

# A simple randomized algorithm for approximating the spectral norm of streaming data

Spencer Shortt

<sup>1</sup>Advisor: Stephen Becker  
University of Colorado Boulder

<sup>2</sup>Committee Member: Kyle Luh  
University of Colorado Boulder

<sup>3</sup>Committee Member: Sean O'Rourke  
University of Colorado Boulder

April 2023

# Motivation: Why approximate the spectral norm?

Spectral norms can be used as an error estimator when trying to approximate matrices.

## Remark 2.1 (Martinson et. al., 2020)

For example, let us consider a variant of the spiked covariance model that is common in statistics applications. Suppose we need to approximate a rank-one matrix contaminated with additive noise:  $A = \vec{u}\vec{u}^* + G \in \mathbb{R}^{n \times n}$ , where  $\|\vec{u}\| = 1$  and  $G \in \mathbb{R}^{n \times n}$  has independent entries from  $\mathcal{N}(0, n^{-1})$  entries. With respect to the Frobenius norm, the zero matrix is almost as good an approximation of  $A$  as the rank-one matrix  $uu^*$ :

$$\mathbb{E}[\|A - \vec{u}\vec{u}^*\|_F^2] = \varepsilon^2 n \text{ and } \mathbb{E}[\|A - 0\|_F^2] = \varepsilon^2 n$$

# An Existing Approach: Power Method

Liberty et. al., (2007)

Suppose  $A$  is an  $m \times n$  complex-valued matrix and  $\vec{\omega}$  is a  $n \times 1$  column vector with i.i.d. entries from a complex gaussian distribution. With  $\vec{\nu} = \frac{\vec{\omega}}{\|\vec{\omega}\|_2}$ , we define

$$p_j(A) = \sqrt{\frac{\|(A^* A)^j \vec{\nu}\|_2}{\|(A^* A)^{j-1} \vec{\nu}\|_2}}.$$

Then  $p_j(A) \geq \|A\|/10$  with probability greater than  $1 - 4\sqrt{n/(j-1)}100^{-j}$ , and  $p_j(A) \leq \|A\|$  for all  $j$ .

# Streaming Data

Let  $A$  be an  $m \times n$  matrix of data, and suppose we went to append it with an  $m \times k$  dataset  $B$ .

$$C = [A \mid B] = [A \mid 0_{m \times k}] + [0_{m \times n} \mid B] = A' + B'$$

By writing  $C = A' + B'$  as above, we run into a potential storage issue. We must store both the old data  $A$  and new data  $B$  in order to calculate the spectral norm approximation:

$$p_j(C) = \sqrt{\frac{\|((A' + B')^*(A' + B'))^j \vec{v}\|_2}{\|((A' + B')^*(A' + B'))^{j-1} \vec{v}\|_2}}$$

# A Different Approach:

## Lemma 4.1 (Halko et. al., 2011)

Let  $A$  be a real  $m \times n$  matrix. Fix a positive integer  $r$  and a real number  $\alpha > 1$ . Draw an independent family  $\{\vec{\omega}_i : i = 1, 2, \dots, N\}$  of standard Gaussian vectors. Then

$$\|A\| \leq \alpha \max_{i=1,2,\dots,N} \|A\vec{\omega}_i\|$$

except with probability  $\alpha^{-N}$

# Efficient Storage for Streaming Data

Let  $\Omega_A = [\vec{\omega}_1 \quad \vec{\omega}_2 \quad \dots \quad \vec{\omega}_N]$  be an  $n \times N$  matrix whose columns are independent standard Gaussian vectors, and define

$$Y_A = A\Omega_A = [A\vec{\omega}_1 \quad A\vec{\omega}_2 \quad \dots \quad A\vec{\omega}_N].$$

To achieve the bound on the previous slide, calculate  $\max_{i=1,2,\dots,N} \|A\vec{\omega}_i\|$

Suppose now that we append the  $m \times k$  matrix  $B$  to  $A$  to get  $C = [A \mid B]$ . We let  $\Omega_B$  be a  $k \times N$  matrix whose columns are independent standard Gaussian vectors, and define  $\Omega_C = \begin{bmatrix} \Omega_A \\ \Omega_B \end{bmatrix}$ . Then

$$Y_C = C\Omega_C = [A \mid B] \begin{bmatrix} \Omega_A \\ \Omega_B \end{bmatrix} = A\Omega_A + B\Omega_B = Y_A + B\Omega_B,$$

implying that we only need to calculate  $B\Omega_B$  after storing the  $m \times N$  matrix  $Y_A$ .

# Frobenius Norm is off by factor of $r^{1/2}$

We can write the Frobenius norm as the  $\ell_2$ -norm of the singular values:

$\|A\|_F = \sqrt{\sum_{i=1}^r \sigma_i^2}$ . Using this and the fact that the spectral norm of  $A$  is the largest singular value of  $A$ , we have

$$\|A\| \leq \|A\|_F \leq r^{1/2} \|A\|$$

since

$$\sigma_{\max} \leq \left( \sum_{j=1}^r \sigma_j^2 \right)^{1/2} \leq r^{1/2} \sigma_{\max}.$$

This tells us that the Frobenius norm can be off from the spectral norm by a factor of  $r^{1/2}$ .

# Estimate is greater than Frobenius norm

We show  $\mathbb{E}[\|A\tilde{\omega}\|^2] = \|A\|_F^2$ :

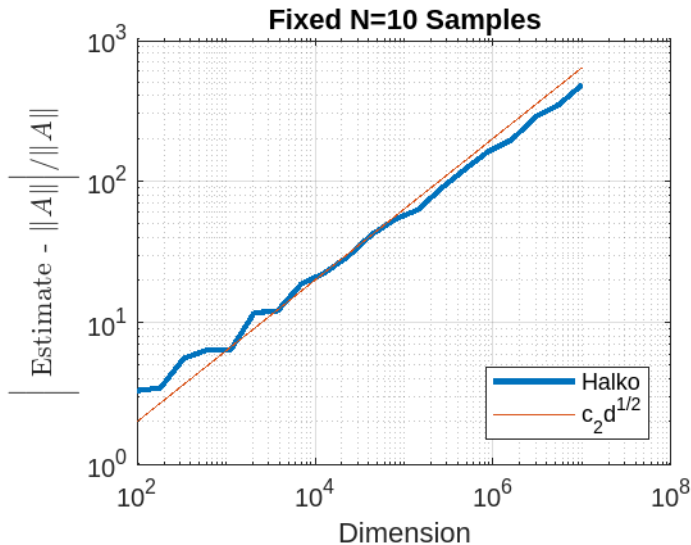
$$\begin{aligned}\mathbb{E}[\|A\tilde{\omega}\|^2] &= \mathbb{E}[\tilde{\omega}^T A^T A \tilde{\omega}] = \text{Tr}(\mathbb{E}[\tilde{\omega}^T A^T A \tilde{\omega}]) = \mathbb{E}[\text{Tr}(\tilde{\omega}^T A^T A \tilde{\omega})] \\ &= \mathbb{E}[\text{Tr}(A^T A \tilde{\omega} \tilde{\omega}^T)] = \text{Tr}(\mathbb{E}[A^T A \tilde{\omega} \tilde{\omega}^T]) = \text{Tr}(A^T A \mathbb{E}[\tilde{\omega} \tilde{\omega}^T]) \\ &= \text{Tr}(A^T A) = \|A\|_F^2\end{aligned}$$

Analyzing the bound given by Halko et. al. (2011), we see

$$\mathbb{E}[\max_{i=1,2,\dots,N} \|A\tilde{\omega}_i\|^2] \geq \mathbb{E}\left[\frac{1}{N} \sum_{i=1}^N \|A\tilde{\omega}_i\|^2\right] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\|A\tilde{\omega}_i\|^2] = \|A\|_F^2$$



# Plot of Error



# The $\ell_4$ norm is better

By the same argument as before,

$$\|A\| \leq \|\vec{\sigma}\|_4 \leq r^{1/4} \|A\|.$$

One idea is to approximate the  $\ell_4$ -norm of the singular values since this is a tighter bound.

Let  $\vec{\omega}, \vec{\nu} \in \mathcal{N}(0, I_r)$  be independent gaussian random vectors. Define the random variable  $X = (A\vec{\omega})^T A\vec{\nu}$ . We will show  $\mathbb{E}[X^2] = \|\vec{\sigma}\|_4^4$

# WLOG, use diagonal matrices

Let  $A = U\Sigma V^T$  be the singular value composition of our  $m \times n$  matrix  $A$ .  
By orthogonality,

$$\begin{aligned}\mathbb{E}[X^2] &= \mathbb{E}[(A\vec{\omega})^T A\vec{\nu}]^2 = \mathbb{E}[(\vec{\omega}^T A^T A\vec{\nu})^2] = \mathbb{E}[(\vec{\omega}^T V\Sigma U^T U\Sigma V^T \vec{\nu})^2] \\ &= \mathbb{E}[(\vec{\omega}^T V\Sigma^2 V^T \vec{\nu})^2] = \mathbb{E}[(V^T \vec{\omega})^T \Sigma^2 V^T \vec{\nu})^2]\end{aligned}$$

Since  $\vec{\omega}, \vec{\nu} \in \mathcal{N}(0, I_n)$ , we have that  $V^T \vec{\omega}, V^T \vec{\nu} \in \mathcal{N}(0, V^T V) = \mathcal{N}(0, I_n)$ .  
Thus,

$$\mathbb{E}[(A\vec{\omega})^T A\vec{\nu}]^2 = \mathbb{E}[(\Sigma\omega)^T \Sigma\nu]^2$$

Furthermore, since  $\Sigma$  only has  $r$  non-zero values along its diagonal, without loss of generality, we can let  $\Sigma$  be an  $r \times r$  diagonal matrix from here on and have  $\vec{\omega}, \vec{\nu} \in \mathcal{N}(0, I_r)$ , and later on we will assume  $A$  to be the same.

## Calculating the $\ell_4$ norm:

$$\begin{aligned}\mathbb{E}[X^2] &= \mathbb{E}[(\Sigma \vec{\omega})^T \Sigma \vec{\nu}]^2 = \mathbb{E}[(\sum_{j=1}^r \sigma_j^2 \omega_j \nu_j)^2] \\ &= \mathbb{E}[\sum_{j=1}^r \sum_{k=1}^r \sigma_j^2 \sigma_k^2 \omega_j \omega_k \nu_j \nu_k] = \sum_{j=1}^r \sum_{k=1}^r \sigma_j^2 \sigma_k^2 \mathbb{E}[\omega_j \omega_k \nu_j \nu_k] \\ &= \sum_{j=1}^r \sum_{k=1}^r \sigma_j^2 \sigma_k^2 \mathbb{E}[\omega_j \omega_k] \mathbb{E}[\nu_j \nu_k] = \sum_{j=1}^r \sigma_j^4 = \|\vec{\sigma}\|_4^4.\end{aligned}$$

Practically speaking, we draw random vectors from a Gaussian distribution to create a sample mean to approximate  $\mathbb{E}[X^2]$ . Thus, we would like to show that the difference  $\left| \frac{1}{N} \sum_{j=1}^N X_j^2 - \mathbb{E}[X^2] \right|$  is small with high probability.

# Sub-Weibull Random Variables

We define  $X$  to be sub-Weibull random variable with tail parameter  $\theta$  if

$$\mathbb{P}(|X| \geq x) \leq a \exp(-bx^{1/\theta}) \text{ for all } x > 0, \text{ for some } \theta, a, b > 0$$

Equivalently, a random variable is a sub-Weibull with tail parameter  $\theta$  if there exists some constant  $K_2 > 0$  such that

$$\|X\|_p := (\mathbb{E}[|X|^p])^{1/p} \leq K_2 p^\theta$$

for all  $p \geq 1$ .

## Examples

Sub-Gaussian random variables have  $\theta = 1/2$

Sub-Exponential have  $\theta = 1$

# $X = (D\vec{\omega})^T D\vec{\nu}$ is sub-exponential

Let  $A$  be a diagonal  $r \times r$  matrix with positive diagonal entries  $\sigma_i$ , and let  $\omega_i, \nu_i \in \mathcal{N}(0, 1)$ . Since  $\omega_i, \nu_i$  are sub-Gaussian, there exists a constant  $k$  such that for all  $p \geq 1$ ,

$$\|\omega_i\|_p \leq kp^{1/2}.$$

Since  $\|\cdot\|_p$  is a norm, we can use the triangle inequality on  $X$ :

$$\|X\|_p = \left\| \sum_{i=1}^r \sigma_i^2 \omega_i \nu_i \right\|_p \leq \sum_{i=1}^r \sigma_i^2 \|\omega_i \nu_i\|_p = \sum_{i=1}^r \sigma_i^2 (\mathbb{E}[|\omega_i|^p |\nu_i|^p])^{1/p}.$$

By independence, the above equals

$$\sum_{i=1}^r \sigma_i^2 (\mathbb{E}[|\omega_i|^p])^{1/p} (\mathbb{E}[|\nu_i|^p])^{1/p} \leq \sum_{i=1}^r \sigma_i^2 (kp^{1/2})(kp^{1/2}) = k^2 p \|A\|_F^2$$

# We care about $X^2$ , but there's a problem

$X^2 = ((D\vec{\omega})^T D\vec{\nu})^2$  is sub-Weibull with parameter  $\theta = 2$ :

$$\begin{aligned}\|X^2\|_p &= (\mathbb{E}[|X^2|^p])^{1/p} = ((\mathbb{E}[|X|^{2p}])^{1/2p})^2 = (\|X\|_{2p})^2 \\ &\leq (\|A\|_F^2 k^2(2p))^2 = 4k^4 \|A\|_F^4 p^2.\end{aligned}$$

We would like to use concentration properties of sub-Weibull random variables to show the difference  $\left| \frac{1}{N} \sum_{j=1}^N X_j^2 - \mathbb{E}[X^2] \right|$  is small with high probability.

## Corollary 3.1 (Vladimirova et. al., 2020)

Let  $X_1, \dots, X_n$  be identically distributed sub-Weibull random variables with tail parameter  $\theta$ . Then, for all  $x \geq NK_\theta$ , we have

$$\mathbb{P}(|\sum_{i=1}^N X_i| \geq x) \leq \exp(-(\frac{x}{NK_\theta}))$$

for some constant  $K_\theta$  dependent on  $\theta$ .

The problem is that for our situation,  $K_\theta$  is proportional to  $1/N$ .



## Theorem 3.1 (Kuchibhotla et. al., 2022)

If  $X_1, \dots, X_n$  are independent mean zero random variables with  $\|X_i\|_{\psi_\alpha} < \infty$  for all  $1 \leq i \leq n$  and some  $\alpha > 0$ , then for any vector  $(a_1, \dots, a_n) \in \mathbb{R}^n$ , then we have

$$\mathbb{P}\left(\left|\sum_{i=1}^n a_i X_i\right| \geq 2eC(\alpha)\|b\|_2\sqrt{t} + 2eL_n^*(\alpha)t^{1/\alpha}\|b\|_{\beta(\alpha)}\right) \leq 2e^{-t}$$

for all  $t \geq 0$ , where  $b = (a_1\|X_1\|_{\psi_\alpha}, \dots, a_n\|X_n\|_{\psi_\alpha}) \in \mathbb{R}^n$ .

# Another attempt

## Theorem:

Let  $A$  be an  $m \times n$  real-valued matrix with rank  $r > 16$ . Draw  $\vec{\omega}_i$  and  $\vec{v}_i$  independently from  $\mathcal{N}(0, I_n)$  for all  $i \in \{1, \dots, N\}$ . If we define  $X_i = (A\vec{\omega}_i)^T A\vec{v}_i$ , then there exists a constant  $K > 0$  such that for any  $t > 0$ ,

$$\left| \frac{1}{N} \sum_{i=1}^N |X_i|^{1/2} - \|A\| \right| \leq (r^{1/4} - 1) \|A\| + t,$$

with probability greater than  $1 - 2 \exp(-\frac{Nt^2}{Kr\|A\|^2})$ .

This theorem is far from ideal.

If  $\|A\| \leq \frac{1}{N} \sum_{i=1}^N |X_i|^{1/2}$ , we have that  $\frac{1}{N} \sum_{i=1}^N |X_i|^{1/2} \leq r^{1/4} \|A\| + t$  and is actually a slightly better approximation than our estimator  $\frac{1}{N} \sum_{i=1}^N X_i^2$ .

However, it is not guaranteed that  $\|A\| \leq \frac{1}{N} \sum_{i=1}^N |X_i|^{1/2}$ .

# (Proof) Concave Jensen

We use the concave version of Jensen's inequality:

$$\mathbb{E}[|X|^{1/2}] = \mathbb{E}[|X|^{2/4}] \leq (\mathbb{E}[X^2])^{1/4} = \|\vec{\sigma}\|_4$$

If  $\|A\| \leq \mathbb{E}[|X|^{1/2}]$ ,

$$\mathbb{E}[|X|^{1/2}] - \|A\| \leq r^{1/4}\|A\| - \|A\| = (r^{1/4} - 1)\|A\|,$$

and if  $\|A\| \geq \mathbb{E}[|X|^{1/2}]$ ,

$$\|A\| - \mathbb{E}[|X|^{1/2}] \leq \|A\| \leq (r^{1/4} - 1)\|A\|$$

Thus we have a bound on the absolute value of the error.

## (Proof) $X^{1/2}$ is sub-Gaussian

The advantage of using  $|X|^{1/2}$  is that it is sub-Gaussian with constant proportional to  $\|A\|_F$ . Using Jensen's inequality again, we see

$$\| |X|^{1/2} \|_p = (\mathbb{E}[|X|^{p/2}])^{1/p} \leq ((\mathbb{E}[|X|^p])^{1/p})^{1/2} = (\|X\|_p)^{1/2} \leq k \|A\|_F p^{1/2}$$

Thus, we will apply general Hoeffding's inequality to show  $\mathbb{E}[|X|^{1/2}]$  can be closely approximated by  $\frac{1}{N} \sum_{j=1}^N |X_j|^{1/2}$  with high probability.

# (Proof) General Hoeffding's Inequality

Given a random variable  $X$ , we define the sub-Gaussian norm of  $X$  to be

$$\|X\|_{\psi_2} = \inf\{t > 0 : \mathbb{E}[\exp(X^2/t^2)] \leq 2\}$$

## General Hoeffding's Inequality (Vershynin, 2018)

Let  $X_1, X_2, \dots, X_N$  be independent, mean zero, sub-gaussian random variables, and  $a = (a_1, a_2, \dots, a_N) \in \mathbb{R}^N$ . Then for every  $t \geq 0$

$$\mathbb{P}\left(\left|\sum_{j=1}^N a_j X_j\right| \geq t\right) \leq 2 \exp\left(-\frac{ct^2}{K^2 \|a\|_2^2}\right)$$

where  $K = \max_j \|X_j\|_{\psi_2}$

# (Proof) Applying Hoeffding

Using the triangle inequality,

$$\begin{aligned}\| |X|^{1/2} - \mathbb{E}[|X|^{1/2}] \|_p &\leq \| |X|^{1/2} \|_p + \| \mathbb{E}[|X|^{1/2}] \|_p \leq k \|A\|_F p^{1/2} + \mathbb{E}[|X|^{1/2}] \\ &\leq k \|A\|_F p^{1/2} + r^{1/4} \|A\| p^{1/2} \leq r^{1/2} (k+1) \|A\| p^{1/2}.\end{aligned}$$

We can assert that  $\| |X|^{1/2} - \mathbb{E}[|X|^{1/2}] \|_{\psi_2} = C r^{1/2} (k+1) \|A\|$  for some constant  $C > 0$ .

# Applying Hoeffding

This lets us apply Hoeffding to the subgaussian random variables  $\tilde{X}_j = |X_j|^{1/2} - \mathbb{E}[|X|^{1/2}]$  with  $a_j = 1/N$  for all  $j$  and  $K = C^2(k+1)^2/c$ :

$$\mathbb{P}\left(\left|\frac{1}{N} \sum_{j=1}^N |X_j|^{1/2} - \mathbb{E}[|X|^{1/2}]\right| \geq t\right) \leq 2 \exp\left(-\frac{Nt^2}{Kr\|A\|^2}\right)$$

# (Proof) Conclusion

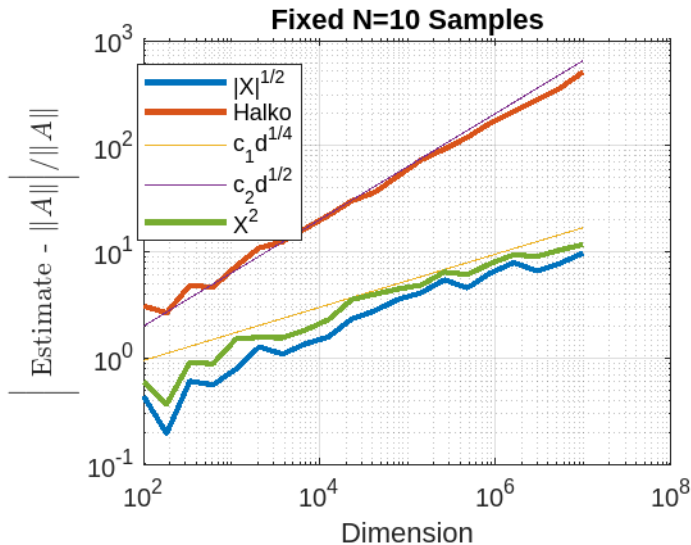
Finally, by the triangle inequality,

$$\begin{aligned} \left| \frac{1}{N} \sum_{i=1}^N |X_i|^{1/2} - \|A\| \right| &\leq \left| \mathbb{E}[|X|^{1/2}] - \|A\| \right| + \left| \frac{1}{N} \sum_{i=1}^N |X_i|^{1/2} - \mathbb{E}[|X|^{1/2}] \right| \\ &\leq (r^{1/4} - 1) \|A\| + t \end{aligned}$$

with probability greater than  $1 - 2 \exp(-\frac{Nt^2}{Kr\|A\|^2})$ .



# Conclusion





N. Halko, P. G. Martinsson, J. Tropp, *Finding Structure with Randomness: Probabilistic Algorithms for Constructing Approximate Matrix Decompositions*, SIAM Review **53**(2), (2011), 217-288.



E. Liberty, F. Woolfe, P.G. Martinsson, V. Rokhlin, M. Tygert, *Randomized algorithms for the low-rank approximation of matrices*, Proc. Natl. Acad. Sci. **104**(51) (2007), 20167–20172.



A. Kuchibhotla, A. Chakraborty, *Moving beyond sub-Gaussianity in high-dimensional statistics: applications in covariance estimation and linear regression*, Inf. Inference **11**(4) (2022), 1389-1456.



P. Martinsson, J. Tropp, Randomized numerical linear algebra: Foundations and algorithms. *Acta Numerica*, **29** (2020), 403–572.



R. Vershynin, *High dimensional probability. An introduction with applications in Data Science*. Cambridge University Press, 2018.



M. Vladimirova, S. Girard, H. Nguyen, J. Arbel *Sub-Weibull distributions: Generalizing sub-Gaussian and sub-Exponential properties to heavier tailed distributions*. *Stat.* **9** (2020)

# Highlighting text

In this slide, some important text will be **highlighted** because it's important. Please, don't abuse it.

## Remark

Sample text

## Important theorem

Sample text in red box

## Examples

Sample text in green box. The title of the block is “Examples”.