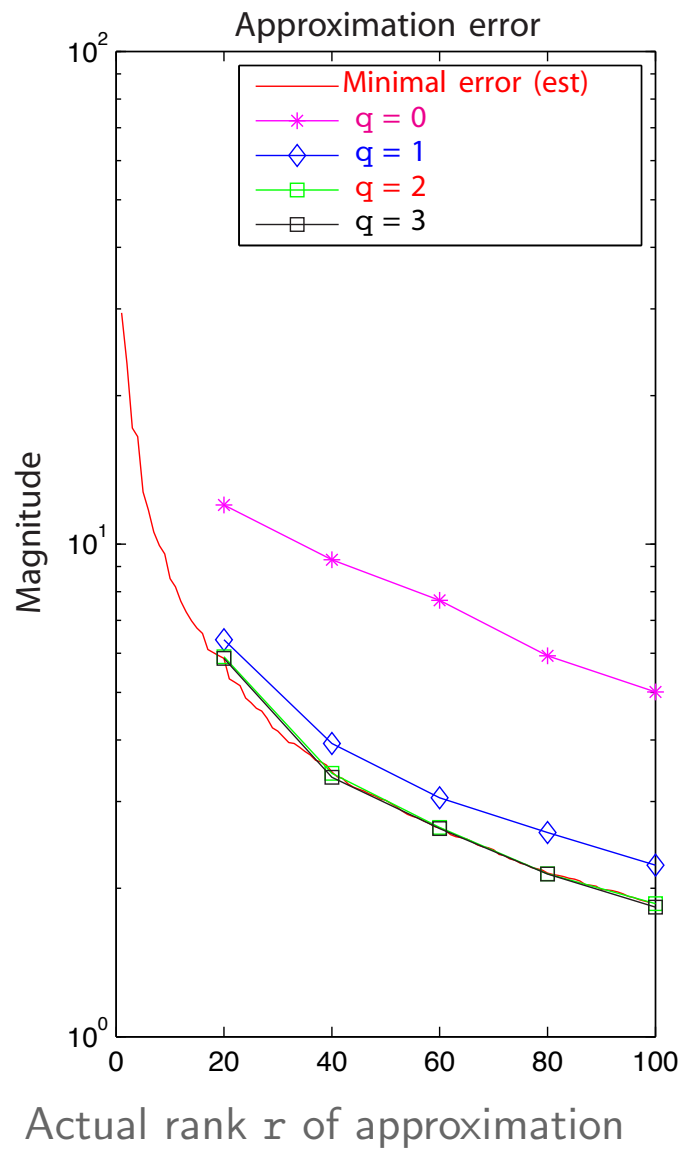


Eigenfaces

- 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 of A are called **eigenfaces**
- Attempt to compute first 100 eigenfaces using power scheme



[Ref] Image from Scholarpedia article “Eigenfaces,” accessed 12 October 2009



[Ref] Halko et al. 2011, §7.3.

Resources

Articles with Broad Scope

- 🐼 M. Mahoney, “Randomized algorithms for matrices and data,” *Foundations & Trends in Machine Learning*, 2011
- 🐼 N. Halko, P.-G. Martinsson, and J. A. Tropp, “Finding structure with randomness...,” *SIAM Rev.*, 2011
- 🐼 J. A. Tropp, “User-Friendly Tools for Random Matrices.” Submitted to *Foundations & Trends in Machine Learning*, 2014

High-Quality Software for Randomized NLA

🐛 [Mark Tygert.](#)

<http://tygert.com/software.html>

🐛 [Haim Avron.](#)

http://researcher.watson.ibm.com/researcher/view_group.php?id=5131

🐛 [Xiangrui Meng.](#)

<http://www.stanford.edu/~mengxr/pub/lsrcn.htm>

To learn more...

E-mail: `jtropp@cms.caltech.edu`

Web: `http://users.cms.caltech.edu/~jtropp`