# Randomized Numerical Linear Algebra

http://bicmr.pku.edu.cn/~wenzw/bigdata2023.html

Acknowledgement: this slides is based on Prof. Petros Drineas and Prof. Michael W. Mahoney's, Prof. Gunnar Martinsson lecture notes

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● のへで

1/49

## Outline



### Randomized Numerical Linear Algebra (RandNLA)

- Approximating Matrix Multiplication
- 3 Approximate SVD
- 4 Random Sampling for SVD
- 5 Single View Algorithm For Matrix Approximation

**Randomization and sampling** allow us to design provably accurate algorithms for problems that are:

#### Massive

(matrices so large that can not be stored at all, or can only be stored in slow memory devices)

#### Computationally expensive or NP-hard

(combinatorial optimization problems such as the Column Subset Selection Problem)

# RandNLA: sampling rows/columns

#### **Randomized algorithms**

• By (carefully) sampling rows/columns of a matrix, we can construct new, smaller matrices that are close to the original matrix (w.r.t. matrix norms) with high probability.

$$\left(\begin{array}{c} & A \\ & A \end{array}\right)\left(\begin{array}{c} & B \\ & B \end{array}\right)\approx\left(\begin{array}{c} & C \\ & C \end{array}\right)\left(\begin{array}{c} & R \end{array}\right)$$

 By preprocessing the matrix using random projections, we can sample rows/columns much less carefully (uniformly at random) and still get nice bounds with high probability.

# RandNLA: sampling rows/columns

#### Matrix perturbation theory

• The resulting smaller matrices behave similarly (in terms of singular values and singular vectors) to the original matrices thanks to the norm bounds.

Structural results that "decouple" the "randomized" part from the "matrix perturbation" part are important in the analyses of such algorithms.

#### Interplay

- Applications in BIG DATA: (Data Mining, Information Retrieval, Machine Learning, Bioinformatics, etc.)
- Numerical Linear Algebra: Matrix computations and linear algebra (ie., perturbation theory)
- Theoretical Computer Science: Randomized and approximation algorithms

### Issues

- Selecting good columns that "capture the structure" of the top principal components
  - Combinatorial optimization problem; hard even for small matrices.
  - Often called the Column Subset Selection Problem (CSSP).
  - Not clear that such columns even exist.

The two issues:

- Fast approximation to the top k singular vectors of a matrix, and
- Selecting columns that capture the structure of the top k singular vectors

are connected and can be tackled using the same framework

### Outline



Randomized Numerical Linear Algebra (RandNLA)



### Approximating Matrix Multiplication

3 Approximate SVD

4 Random Sampling for SVD

5 Single View Algorithm For Matrix Approximation

# Approximating Matrix Multiplication

#### **Problem Statement**

Given an m-by-n matrix A and an n-by-p matrix B, approximate the product AB, Or equivalently, Approximate the sum of n rank-one matrices

$$AB = \sum_{i=1}^{n} \underbrace{\left(A^{(i)}\right)\left(B_{(i)}\right)}_{\in \mathbb{R}^{m \times p}}$$

- $A^{(i)}$  the *i*-th column of A
- $B_{(i)}$  the *i*-th row of B
- Each term in the summation is a rank-one matrix

# A sampling approach

$$AB = \sum_{i=1}^{n} \underbrace{\left(A^{(i)}\right)\left(B_{(i)}\right)}_{\in \mathbb{R}^{m \times p}}$$

Algorithm

• Fix a set of probabilities  $p_i$ , i = 1, ..., n, summing up to 1.

 Approximate the product AB by summing the c terms, after scaling.

### Generate Discrete Distributions

Consider a discrete random variable with possible values  $c_1 < \ldots < c_n$ . The probability attached to  $c_i$  is  $p_i$ . Let

$$q_0 = 0, \quad q_i = \sum_{j=1}^i p_j.$$

They are the cumulative probabilities associated with  $c_i$ , i.e.,  $q_i = F(c_i)$ .

#### To sample this distribution

- generate a uniform U
- find  $K \in \{1, ..., n\}$  such that  $q_{K-1} < U < q_K$

• set 
$$X = c_K$$

#### • Sampling with replacement:

Each data unit in the population is allowed to appear in the sample more than once.

It is easy to analyze mathematically.

#### • Sampling without replacement: Each data unit in the population is allowed to appear in the sample no more than once.

# A sampling approach

$$AB = \sum_{i=1}^{n} \underbrace{\left(A^{(i)}\right)\left(B_{(i)}\right)}_{\in \mathbb{R}^{m \times p}}$$
$$\approx \frac{1}{c} \sum_{t=1}^{c} \frac{1}{p_{j_{t}}} \underbrace{\left(A^{(j_{t})}\right)\left(B_{(j_{t})}\right)}_{\in \mathbb{R}^{m \times p}}$$

Keeping the terms  $j_1, j_2, \ldots, j_c$ 

# The algorithm (matrix notation)



Algorithm:

- Pick c columns of *A* to form an m-by-c matrix *C* and the corresponding c rows of *B* to form a c-by-p matrix *R*.
- Approximate *AB* by *CR*.

Note

- We pick the columns and rows with non-uniform probabilities.
- We scale the columns (rows) prior to including them in C(R).

# The algorithm (matrix notation)



Algorithm:

- Create C and R by performing c i.i.d. trials, with replacement.
- For t = 1, ..., c, pick a column  $A^{(j_t)}$  and a row  $B_{(j_t)}$  with probability

$$\mathbf{P}(j_t = i) = \frac{\|A^{(i)}\|_2 \|B_{(i)}\|_2}{\sum_{j=1}^n \|A^{(j)}\|_2 \|B_{(j)}\|_2}$$

• Include  $A^{(j_t)}/(cp_{j_t})^{1/2}$  as a column of *C*, and  $B_{(j_t)}/(cp_{j_t})^{1/2}$  as a row of *R* 

# The algorithm (matrix notation)

 Let S be an n-by-c matrix whose t-th column (for t = 1,...,c) has a single non-zero entry, namely

$$S_{j_t t} = \frac{1}{\sqrt{cp_{j_t}}}$$

Clearly:

$$AB \approx CR = (AS)(S^TB)$$

Note: *S* is sparse (has exactly c non-zero elements, one per column).

• It is easy to implement this particular sampling in two passes.

# A bound for the Frobenius norm

For the above algorithm,

$$\mathbf{E}[\|AB - CR\|_{F}] = \mathbf{E}[\|AB - ASS^{T}B\|_{F}] \le \frac{1}{c}\|A\|_{F}\|B\|_{F}$$

- The expectation of CR (element-wise) is AB (unbiased estimator), regardless of the sampling probabilities.
- Our particular choice of sampling probabilities minimizes the variance of the estimator (w.r.t. the Frobenius norm of the error AB-CR).
- prove using elementary manipulations of expectation
- Measure concentration follows from a martingale argument.
- The above bound also implies an upper bound for the spectral norm of the error AB CR.

Let  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{p \times p}$ ,  $1 \le c \le n$ , and  $p_i \ge 0$ ,  $\sum_i p_i = 1$ . Then

$$\mathbf{E}[(CR)_{ij}] = (AB)_{ij}, \quad \text{Var}[(CR)_{ij}] = \frac{1}{c} \sum_{i=1}^{n} \frac{A_{ik}^2 B_{kj}^2}{p_k} - \frac{1}{c} (AB)_{ij}^2$$

• Define 
$$X_t = \left(\frac{A^{(i_t)}B_{(i_t)}}{cp_{i_t}}\right)_{ij} = \frac{A_{ii_t}B_{i_tj}}{cp_{i_t}}$$
. Then  

$$\mathbf{E}[X_t] = \sum_{k=1}^n p_k \frac{A_{ik}B_{kj}}{cp_k} = \frac{1}{c}(AB)_{ij} \text{ and } \mathbf{E}[X_t^2] = \sum_{k=1}^n \frac{A_{ik}^2B_{kj}^2}{c^2p_k}$$

•  $\mathbf{E}[(CR)_{ij}] = \sum_{t=1}^{c} \mathbf{E}[X_t] = (AB)_{ij}$ 

$$\operatorname{Var}[X_t] = E[X_t^2] - E[X_t]^2 = \sum_{k=1}^n \frac{A_{ik}^2 B_{kj}^2}{c^2 p_k} - \frac{1}{c^2} (AB)_{ij}^2$$

Lemma:

$$\mathbf{E}[\|AB - CR\|_F^2] = \sum_{k=1}^n \frac{\|A^{(k)}\|^2 \|B_{(k)}\|^2}{cp_k} - \frac{1}{c} \|AB\|_F^2$$

Proof:

$$\mathbf{E}[\|AB - CR\|_{F}^{2}] = \sum_{i=1}^{n} \sum_{j=1}^{p} \mathbf{E}[(AB - CR)_{ij}^{2}] = \sum_{i=1}^{n} \sum_{j=1}^{p} \operatorname{Var}[(CR)_{ij}]$$
$$= \frac{1}{c} \sum_{k=1}^{n} \frac{1}{p_{k}} \left(\sum_{i} A_{ik}^{2}\right) \left(\sum_{i} B_{kj}^{2}\right) - \frac{1}{c} \|AB\|_{F}^{2}$$
$$= \frac{1}{c} \sum_{k=1}^{n} \frac{1}{p_{k}} \|A^{(k)}\|^{2} \|B_{(k)}\|^{2} - \frac{1}{c} \|AB\|_{F}^{2}$$

• Find  $p_k$  to minimize  $\mathbf{E}[||AB - CR||_F^2]$ :

$$\min_{\sum_{k=1}^{n} p_{k}=1} f(p_{1}, \dots, p_{n}) = \sum_{k=1}^{n} \frac{1}{p_{k}} \|A^{(k)}\|^{2} \|B_{(k)}\|^{2}$$

• Introduce 
$$L = f(p_1, \dots, p_n) + \lambda(\sum_{k=1}^n p_k - 1)$$
 and solve  $\frac{\partial L}{\partial p_i} = 0$ 

• It gives 
$$p_k = \frac{\|A^{(k)}\| \|B_{(k)}\|}{\sum_{k'=1}^n |A^{(k')}| |B_{(k')}|}$$
. Then

$$\mathbf{E}[\|AB - CR\|_{F}^{2}] = \frac{1}{c} \left( \sum_{k=1}^{n} \|A^{(k)}\| \|B_{(k)}\| \right)^{2} - \frac{1}{c} \|AB\|_{F}^{2} \\ \leq \frac{1}{c} \|A\|_{F}^{2} \|B\|_{F}^{2}$$

#### If $B = A^T$ , then the sampling probabilities are

$$\mathbf{P}(j_t = i) = \frac{\|A^{(i)}\|_2^2}{\|A\|_F^2}$$

Also,  $R = C^T$ , and the error bounds are:

$$\mathbf{E}[\|AA^{T} - CC^{T}\|_{F}] = \mathbf{E}[\|AA^{T} - ASS^{T}A^{T}\|_{F}] \le \frac{1}{c} \|A\|_{F}^{2}$$

(ロ)、(型)、(目)、(目)、(目)、(ロ)、(20/49)

.

### Outline



Approximating Matrix Multiplication

### 3 Approximate SVD

- 4 Random Sampling for SVD
- 5 Single View Algorithm For Matrix Approximation

# Low-Rank Matrix Approximation

#### Problem Statement:

Given: mxn matrix *A*, and  $0 < k < \min(m, n) = n$ . Goal: Compute a rank-k approximation to *A*.

- Fast low-rank matrix approximation is key to efficiency of superfast direct solvers for integral equations and many large sparse linear systems.
- Indispensable tool in mining large data sets.
- Randomized algorithms compute accurate truncated SVD.
- Minimum work and communication/Exceptionally high success rate.

## Low-rank Approximation

seek to compute a rank-k approximation with  $k \ll n$ 



- Eigenvectors corresponding to leading eigenvalues.
- Singular Value Decomposition (SVD) / Principal Component Analysis (PCA).
- Spanning columns or rows.

The problem being addressed is ubiquitous in applications.

# Review of existing methods: dense matrix

For a dense  $n \times n$  matrix that fits in RAM, excellent algorithms are already part of LAPACK (and incorporated into Matlab, Mathematica, etc).

- Double precision accuracy.
- Very stable.
- $O(n^3)$  asymptotic complexity. Reasonably small constants.
- Require extensive random access to the matrix.
- When the target rank k is much smaller than n, there also exist  $O(n^2k)$  methods with similar characteristics (the well-known Golub-Businger method, RRQR by Gu and Eisentstat, etc).
- For small matrices, the state-of-the-art is quite satisfactory. (By "small" we mean something like n ≤ 10000 on today's computers.)

# Review of existing methods: structured matrix

If the matrix is large, but can rapidly be applied to a vector (if it is sparse, or sparse in Fourier space, or amenable to the FMM, etc.), so called Krylov subspace methods often yield excellent accuracy and speed.

#### Lanczos-based methods:

• From  $v \in \mathbb{R}^n$ , computes orthonormal basis *V* for

$$\mathcal{K}(A, \mathbf{v}) = span\left\{\mathbf{v}, A\mathbf{v}, A^2\mathbf{v}, \cdots, A^{k-1}\mathbf{v}\right\}$$

- **2** Rayleigh-Ritz:  $eig(V^TAV) \Rightarrow$  Ritz pairs  $\approx$  eigenpairs
- If "not converged", update v and go to Step 1.

#### Strength and weakness:

- Most efficient in terms of the number of Av (or SpMv)
- Fast and reliable for computing "not too many" eigenpairs
- Lower concurrency and unable to be warm-started

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

# "New" challenges in algorithmic design

The existing state-of-the-art methods of numerical linear algebra that we have very briefly outlined were designed for an environment where the matrix fits in RAM and the key to performance was to minimize the number of floating point operations required. Currently, communication is becoming the real bottleneck:

- While clock speed is hardly improving at all anymore, the cost of a flop keeps going down rapidly. (Multi-core processors, GPUs, cloud computing, etc.)
- The cost of slow storage (hard drives, flash memory, etc.) is also going down rapidly.
- Communication costs are decreasing, but not rapidly. Moving data from a hard-drive. Moving data between nodes of a parallel machine. (Or cloud computer ...) The amount of fast cache memory close to a processor is not improving much. (In fact, it could be said to be shrinking — GPUs, multi-core, etc.)
- "Deluge of data". Driven by ever cheaper storage and acquisition techniques. Web search, data mining in archives of documents 26/49

#### 两阶段计算框架

- 阶段1 —— 寻找近似空间:构造矩阵H<sub>k</sub> ∈ ℝ<sup>m×k</sup>,其列是正交的, 并使得A ≈ H<sub>k</sub>H<sup>T</sup><sub>k</sub>A.
- 阶段2 构造特定分解:
   构造一个k×n的矩阵B = H<sub>k</sub><sup>T</sup>A.
   2 对小矩阵B进行SVD分解,使得B = Û<sub>k</sub>Σ<sub>k</sub>V<sub>k</sub><sup>T</sup>.
   3 计算U<sub>k</sub> = H<sub>k</sub>Û<sub>k</sub>.

具体地说:

$$H_k H_k^{\mathrm{T}} A = H_k \underbrace{\mathcal{B}}_{=\hat{U}_k \Sigma_k V_k^{\mathrm{T}}} = H_k \hat{U}_k \Sigma_k V_k^{\mathrm{T}} = \underbrace{\mathcal{U}}_{=H_k \hat{U}_k} \Sigma_k V_k^{\mathrm{T}}.$$

如果在第一阶段得到的矩阵 $H_k$ 满足条件 $||A - H_k H_k^T A|| \le \varepsilon$ ,则最终得到的秩k逼近满足:

$$\|A - U_k \Sigma_k V_k^{\mathrm{T}}\| = \|A - H_k H_k^{\mathrm{T}} A\| \leq \varepsilon.$$

# Linear Time SVD Algorithm

- Input: m-by-n matrix A,  $1 \le k \le c \le n$ ,  $\{p_i\}_{i=1}^n$  such that  $p_i \ge 0$ and  $\sum_i p_i = 1$
- Sampling:
  For t =

For 
$$t = 1$$
 to  $c$   
pick  $i_t \in \{1, \dots, n\}$  with  $\mathcal{P}(i_t = \alpha) = p_{\alpha}$ .  
Set  $C^{(t)} = \frac{A^{(i_t)}}{\sqrt{cp_{i_t}}}$ 

- Compute  $C^T C$  and its eigenvalue decomposition, say  $C^T C = \sum_{t=1}^{c} \sigma_t(C)^2 y_t y_t^T$
- Compute  $h_t = \frac{Cy_t}{\sigma_t(C)}$  for t = 1, ..., k. (Note:  $A = U\Sigma V^T$  and  $C = H\Sigma_C Y^T \Longrightarrow H = CY\Sigma_C^{-1}$ )

• Return  $H_k$  where  $H_k^{(t)} = h_t$  and  $\sigma_t(C)$  for  $t = 1, \ldots, k$ 

The left singular vectors of C are with high probability approximations to the left singular vectors of A

### Main theoretical results

Let  $H_k$  be constructed the linear Time SVD

$$\mathbf{E}[\|A - H_k H_k^T A\|_F^2] \le \|A - A_k\|_F^2 + \epsilon \|A\|_F^2$$

- Exact SVD of  $A = U\Sigma V^T$ ,  $A_k = U_k \Sigma_k V_k^T = U_k U_k^T A = AV_k V_k^T$ .  $\min_{\operatorname{rank}(B) \le k} ||A - B||_2 = ||A - A_k||_2 = \sigma_{k+1}(A)$  $\min_{\operatorname{rank}(B) \le k} ||A - B||_F^2 = ||A - A_k||_F^2 = \sum_{t=k+1}^r \sigma_t^2(A)$
- perturbation theory of matrices

$$\max_{1 \le t \le n} |\sigma_t(A + E) - \sigma_t(A)| \le ||E||_2, \quad \sum_{k=1}^n (\sigma_k(A + E) - \sigma_k(A))^2 \le ||E||_F^2$$

the latter is known as Hoffman-Wielandt inequality

• Exact SVD of 
$$C = H \Sigma_C Y^T$$

Lemma:

$$\begin{aligned} \|A - H_k H_k^T A\|_F^2 &\leq \|A - A_k\|_F^2 + 2\sqrt{k} \|AA^T - CC^T\|_F \\ \|A - H_k H_k^T A\|_2^2 &\leq \|A - A_k\|_2^2 + 2\|AA^T - CC^T\|_2 \end{aligned}$$

Proof of the first inequality

• 
$$||X||_F^2 = \operatorname{Tr}(X^T X)$$
 and  $\operatorname{Tr}(X + Y) = \operatorname{Tr}(X) + \operatorname{Tr}(Y)$   
 $||A - H_k H_k^T A||_F^2 = \operatorname{Tr}((A - H_k H_k^T A)^T (A - H_k H_k^T A)))$   
 $= \operatorname{Tr}(A^T A) - \operatorname{Tr}(A^T H_k H_k^T A) = ||A||_F^2 - ||A^T H_k||_F^2$ 

• Using Cauchy-Schwartz inequality:

$$\left| \|A^{T}H_{k}\|_{F}^{2} - \sum_{t=1}^{k} \sigma_{t}^{2}(C) \right| \leq \sqrt{k} \left( \sum_{t=1}^{k} (\|A^{T}h_{t}\|^{2} - \sigma_{t}^{2}(C))^{2} \right)^{1/2}$$

$$= \sqrt{k} \left( \sum_{t=1}^{k} (\|A^{T}h_{t}\|^{2} - \|C^{T}h_{t}\|^{2})^{2} \right)^{1/2} = \sqrt{k} \left( \sum_{t=1}^{k} ((h_{t})^{T}(AA^{T} - CC^{T})h_{t})^{2} \right)$$

$$\leq \sqrt{k} \|AA^{T} - CC^{T}\|_{F}$$

$$= \sqrt{k} \left( \sum_{t=1}^{k} (\|A^{T}h_{t}\|^{2} - \|C^{T}h_{t}\|^{2})^{2} \right)^{1/2} = \sqrt{k} \left( \sum_{t=1}^{k} (\|A^{T}h_{t}\|^{2} - \|C^{T}h_{t}\|^{2})^{2} \right)^{1/2}$$

• by Hoffman-Wielandt inequality

$$\begin{aligned} \left| \sum_{t=1}^{k} \sigma_{t}^{2}(C) - \sum_{t=1}^{k} \sigma_{t}^{2}(A) \right| &\leq \sqrt{k} \left( \sum_{t=1}^{k} (\sigma_{t}^{2}(C) - \sigma_{t}^{2}(A))^{2} \right)^{1/2} \\ &= \sqrt{k} \left( \sum_{t=1}^{k} (\sigma_{t}(CC^{T}) - \sigma_{t}(AA^{T}))^{2} \right)^{1/2} \\ &\leq \sqrt{k} \left( \sum_{t=1}^{m} (\sigma_{t}(CC^{T}) - \sigma_{t}(AA^{T}))^{2} \right)^{1/2} \leq \sqrt{k} \|CC^{T} - AA^{T}\|_{F} \end{aligned}$$

• Therefore

$$\left| \|A^{T}H_{k}\|_{F}^{2} - \sum_{t=1}^{k} \sigma_{t}^{2}(A) \right| \leq 2\sqrt{k} \|AA^{T} - CC^{T}\|_{F}$$

٠

matrix approximation gives

$$\mathbf{E}[\|AB - CR\|_{F}^{2}] \leq \frac{1}{c} \|A\|_{F}^{2} \|B\|_{F}^{2}$$

which yields

$$2\sqrt{k}\mathbf{E}[\|AA^{T} - CC^{T}\|_{F}] \leq \left(\frac{4k}{c}\right)^{1/2} \|A\|_{F}^{2}$$

$$||A^{T}H_{k}||_{F}^{2} \ge \sum_{t=1}^{k} \sigma_{t}^{2}(A) - 2\sqrt{k}||AA^{T} - CC^{T}||_{F}$$

• If  $c \ge 4k/\epsilon^2$ , then

$$\mathbf{E}[\|A - H_k H_k^T A\|_F^2] \leq \|A\|_F^2 - \sum_{t=1}^k \sigma_t^2(A) + 2\sqrt{k} \mathbf{E}[\|AA^T - CC^T\|_F]$$

$$\leq \|A - A_k\|_F^2 + \epsilon \|A\|_F^2$$

$$= \sum_{t=1}^k \sigma_t^2(A) + 2\sqrt{k} \mathbf{E}[\|AA^T - CC^T\|_F]$$

$$\leq \|A - A_k\|_F^2 + \epsilon \|A\|_F^2$$

$$= \sum_{t=1}^k \sigma_t^2(A) + 2\sqrt{k} \mathbf{E}[\|AA^T - CC^T\|_F]$$

$$\leq \|A - A_k\|_F^2 + \epsilon \|A\|_F^2$$

$$= \sum_{t=1}^k \sigma_t^2(A) + 2\sqrt{k} \mathbf{E}[\|AA^T - CC^T\|_F]$$

$$\leq \|A - A_k\|_F^2 + \epsilon \|A\|_F^2$$

$$= \sum_{t=1}^k \sigma_t^2(A) + 2\sqrt{k} \mathbf{E}[\|AA^T - CC^T\|_F]$$

$$\leq \|A - A_k\|_F^2 + \epsilon \|A\|_F^2$$

$$= \sum_{t=1}^k \sigma_t^2(A) + 2\sqrt{k} \mathbf{E}[\|AA^T - CC^T\|_F]$$

$$\leq \|A - A_k\|_F^2 + \epsilon \|A\|_F^2$$

$$= \sum_{t=1}^k \sigma_t^2(A) + 2\sqrt{k} \mathbf{E}[\|AA^T - CC^T\|_F]$$

$$\leq \|A - A_k\|_F^2 + \epsilon \|A\|_F^2$$

$$= \sum_{t=1}^k \sigma_t^2(A) + 2\sqrt{k} \mathbf{E}[\|AA^T - CC^T\|_F]$$

### Outline



- Approximating Matrix Multiplication
- 3 Approximate SVD
- 4

#### Random Sampling for SVD

5 Single View Algorithm For Matrix Approximation

# Range finding problem

Given an  $m \times n$  matrix A and an integer  $k < \min(m, n)$ , find an orthonormal  $m \times k$  matrix Q such that

$$A \approx Q Q^T A$$

Solving the primitive problem via randomized sampling — intuition

- Draw random vectors  $r_1, r_2, \ldots, r_k \in \mathbb{R}^n$ .
- Form "sample" vectors  $y_1 = Ar_1, y_2 = Ar_2, \dots, y_k = Ar_k \in \mathbb{R}_m$ .
- Form orthonormal vectors  $q_1, q_2, \ldots, q_k \in \mathbb{R}^m$  such that

$$\operatorname{span}\{q_1, q_2, \ldots, q_k\} = \operatorname{span}\{y_1, y_2, \ldots, y_k\}$$

Almost always correct if A has exact rank k

# Low-Rank Approximation: Randomized Sampling

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ● □ ● ○ ○ ○

35/49

### Algorithm RandSam0

Input: mxn matrix A, int k, p.

- Draw a random  $n \times (k + p)$  matrix  $\Omega$
- Compute  $QR = A\Omega$
- and SVD:  $Q^T A = \hat{U} \hat{\Sigma} \hat{V}^T$

Fruncate SVD:  $\hat{U}_k \hat{\Sigma}_k \hat{V}_k^T$ Output:  $B = (Q\hat{U}_k)\hat{\Sigma}_k \hat{V}_k^T$ 

- Easy to implement.
- Very efficient computation.
- Minimum communication.

### error for Gaussian test matrices

Ref: Halko, Martinsson, Tropp, 2009 & 2011; Martinsson, Rokhlin, Tygert (2006)

- Let A denote an  $m \times n$  matrix with singular values  $\{\sigma_j\}_{j=1}^{\min(m,n)}$
- Let k denote a target rank and let p denote an over-sampling parameter.
- Let  $\Omega$  denote an  $n \times (k + p)$  Gaussian matrix.

• Let Q denote the  $m \times (k + p)$  matrix  $Q = orth(A\Omega)$ . If  $p \ge 4$ , then

$$\|A - QQ^*A\|_2 \le \left(1 + 6\sqrt{(k+p)p\log p}\right)\sigma_{k+1} + 3\sqrt{k+p}\left(\sum_{j>k}\sigma_j^2\right)^{1/2}$$

except with probability at most  $3p^{-p}$ .

# Improved Randomized Sampling

### Algorithm RandSam1

Input: mxn matrix A, int k, p, c.

- Draw a random  $n \times (k + p + c)$  matrix  $\Omega$
- Compute  $QR = A\Omega$
- and SVD:  $Q^T A = \hat{U} \hat{\Sigma} \hat{V}^T$

Fruncate SVD:  $\hat{U}_k \hat{\Sigma}_k \hat{V}_k^T$ Output:  $B = (Q\hat{U}_k)\hat{\Sigma}_k \hat{V}_k^T$ 

- Only change from RandSam0: p becomes p + c
- Smallest modification of any algorithm.
- c allows a drastically different error bound, controls accuracy.
- p remains in control of failure chance.

# Randomized Power Method

### Algorithm RandSam2

Input: mxn matrix A, int k, p, c, q

- Draw a random  $n \times (k + p + c)$  matrix  $\Omega$
- Compute  $QR = (AA^T)^q A\Omega$
- and SVD:  $Q^T A = \hat{U} \hat{\Sigma} \hat{V}^T$

Fruncate SVD:  $\hat{U}_k \hat{\Sigma}_k \hat{V}_k^T$ 

Output:  $B = (Q\hat{U}_k)\hat{\Sigma}_k\hat{V}_k^T$ 

- QR needs done carefully for numerical accuracy.
- Algorithm is old one when q = 0; but q = 1 far more accurate.
- Should converge faster when singular values do not decay very fast.

### Example 1

We consider a  $1000 \times 1000$  matrix A whose singular values are shown below:



A is a discrete approximation of a certain compact integral operator normalized so that ||A|| = 1. Curiously, the nature of *A* is in a strong sense irrelevant: the error distribution depends only on  $\{\sigma_j\}_{j=1,...,n}^{\min(m,n)}$ .

### Example 2

We consider a  $1000 \times 1000$  matrix A whose singular values are shown below:



A is a discrete approximation of a certain compact integral operator normalized so that ||A|| = 1. Curiously, the nature of *A* is in a strong sense irrelevant: the error distribution depends only on  $\{\sigma_j\}_{j=1,...,n}^{\min(m,n)}$ .

### Example 3

The matrix A being analyzed is a  $9025 \times 9025$  matrix arising in a diffusion geometry approach to image processing.

To be precise, A is a graph Laplacian on the manifold of  $3 \times 3$  patches.





The pink lines illustrates the performance of the basic random sampling scheme. The errors for q = 0 are huge, and the estimated eigenvalues are much too small.

## Outline

- Randomized Numerical Linear Algebra (RandNLA)
- Approximating Matrix Multiplication
- 3 Approximate SVD
- 4 Random Sampling for SVD

#### 5 Single View Algorithm For Matrix Approximation

### Low-rank reconstruction

Given  $A \in \mathbb{R}^{m \times n}$  and a target rank *r*. Select *k* and  $\ell$ . Given random matrices  $\Omega \in \mathbb{R}^{n \times k}$  and  $\Psi \in \mathbb{R}^{\ell \times m}$ . Compute

 $Y = A\Omega, \qquad W = \Psi A,$ 

Then an approximation  $\hat{A}$  is computed:

- Form an orthogonal-triangular factorization Y = QR where  $Q \in \mathbb{R}^{m \times k}$ .
- Solve a least-squares problem to obtain  $X = (\Psi Q)^{\dagger} W \in \mathbb{R}^{k \times n}$
- Construct the rank-k approximation  $\hat{A} = QX$

Suppose k = 2r + 1 and  $\ell = 4r + 2$ , then

$$\mathbf{E} \|A - \hat{A}\|_F \le 2 \min_{\operatorname{rank}(Z) \le r} \|A - Z\|_F$$

### Linear update of A

• Suppose that *A* is sent as a sequence of additive updates:

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

Then one compute

$$Y \leftarrow Y + H\Omega, \quad W = W + \Psi H$$

• Suppose that *A* is sent as a sequence of additive updates:

$$A = \theta A + \eta H$$

Then one compute

$$Y \leftarrow \theta Y + \eta H\Omega, \quad W = \theta W + \eta \Psi H$$

### Intuition

Suppose

$$A \approx QQ^*A.$$

We want to form the rank-k approximation  $Q(Q^*A)$ , but we cannot compute the factor  $Q^*A$  without revisiting the target matrix A.

Note

$$W = \Psi(QQ^*A) + \Psi(A - QQ^*A) \approx (\Psi Q)(Q^*A)$$

• The construction of *X*:

$$X = (\Psi Q)^{\dagger} W \approx Q^* A$$

Hence

$$\hat{A} = QX \approx QQ^*A \approx A$$

(ロ) (型) (主) (主) (主) (A)(2) (46/49)

### Projection onto a Convex Set.

Let *C* be a closed and convex set. Define the projection:

$$\Pi_C(M) = \arg\min_X \quad \|X - M\|_F^2, \text{ s.t. } X \in C$$

• Suppose  $A \in C$ . Let  $\hat{A}_{in}$  be an initial approximation of A,

$$||A - \Pi_C(\hat{A}_{in})||_F \le ||A - \hat{A}_{in}||_F$$

Conjugate Symmetric Approximation

$$H^n = \{X \in \mathcal{C}^{n \times n} | X = X^*\}$$

The projection

$$\Pi_{H^n}(M) = \frac{1}{2}(M+M^*)$$

### Conjugate Symmetric Approximation.

Let 
$$A \in H^n$$
. Let  $\hat{A} = QX$ .  
•  $\Pi_{H^n}(\hat{A}) = \frac{1}{2}(QX + X^*Q^*) = \frac{1}{2}[Q, X^*] \begin{pmatrix} 0 & I \\ I & 0 \end{pmatrix} [Q, X^*]^*$   
• Let  $[Q, X^*] = U[T_1, T_2]$ . Then  
 $S = \frac{1}{2}(T_1T_2^* + T_2T_1^*)$   
• Construct

$$\hat{A}_{sym} = USU^*$$

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ○ □ ○ ○ ○ ○

48/49

# **PSD** Approximation

Let *A* be positive semidefinite (PSD). Let  $\hat{A} = QX$ .

• Form eigenvalue decomposition

 $S = VDV^*$ 

• Compute  $\hat{A}_{sym} = (UV)D(UV)^*$ • Construct  $\hat{i}$   $\Pi_{sym} = (\hat{i}) - (UU)D(UV)^*$ 

$$\hat{A}_{+} = \prod_{H^{n}_{+}}(\hat{A}) = (UV)D_{+}(UV)^{*}$$