# SketchySVD

## Joel A. Tropp

Computing + Mathematical Sciences

California Institute of Technology

`jtropp@cms.caltech.edu`

Coauthors: **Alp Yurtsever** (EPFL), **Volkan Cevher** (EPFL);

**Yiming Sun** (Cornell), **Yang Guo** (UWisc), **Charlene Luo** (Columbia), **Madeleine Udell** (Cornell)

Thanks: **Gunnar Martinsson** (UT–Austin), **Mark Tygert** (Facebook)

# SIAM Journal on Mathematics of Data Science (SIMODS)

https://simods.siam.org/

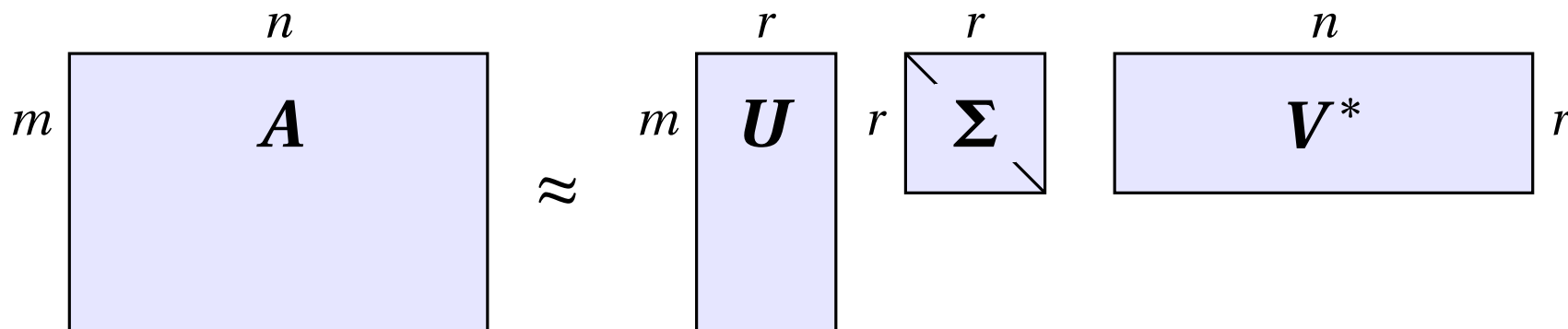**Editors:** Tammy Kolda (EIC); Al Hero, Mike Jordan, Rob Nowak, Joel Tropp

# First SIAM Conference on Mathematics of Data Science (MDS 2020)

## 5–7 May 2020
## Cincinnati, Ohio, USA

**Co-Chairs:** Gitta Kutyniok, Ali Pinar, Joel A. Tropp

# The Famous Truncated SVD

# Truncated Singular Value Decomposition (TSVD)



$$A \approx U \Sigma V^*$$

- $U, V$ have orthonormal columns and $\Sigma$ is nonnegative diagonal
- Approximately $r(m+n)$ degrees of freedom

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

**Applications:**

- Least-squares computations (linear regression)
- Principal component analysis (orthogonal regression; total least squares)
- Summarization, data reduction, visualization, ...

# A Paean to the Truncated SVD

"[Truncated SVD] is one of the few methods that has solid foundations and can be trusted, provided that the computations are correct. Having something reliable in machine learning is worth its weight in gold — there are almost no gold standards in the field, precluding rapid progress. Think of what happens to machine learning when the routine for matrix–vector multiplication doesn't always work right. You can't debug any code, much less a big system consisting of many algorithms thrown together."

**–Nemo**

# Modern Numerical Linear Algebra

# What's Wrong with Classical TSVD Algorithms?

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

**Climate Change:**

❧ Medium- to large-scale data (Megabytes+)
❧ New architectures (multi-core, distributed, data centers, ...)

❧ **Today:** New data presentations (dynamic, off-core, streaming)

**Engineering:**

❧ Theoretically, we already know how to do streaming TSVD, but...
❧ Current "algorithms" are not ready for implementation
❧ For scientific applications, high accuracy is essential!

❧ **Today:** First practical algorithms for streaming TSVD

# Streaming Linear Algebra

**The Turnstile Model:**

$$A \quad = \quad H_1 + H_2 + H_3 + H_4 + \cdots$$

- ✏️ Huge input matrix $A$ is presented as a sum of innovations $H_i$
- ✏️ Must discard each innovation $H_i$ after it is processed

- ✏️ **Goal:** Without storing $A$ in full, return TSVD after seeing all updates

**Applications:**

- ✏️ Scientific simulation and data collection
- ✏️ One-pass approximation of matrix stored out-of core
- ✏️ Large-scale semidefinite programming algorithms

# Randomized Linear Sketches

$$\text{sketch} \quad = \quad \mathcal{L}(\boldsymbol{A}) \quad = \quad \sum_i \mathcal{L}(\boldsymbol{H}_i)$$

- �expl] Select a **linear** map $\mathcal{L}$ without reference to $\boldsymbol{A}$
- ✐ Sketch is much smaller than input matrix
- ✐ Use **randomness** so sketch works for an arbitrary input
- ✐ **[LNW14] Essentially the only way to handle the turnstile model!**

**Examples:**

- ✐ Left multiply: $\mathcal{L}(\boldsymbol{A}) = \boldsymbol{\Upsilon} \boldsymbol{A}$
- ✐ Right multiply: $\mathcal{L}(\boldsymbol{A}) = \boldsymbol{A} \boldsymbol{\Omega}^*$
- ✐ Select some entries: $\mathcal{L}(\boldsymbol{A}) = \{a_{ij} : (i, j) \in E\}$

**Sources:** Alon et al. 1996; Sarlós 2006; Muthukrishnan 2008; Woolfe et al. 2008; Clarkson & Woodruff 2009; HMT 2011; Mahoney 2012; Woodruff 2014; Li et al. 2014; Drineas & Mahoney 2016; TYUC 2016–2019; ....

# Overview of SketchySVD

# Sketch *est Omnis Divisa in Partes Tres*

- Let $A \in \mathbb{C}^{m \times n}$ be an input matrix (presented in turnstile model)

- Fix sketch size parameters $(k, s)$ with $r \leq k \leq s \ll \min\{m, n\}$

- Draw linear dimension reduction maps independently at random:

$$\mathbf{\Upsilon} : \mathbb{C}^m \to \mathbb{C}^k \quad \text{and} \quad \mathbf{\Omega} : \mathbb{C}^n \to \mathbb{C}^k$$

$$\mathbf{\Phi} : \mathbb{C}^m \to \mathbb{C}^s \quad \text{and} \quad \mathbf{\Psi} : \mathbb{C}^n \to \mathbb{C}^s$$

- **Co-range and range sketches:**

$$X = \mathbf{\Upsilon} A \in \mathbb{C}^{k \times n} \quad \text{and} \quad Y = A\mathbf{\Omega}^* \in \mathbb{C}^{m \times k}$$

- **Core sketch:**

$$Z = \mathbf{\Phi} A \mathbf{\Psi}^* \in \mathbb{C}^{s \times s}$$

**Sources:** J. Caesar ca. 50 BCE; Vempala et al. 1998–2000; Drineas et al. 2004–2006; Martinsson et al. 2004–2012; Sarlós 2006; Woolfe et al. 2008; Clarkson & Woodruff 2009, 2011; Boutsidis et al. 2011, 2016; HMT 2011; Nelson et al. 2012–2015; Woodruff 2014; Gu 2015; Cohen et al. 2015; Upadhyay 2016; TYUC 2016–2019....

# The SKETCHYSVD Procedure

1. Use range sketches $X, Y$ to find orthonormal $Q \in \mathbb{C}^{m \times k}$ and $P \in \mathbb{C}^{n \times k}$ where

$$A \quad \approx \quad QQ^* APP^*$$

2. Use core sketch $Z \in \mathbb{C}^{s \times s}$ to find core approximation $C \in \mathbb{C}^{k \times k}$ such that

$$C \quad \approx \quad Q^* AP$$

3. For $r \leq k$, apply classical or randomized TSVD algorithm to form

$$[\![C]\!]_r \quad = \quad U\Sigma V^*$$

4. Obtain approximate $r$-truncated SVD $\hat{A}_r$ in factored form:

$$\hat{A}_r \quad := \quad (QU)\Sigma(PV)^* \quad = \quad Q[\![C]\!]_r P^*$$

$$\approx \quad QCP^* \quad \approx \quad QQ^* APP^* \quad \approx \quad A$$

**Sources:** Vempala et al. 1998–2000; Drineas et al. 2004–2006; Martinsson et al. 2004–2012; Sarlós 2006; Woolfe et al. 2008; Clarkson & Woodruff 2009; Boutsidis et al. 2011, 2016; HMT 2011; Mahoney 2012; Nelson et al. 2012–2015; Woodruff 2014; Gu 2015; Cohen et al. 2015; Upadhyay 2016; TYUC 2016–2019....

# SKETCHYSVD: Pseudocode

**Input:** Sketch size parameters; input matrix $A \in \mathbb{C}^{m \times n}$ as a turnstile stream

**Output:** Rank-$r$ approximation $\hat{A}_r = U \Sigma V^*$

1  **function** INITIALIZE($m, n, k, s$)                                      ▷ Set up the sketch
2      Draw random dimension reduction maps $\Upsilon, \Omega, \Phi, \Psi$
3      $X \leftarrow 0$ and $Y \leftarrow 0$ and $Z \leftarrow 0$
4  **function** LINEARUPDATE($H$)                                   ▷ Process $A \leftarrow A + H$
5      $X \leftarrow X + \Upsilon H$
6      $Y \leftarrow Y + H\Omega^*$
7      $Z \leftarrow Z + \Phi H \Psi^*$
8  **function** SKETCHYSVD($r$)                               ▷ Compute $r$-truncated SVD
9      $Q \leftarrow \text{economy\_qr}(Y)$                               ▷ Basis for range
10     $P \leftarrow \text{economy\_qr}(X^*)$                       ▷ Basis for co-range
11     $C \leftarrow ((\Phi Q) \backslash Z) / (\Psi P)$                               ▷ Core matrix
12     $(U, \Sigma, V) \leftarrow \text{randsvd}(C; r)$                               ▷ Use [HMT11]
13     $U \leftarrow QU$ and $V \leftarrow PV$                       ▷ Consolidate unitary factors
14     **return** $(U, \Sigma, V)$

# SKETCHYSVD: Analysis

**Theorem 1** (TYUC 2018). **Assume**

- *The input matrix $A \in \mathbb{C}^{m \times n}$ and the sketch parameters satisfy $s \geq 2k$*
- *The dimension reduction maps are independent complex standard normal*

**Then** *SKETCHYSVD computes a rank-$k$ approximation $\hat{A}_k$ for which*

$$\mathbb{E}\left\| A - \hat{A}_k \right\|_F^2 \quad \leq \quad \frac{s}{s-k} \quad \cdot \quad \min_{\varrho < k} \quad \frac{k+\varrho}{k-\varrho} \cdot \left\| A - [\![A]\!]_\varrho \right\|_F^2$$

*In particular, when $k = (1 + \varepsilon^{-1})r$ and $s = (1 + \varepsilon^{-1})k$ for $\varepsilon \in (0, 1]$,*

$$\mathbb{E}\left\| A - \hat{A}_k \right\|_F^2 \quad \leq \quad (1 + 5\varepsilon) \cdot \left\| A - [\![A]\!]_r \right\|_F^2$$

- **Key Fact:** Approximation exploits spectral decay
- Related results hold for rank-$r$ approximation + with high probability

# Resource Usage with Sparse Maps, $s = 2k$

## Storage:

- ❧ Dimension reduction maps: $O(m + n)$
- ❧ Sketches: $O(k(m + n))$

## Arithmetic:

- ❧ Linear update: Depends on structure of update (cheap!)
- ❧ SKETCHYSVD: $O(k^2(m + n))$
  - ❧ Computation of range and co-range: $O(k^2(m + n))$
  - ❧ Computation of core: $O(k(m + n) + k^3)$
  - ❧ Truncated SVD of core: $O(k^3)$
  - ❧ Consolidation: $O(k^2(m + n))$

## Communication:

- ❧ One pass over data

**Sources:** Charikar et al. 2002; Cormode & Muthukrishnan 2005; Sarlós 2006; Woolfe et al. 2008; Clarkson & Woodruff 2009, 2011; HMT 2011; Nelson et al. 2012–2015; Meng & Mahoney 2013; Cohen 2015; TYUC 2016–2019....

# Empirical Performance

# Important Things I'm Not Going to Show You

- SKETCHYSVD is insensitive to the choice of dimension reduction map

- Theory gives parameter choices $(k, s)$ that are nearly optimal in practice

- SKETCHYSVD beats earlier techniques for synthetic and real data

- Methodology for estimating errors and selecting the truncation rank $r$

- Sampling distribution of approximation error and error estimator

- Other structured approximations via re-factorization or matrix nearness

- Extension to low-rank Tucker approximation of a tensor

**Sources:** HMT 2011; TYUC 2016–2019; SGLTU 2018–2019.

# Why Truncate?



(A) Tail Energy $\tau_r(\hat{\boldsymbol{A}}_k)$

(B) Error in Rank-$r$ Truncation $\hat{\boldsymbol{A}}_r$

$$\tau_{r+1}(\boldsymbol{M}) = \|\boldsymbol{M} - [\![\boldsymbol{M}]\!]_r\|_{\mathrm{F}}$$

$$\mathrm{relerr}(\hat{\boldsymbol{A}}) = \frac{\|\boldsymbol{A} - \hat{\boldsymbol{A}}\|_{\mathrm{F}}}{\tau_r(\boldsymbol{A})} - 1$$

**Comments:** StreamVel Data: $m = 10,738$; $n = 5,001$; 430 MB. Algorithm: Sparse maps; $s = 2k + 1$.

# Reconstruction of von Kármán Street

# Left Singular Vectors of von Kármán Street

Approximate [TYUC19]                    Exact



**Comments:** Data: $m = 10,738$; $n = 5,001$; 430 MB. Algorithm: Sparse maps; rank $r = 5$; storage $T = 48(m + n)$. Compression: $71\times$.

# Left Singular Vectors of von Kármán Street

Approximate [HMT11]           Exact



**Comments:** Data: $m = 10,738$; $n = 5,001$; 43O MB. Algorithm: Sparse maps; rank $r = 5$; storage $T = 48(m + n)$. Compression: $71\times$.

# Singular Vectors of Sea Surface Temperature Data



Spatiotemporal Avg.

Austral / Boreal

Madden–Julian Osc.

Madden–Julian Osc.

La Niña / El Niño

**Comments:** Data: $m = 691,150$; $n = 13,670$; 75 GB. Algorithm: Sparse maps; $k = 48$; $s = 839$. Compression ratio: 222×.

# Contributions

1. The first practical TSVD algorithms for streaming data

2. Rigorous *a priori* error bounds and parameter recommendations

3. *A posteriori* error estimation and methodology for rank truncation

4. Extensions to low-rank Tucker approximation for tensors

5. Validation on applications in scientific computing and optimization

# To learn more...

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

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

**Papers:**

- Halko, Martinsson, & Tropp, "Finding structure with randomness: Probabilistic algorithms for computing approximate matrix decompositions," *SIREV*, 2011.
- Tropp, Yurtsever, Udell, & Cevher, "Sketchy decisions: Low-rank matrix optimization with optimal storage," AISTATS 2017.
- Tropp, Yurstever, Udell, & Cevher, "Practical sketching algorithms for low-rank matrix approximation," *SIMAX*, 2017.
- Tropp, Yurtsever, Udell, & Cevher, "Fixed-rank approximation of a positive-semidefinite matrix from streaming data," *NeurIPS*, 2017.
- Tropp, Yurtsever, Udell, & Cevher, "Streaming low-rank matrix approximation with an application to scientific simulation," arXiv cs.NA 1902.08651.
- Sun, Guo, Tropp, & Udell, "Tensor random projections for low-memory dimension reduction," *NeurIPS Relational Databases Workshop*, 2018.
- Sun, Guo, Luo, Tropp, & Udell, "Low-rank Tucker approximation of a tensor from streaming data." Coming soon!
- Cevher, Tropp, & Yurtsever, "Scalable semidefinite programming." Coming soon!

# Supplementary Materials

# Insensitivity to Dimension Reduction Map



LowRankHiNoise
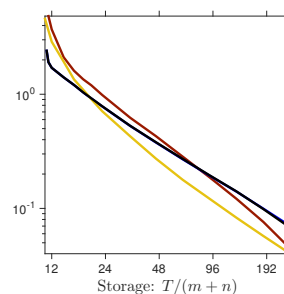LowRankMedNoise
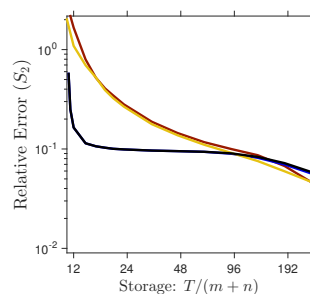LowRankLowNoise

PolyDecaySlow
PolyDecayMed
PolyDecayFast

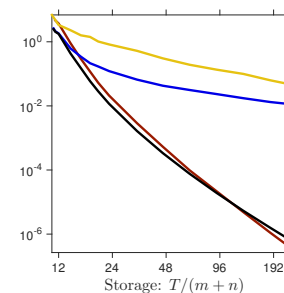ExpDecaySlow
ExpDecayMed
ExpDecayFast

**Comments:** Effective rank $R = 10$, approximation rank $r = 10$, Schatten 2-norm.
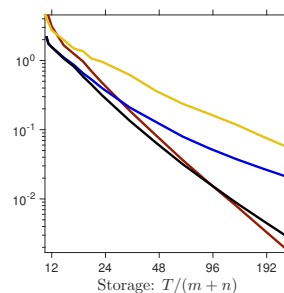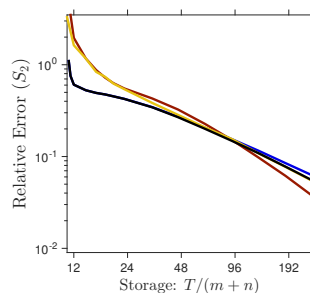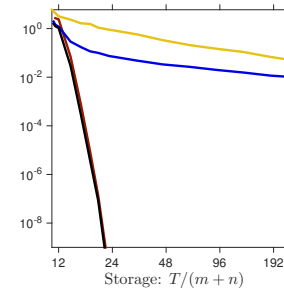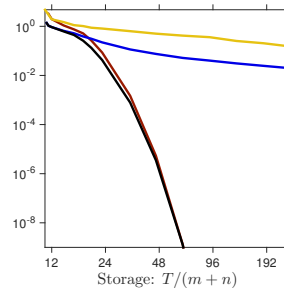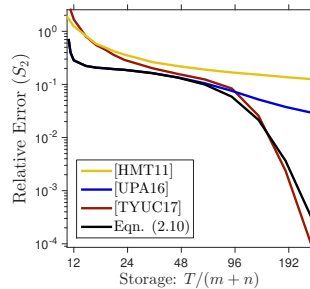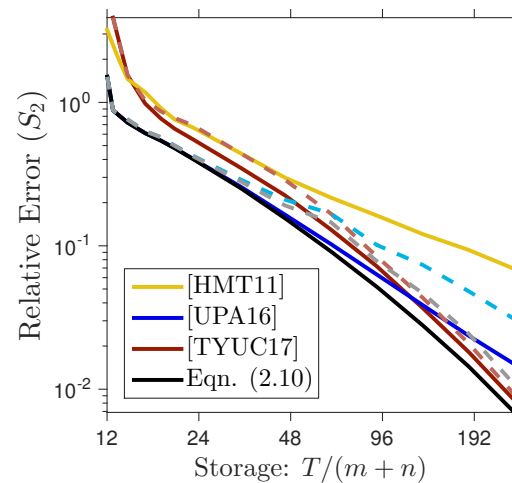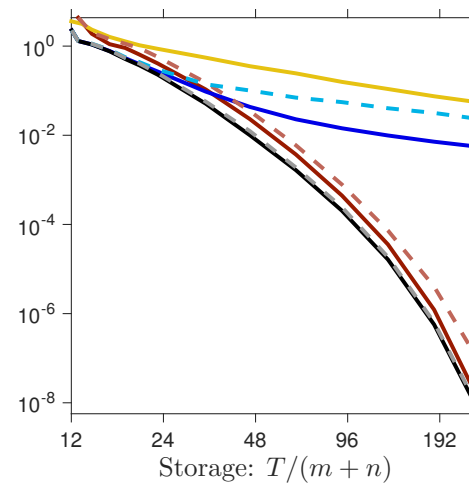
# Performance with Theoretical Parameter Choices



**Comments:** Gaussian maps, effective rank $R = 10$, approximation rank $r = 10$, Schatten 2-norm.

# Method Comparison: Synthetic Data

LowRankHiNoise
LowRankMedNoise
LowRankLowNoise

PolyDecaySlow
PolyDecayMed
PolyDecayFast

ExpDecaySlow
ExpDecayMed
ExpDecayFast



**Comments:** Gaussian maps, effective rank $R = 10$, oracle parameters, approximation rank $r = 10$, Schatten 2-norm.)
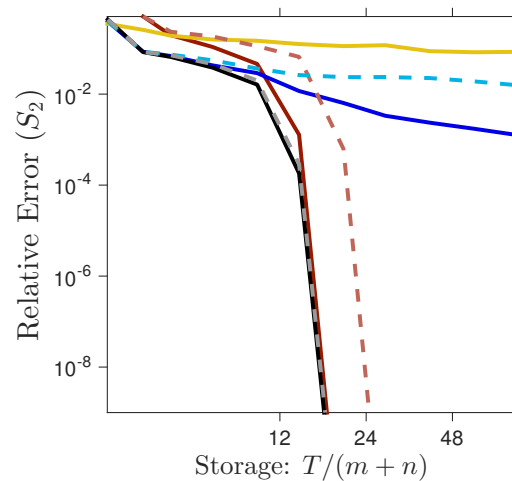
# Method Comparison: Real Data

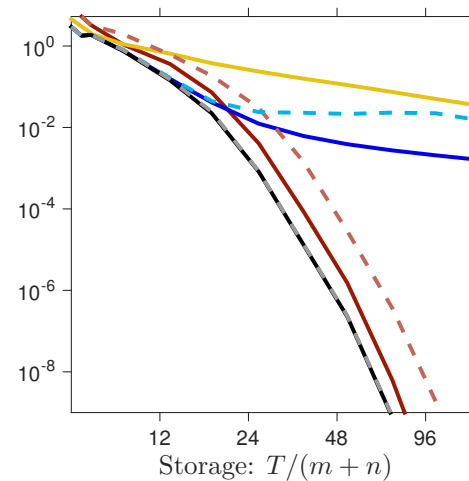

**MinTemp**
$m = 19,264$
$n = 7,305$
$r = 10$

**StreamVel**
$m = 10,738$
$n = 5,001$
$r = 10$

**MaxCut**
$n = 2,000$
$r = 1$

**PhaseRetrieval**
$n = 25,921$
$r = 5$

Legend: [HMT11], [UPA16], [TYUC17], Eqn. (2.10)

Axes: Relative Error ($S_2$); Storage: $T/(m+n)$

**Comments:** Sparse maps, Schatten 2-norm. Solid lines are errors with oracle parameters; dashed lines are *a priori* parameter choices.
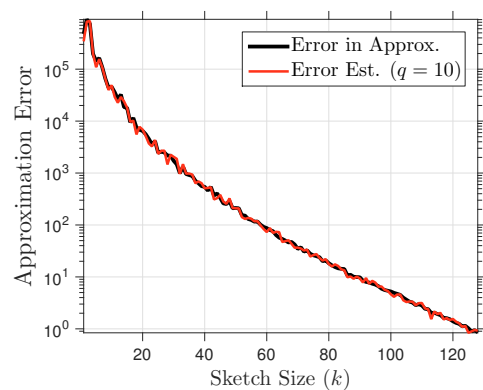
# A *Posteriori* Error Estimation

- Fix a sketch size parameter $q$
- Draw a random Gaussian dimension reduction map $\boldsymbol{\Theta} \in \mathbb{C}^{m \times q}$
- Maintain an error sketch $\boldsymbol{S} = \boldsymbol{\Theta} \boldsymbol{A}$

- Given an approximation $\hat{\boldsymbol{A}}$, compute the error estimator

$$\mathrm{err}_2^2(\hat{\boldsymbol{A}}) = \left\| \boldsymbol{S} - \boldsymbol{\Theta} \hat{\boldsymbol{A}} \right\|_{\mathrm{F}}^2$$

- The error estimator is unbiased and concentrates sharply
- We can also compute an empirical upper bound on the scree curve as

$$\overline{\mathrm{scree}}(r) = \left[ \frac{\tau_{r+1}(\hat{\boldsymbol{A}}) + \mathrm{err}_2(\hat{\boldsymbol{A}})}{\mathrm{err}_2(\boldsymbol{0})} \right]^2 .$$
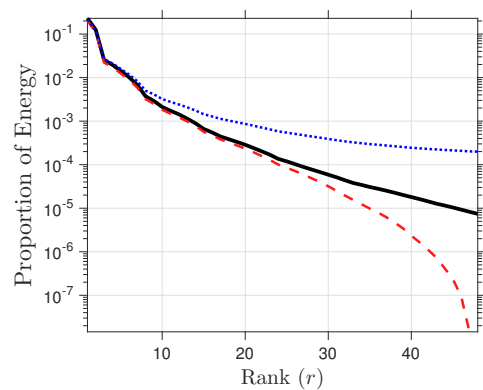
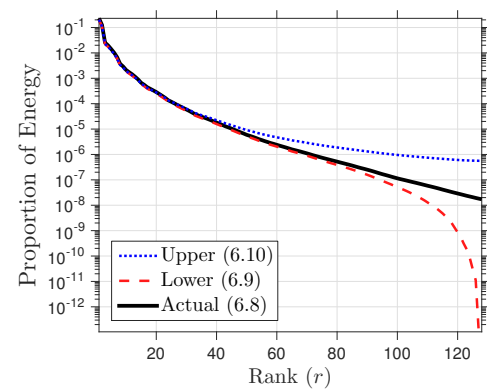# Error Estimates and Empirical Scree Curves



(A) Error Estimates for $\hat{A}$
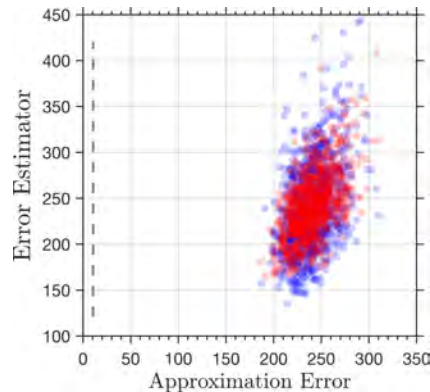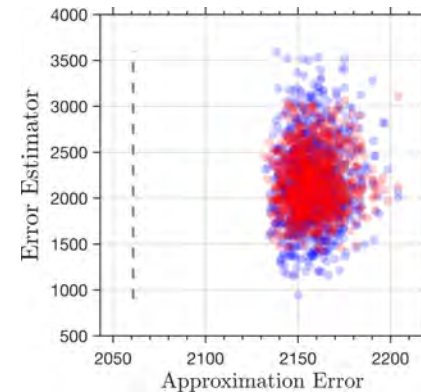
(B) Scree Plot ($k = 16$)

(C) Scree Plot ($k = 48$)

(D) Scree Plot ($k = 128$)

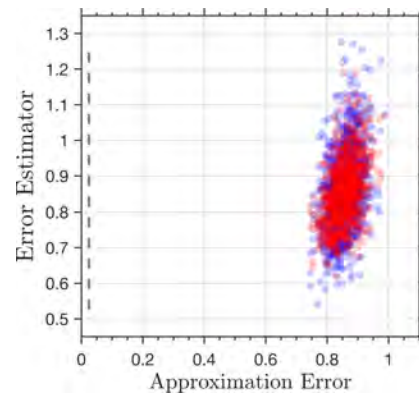**Comments:** `StreamVel`, sparse maps, $s = 2k + 1$, $q = 10$.
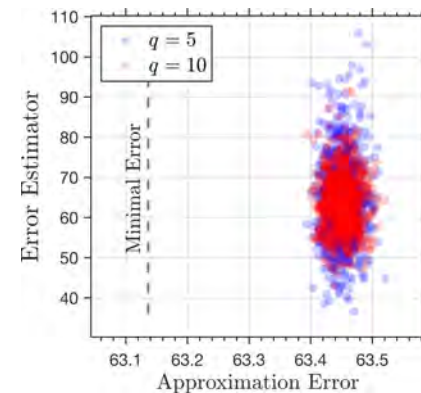
# Sampling Distribution of Error and Estimator



(A) Rank-$k$ Approximation ($k = 48$)

(B) Rank-$r$ Truncation ($k = 48$, $r = 12$)

(C) Rank-$k$ Approximation ($k = 128$)

(D) Rank-$r$ Truncation ($k = 128$, $r = 32$)

Comments: StreamVel, sparse maps, $s = 2k + 1$.