

## Lecture 26: Randomized Numerical Linear Algebra III

Lecturer: Yudong Chen

Scribe: Matthew Zurek

In this lecture,<sup>1,2</sup> we will discuss randomized algorithms for low-rank matrix approximation. Making use of results for approximate matrix multiplication and least squares from the previous two lectures, we give several algorithms and analyze their differing error guarantees.

## 1 Recap

In the previous two lectures, we discussed randomized algorithms for matrix multiplication and least squares problems which beat the runtimes of the standard approaches for these problems, at the cost of having worse accuracy. First we review the high-level strategies of these algorithms.

- Matrix multiplication
  - We randomly subsample columns of  $\mathbf{A}$  and rows of  $\mathbf{B}$ , and use the outer products  $\mathbf{A}_{:,i}\mathbf{B}_{i,:}$  to approximate  $\mathbf{AB}$ .
  - Ideally we should sample each row/column pair  $\mathbf{A}_{:,i}, \mathbf{B}_{i,:}$  with probability proportional to the product of their norms,  $\|\mathbf{A}_{:,i}\|_2 \|\mathbf{B}_{i,:}\|_2$ .
- Least squares problems
  - Given the least squares problem  $\text{minimize}_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{Ax} - \mathbf{b}\|_2^2$  with  $\mathbf{A} \in \mathbb{R}^{n \times d}$  and  $n \gg d$ , we construct  $\Phi \in \mathbb{R}^{r \times n}$  to reduce the dimensionality of  $\mathbf{A}$  and solve the smaller problem of  $\text{minimize}_{\mathbf{x} \in \mathbb{R}^d} \|\Phi(\mathbf{Ax} - \mathbf{b})\|_2^2$ .
  - We covered two approaches to choosing  $\Phi$  which are based on the uniform and non-uniform subsampling strategies, respectively.
    1. Subsampled randomized Hadamard transform (SRHT): First transform  $\mathbf{A}$  to make the leverage scores approximately uniform, and then we can subsample the rows uniformly at random.
    2. We estimate the leverage scores of  $\mathbf{A}$  and sample each row with probability proportional to its leverage score. Naively trying to compute the leverage scores exactly would take SVD time  $\Omega(nd^2)$ , but we can still make this idea work by approximating the leverage scores in a way which made essential use of the SRHT. Ultimately  $\gtrsim (d \text{ polylog } n)/\varepsilon$  rows are needed for an  $\varepsilon$ -approximate answer. If  $\mathbf{A}$  is low-rank then a more careful analysis than what we saw in class allows us to replace  $d$  with  $\text{rank}(\mathbf{A})$ .

## 2 Low-Rank Matrix Approximation

Given a matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , our goal is to compute a rank- $k$  matrix which is a good approximation to  $\mathbf{A}$ . First we review the baseline approach.

- One can compute the SVD of  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^\top$  and then truncate all singular values besides the top  $k$ . The matrix we return is

$$\mathbf{A}_k := \mathbf{U}_k \mathbf{U}_k^\top \mathbf{A}$$

where  $\mathbf{U}_k$  consists of top  $k$  singular vectors (note  $\mathbf{U}_k \mathbf{U}_k^\top$  projects onto the subspace spanned by the top  $k$  singular vectors).

<sup>1</sup>Reading: [Mahoney, 2016]

<sup>2</sup>We thank Yuxin Chen for allowing us to draw from his slides for ELE 520: Mathematics of Data Science, Princeton University, Fall 2020.

- This has optimal error with respect to the Frobenius norm:  $\|\mathbf{A}_k - \mathbf{A}\|_F = \min_{\mathbf{Z}: \text{rank}(\mathbf{Z}) \leq k} \|\mathbf{Z} - \mathbf{A}\|_F$ .  $\mathbf{A}_k$  is also optimal with respect to the spectral norm.
- In general this takes time  $O(n^3)$  (by direct methods) or  $O(n^2k)$  (by power methods). Analogous to least squares problems, direct methods have very high accuracy, while power methods only rely on matrix vector multiplications (and thus can be especially fast for sparse  $\mathbf{A}$ ) but suffer from dependencies on the spectrum of  $\mathbf{A}$ . Power methods also often depend like  $O(\log(1/\varepsilon))$  on the accuracy  $\varepsilon$ .

Our goal is to find a faster randomized algorithm with a time complexity of  $O(n^2 \log k)$  and a worse accuracy dependence like  $1/\varepsilon$ .

## 2.1 Basic Algorithmic Strategy

Instead of directly taking the SVD of  $\mathbf{A}$ , we first subsample it. In our subsampling procedure, we sample  $r$  columns of  $\mathbf{A}$  into  $\mathbf{C} \in \mathbb{R}^{n \times r}$ , and then return

$$\mathbf{C}\mathbf{C}^\dagger\mathbf{A}.$$

Note this is the projection of  $\mathbf{A}$  onto the column space of  $\mathbf{C}$ . It can be computed quickly (relative to taking the SVD of  $\mathbf{A}$ ) because  $\mathbf{C}$  is  $n \times r$  ( $r < n$ ) and thus

- the SVD/pseudoinverse of  $\mathbf{C}$  can be computed in  $O(nr^2)$  time
- the product  $\underbrace{\mathbf{C} \mathbf{C}^\dagger}_{n \times r \quad r \times n} \mathbf{A}$  can be computed in  $O(n^2r)$  time.

While  $\mathbf{C}\mathbf{C}^\dagger\mathbf{A}$  is an approximation of  $\mathbf{A}$ , it is not (necessarily) rank  $k$ . Therefore instead of projecting onto the column space of  $\mathbf{C}$ , we project onto the subspace spanned by the top  $k$  left singular vectors of  $\mathbf{C}$ . This gives our first generic randomized low-rank approximation algorithm.

**Algorithm 1:** Generic low-rank approximation algorithm

- 1: **input:** data matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , subsampled matrix  $\mathbf{C} \in \mathbb{R}^{n \times r}$
  - 2: Compute top  $k$  left singular vectors  $\mathbf{H}_k$  of  $\mathbf{C}$
  - 3: **return**  $\mathbf{H}_k\mathbf{H}_k^\top\mathbf{A}$

Notice that  $\mathbf{C}\mathbf{C}^\dagger\mathbf{A}$  has the form  $\mathbf{C}\mathbf{X}$  and thus can be interpreted as using linear combinations of a few columns of  $\mathbf{A}$  to approximate each column of  $\mathbf{A}$ . On the other hand,  $\mathbf{H}_k\mathbf{H}_k^\top\mathbf{A}$  is actually rank  $k$  but is less interpretable. Nonetheless, the approximation error of  $\mathbf{H}_k\mathbf{H}_k^\top\mathbf{A}$  will still depend on  $\mathbf{C}$ , specifically on the norm of  $\mathbf{A}\mathbf{A}^\top - \mathbf{C}\mathbf{C}^\top$  as we will soon see.

From here, our approach to developing the generic algorithm into a concrete approach resembles our steps from the least squares problem:

1. First we find a deterministic condition on  $\mathbf{C}$  which guarantees good approximation.
2. Then we construct  $\mathbf{C}$  to satisfy this condition.

## 2.2 Analysis

### 2.2.1 Frobenius Approximation Error

To find desirable conditions for  $\mathbf{C}$ , first we consider the Frobenius norm approximation error.

**Lemma 1.** *The output of Algorithm 1 satisfies*

$$\|\mathbf{A} - \mathbf{H}_k \mathbf{H}_k^\top \mathbf{A}\|_{\mathbb{F}}^2 \leq \underbrace{\|\mathbf{A} - \mathbf{U}_k \mathbf{U}_k^\top \mathbf{A}\|_{\mathbb{F}}^2}_{\text{best rank-}k \text{ approx. error}} + 2\sqrt{k} \underbrace{\|\mathbf{A}\mathbf{A}^\top - \mathbf{C}\mathbf{C}^\top\|_{\mathbb{F}}}_{\text{excess error}}$$

where  $\mathbf{U}_k \in \mathbb{R}^{n \times k}$  contains top- $k$  left singular vectors of  $\mathbf{A}$ .

For the high-level idea behind the proof, imagine that we could apply the triangle inequality to write

$$\|\mathbf{A} - \mathbf{H}_k \mathbf{H}_k^\top \mathbf{A}\|_{\mathbb{F}} \leq \|\mathbf{A} - \mathbf{U}_k \mathbf{U}_k^\top \mathbf{A}\|_{\mathbb{F}} + \|\mathbf{U}_k \mathbf{U}_k^\top \mathbf{A} - \mathbf{H}_k \mathbf{H}_k^\top \mathbf{C}\|_{\mathbb{F}} + \|\mathbf{H}_k \mathbf{H}_k^\top \mathbf{C} - \mathbf{H}_k \mathbf{H}_k^\top \mathbf{A}\|_{\mathbb{F}}.$$

The third term is  $\|\mathbf{H}_k \mathbf{H}_k^\top (\mathbf{A} - \mathbf{C})\|_{\mathbb{F}}$  (where  $\mathbf{H}_k \mathbf{H}_k^\top$  is a projection onto a rank  $k$  subspace), and the second term could hopefully be put into a similar form after applying a singular vector perturbation theorem to relate  $\mathbf{U}_k$  to  $\mathbf{H}_k$ . Of course this is not a correct proof, as  $\mathbf{C}$  and  $\mathbf{A}$  do not have the same dimensions. Now we provide a proof.

**Proof of Lemma 1** To begin with, since  $\mathbf{H}_k$  is orthonormal, one has

$$\|\mathbf{A} - \mathbf{H}_k \mathbf{H}_k^\top \mathbf{A}\|_{\mathbb{F}}^2 = \|\mathbf{A}\|_{\mathbb{F}}^2 - \|\mathbf{H}_k^\top \mathbf{A}\|_{\mathbb{F}}^2.$$

Next, letting  $\mathbf{h}_i = (\mathbf{H}_k)_{:,i}$  yields

$$\begin{aligned} \left| \|\mathbf{H}_k^\top \mathbf{A}\|_{\mathbb{F}}^2 - \sum_{i=1}^k \sigma_i^2(\mathbf{C}) \right| &= \left| \sum_{i=1}^k \|\mathbf{A}^\top \mathbf{h}_i\|_2^2 - \sum_{i=1}^k \|\mathbf{C} \mathbf{h}_i\|_2^2 \right| \\ &= \left| \sum_{i=1}^k \langle \mathbf{h}_i \mathbf{h}_i^\top, \mathbf{A}\mathbf{A}^\top - \mathbf{C}\mathbf{C}^\top \rangle \right| \\ &= |\langle \mathbf{H}_k \mathbf{H}_k^\top, \mathbf{A}\mathbf{A}^\top - \mathbf{C}\mathbf{C}^\top \rangle| \\ &\leq \|\mathbf{H}_k \mathbf{H}_k^\top\|_{\mathbb{F}} \|\mathbf{A}\mathbf{A}^\top - \mathbf{C}\mathbf{C}^\top\|_{\mathbb{F}} \\ &\leq \sqrt{k} \|\mathbf{A}\mathbf{A}^\top - \mathbf{C}\mathbf{C}^\top\|_{\mathbb{F}}. \end{aligned}$$

In addition,

$$\begin{aligned} \left| \sum_{i=1}^k \sigma_i^2(\mathbf{C}) - \sum_{i=1}^k \sigma_i^2(\mathbf{A}) \right| &= \left| \sum_{i=1}^k \{\sigma_i(\mathbf{C}\mathbf{C}^\top) - \sigma_i(\mathbf{A}\mathbf{A}^\top)\} \right| \\ &\leq \sqrt{k} \sqrt{\sum_{i=1}^n \{\sigma_i(\mathbf{C}\mathbf{C}^\top) - \sigma_i(\mathbf{A}\mathbf{A}^\top)\}^2} \quad (\text{Cauchy-Schwarz}) \\ &\leq \sqrt{k} \|\mathbf{C}\mathbf{C}^\top - \mathbf{A}\mathbf{A}^\top\|_{\mathbb{F}} \quad (\text{Wielandt-Hoffman inequality}). \end{aligned}$$

Finally, one has  $\|\mathbf{A} - \mathbf{U}_k \mathbf{U}_k^\top \mathbf{A}\|_{\mathbb{F}}^2 = \|\mathbf{A}\|_{\mathbb{F}}^2 - \sum_{i=1}^k \sigma_i^2(\mathbf{A})$ . Combining these results establishes the claim.  $\square$

## 2.2.2 Spectral Approximation Error

A similar bound holds for approximation error in the spectral norm.

**Lemma 2.** *The output of Algorithm 1 satisfies*

$$\|\mathbf{A} - \mathbf{H}_k \mathbf{H}_k^\top \mathbf{A}\|^2 \leq \|\mathbf{A} - \mathbf{U}_k \mathbf{U}_k^\top \mathbf{A}\|^2 + 2\|\mathbf{A}\mathbf{A}^\top - \mathbf{C}\mathbf{C}^\top\|$$

where  $\mathbf{U}_k \in \mathbb{R}^{n \times k}$  contains the top  $k$  left singular vectors of  $\mathbf{A}$ .

**Proof of Lemma 2** First of all,

$$\begin{aligned}\|\mathbf{A} - \mathbf{H}_k \mathbf{H}_k^\top \mathbf{A}\| &= \max_{\mathbf{x}: \|\mathbf{x}\|_2=1} \|\mathbf{x}^\top (\mathbf{I} - \mathbf{H}_k \mathbf{H}_k^\top) \mathbf{A}\|_2 \\ &= \max_{\mathbf{x}: \|\mathbf{x}\|_2=1, \mathbf{x} \perp \mathbf{H}_k} \|\mathbf{x}^\top \mathbf{A}\|_2\end{aligned}$$

Additionally, for any  $\mathbf{x} \perp \mathbf{H}_k$ ,

$$\begin{aligned}\|\mathbf{x}^\top \mathbf{A}\|_2^2 &= |\mathbf{x}^\top \mathbf{C} \mathbf{C}^\top \mathbf{x} + \mathbf{x}^\top (\mathbf{A} \mathbf{A}^\top - \mathbf{C} \mathbf{C}^\top) \mathbf{x}| \\ &\leq |\mathbf{x}^\top \mathbf{C} \mathbf{C}^\top \mathbf{x}| + |\mathbf{x}^\top (\mathbf{A} \mathbf{A}^\top - \mathbf{C} \mathbf{C}^\top) \mathbf{x}| \\ &\leq \sigma_{k+1}(\mathbf{C} \mathbf{C}^\top) + \|\mathbf{A} \mathbf{A}^\top - \mathbf{C} \mathbf{C}^\top\| \\ &\leq \sigma_{k+1}(\mathbf{A} \mathbf{A}^\top) + 2\|\mathbf{A} \mathbf{A}^\top - \mathbf{C} \mathbf{C}^\top\| \\ &= \|\mathbf{A} - \mathbf{U}_k \mathbf{U}_k^\top \mathbf{A}\|^2 + 2\|\mathbf{A} \mathbf{A}^\top - \mathbf{C} \mathbf{C}^\top\|.\end{aligned}$$

This concludes the proof.  $\square$

## 2.3 Algorithms

### 2.3.1 Non-Uniform Sampling Algorithm

Now we proceed to devise a concrete algorithm for finding  $\mathbf{C}$ . As revealed by Lemmas 1 and 2, we need to make  $\mathbf{A} \mathbf{A}^\top - \mathbf{C} \mathbf{C}^\top$  small. To approximate the product  $\mathbf{A} \mathbf{A}^\top$ , we can use the procedures we developed for matrix multiplication! Concretely, one approach is to use the randomized matrix multiplication algorithm based on non-uniform sampling (proportional to  $\|\mathbf{A}_{:,k}\|_2^2$ ):

**Algorithm 2:** One-pass randomized SVD

- 1: Set  $p_k \geq \frac{\beta \|\mathbf{A}_{:,k}\|_2^2}{\|\mathbf{A}\|_F^2}, 1 \leq k \leq n$
- 2: **for**  $l = 1, \dots, r$  **do**
- 3:   Pick  $i_l \in \{1, \dots, n\}$  i.i.d. with prob.  $\mathbb{P}\{i_l = k\} = p_k$
- 4:   Set  $\mathbf{C}_{:,l} = \frac{1}{\sqrt{r p_{i_l}}} \mathbf{A}_{:,i_l}$
- 5: **return**  $\mathbf{H}_k$  as top- $k$  left singular vectors of  $\mathbf{C}$

As stated, Algorithm 2 returns  $\mathbf{H}_k$ , which are approximate top  $k$  left singular vectors. We could repeat to also produce right singular vectors, and hence the above procedure can be viewed as a *randomized SVD* algorithm. Of course we can also use  $\mathbf{H}_k$  to construct the low-rank approximation  $\mathbf{H}_k \mathbf{H}_k^\top \mathbf{A}$ .

Invoking our theorems from Lecture 24 on the Frobenius and spectral error of non-uniform sampling for matrix multiplication, we have that with high probability:

- If  $r \gtrsim \frac{k \log n}{\beta \varepsilon^2}$ , then

$$\|\mathbf{A} - \mathbf{H}_k \mathbf{H}_k^\top \mathbf{A}\|_F^2 \leq \|\mathbf{A} - \mathbf{U}_k \mathbf{U}_k^\top \mathbf{A}\|_F^2 + \varepsilon \|\mathbf{A}\|_F^2. \quad (1)$$

- If  $r \gtrsim \frac{\|\mathbf{A}\|_F^2 \log n}{\|\mathbf{A}\|_2^2 \beta \varepsilon^2}$ , then

$$\|\mathbf{A} - \mathbf{H}_k \mathbf{H}_k^\top \mathbf{A}\|^2 \leq \|\mathbf{A} - \mathbf{U}_k \mathbf{U}_k^\top \mathbf{A}\|^2 + \varepsilon \|\mathbf{A}\|^2. \quad (2)$$

### 2.3.2 An Improved Multi-Pass Algorithm

If we make more passes over the matrix  $\mathbf{A}$ , we can improve upon these error guarantees. Instead of stopping after forming our first low-rank approximation  $\hat{\mathbf{A}}$  to  $\mathbf{A}$ , we can recursively apply Algorithm 2 to the residual  $\mathbf{A} - \hat{\mathbf{A}}$ . This will cause the errors to decrease geometrically in the number  $t$  of recursive steps and will use a total of  $rt$  sampled columns.

**Algorithm 3:** Multi-pass randomized SVD

```

1:  $\mathcal{S} = \{\}$ 
2: for  $l = 1, \dots, t$  do
3:    $\mathbf{E}_l = \mathbf{A} - \mathbf{A}_{\mathcal{S}} \mathbf{A}_{\mathcal{S}}^\dagger \mathbf{A}$ 
4:   Set  $p_k \geq \frac{\beta \|(\mathbf{E}_l)_{:,k}\|_2^2}{\|\mathbf{E}_l\|_F^2}$ ,  $1 \leq k \leq n$ 
5:   Randomly select  $r$  column indices with sampling prob.  $\{p_k\}$  and append to  $\mathcal{S}$ 
6: return  $\mathbf{C} = \mathbf{A}_{\mathcal{S}}$ 

```

**Theorem 3.** Suppose  $r \gtrsim \frac{k \log n}{\beta \varepsilon^2}$ . With high probability,

$$\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger \mathbf{A}\|_F^2 \leq \frac{1}{1 - \varepsilon} \|\mathbf{A} - \mathbf{U}_k \mathbf{U}_k^\top \mathbf{A}\|_F^2 + \varepsilon^t \|\mathbf{A}\|_F^2.$$

First we provide a sketch of the proof.

**Sketch of Proof** After constructing  $\mathbf{C}$  the first time we have  $\mathbf{A} = \mathbf{C}\mathbf{C}^\dagger \mathbf{A} + \mathbf{E}$ , where  $\mathbf{E}$  satisfies the error guarantee

$$\|\mathbf{E}\|_F^2 \leq \|\mathbf{A} - \mathbf{A}_k\|_F^2 + \varepsilon \|\mathbf{A}\|_F^2.$$

After choosing the next set of columns to form  $\mathbf{C}'$ , we have that  $\mathbf{E} = \mathbf{C}'(\mathbf{C}')^\dagger \mathbf{E} + \mathbf{E}'$ , where  $\mathbf{E}'$  satisfies the error guarantee

$$\|\mathbf{E}'\|_F^2 \leq \|\mathbf{E} - \mathbf{E}_k\|_F^2 + \varepsilon \|\mathbf{E}\|_F^2$$

and  $\mathbf{E}_k$  is the best rank- $k$  approximation to  $\mathbf{E}$ . We can show that  $\|\mathbf{E} - \mathbf{E}_k\|_F^2 \leq \|\mathbf{A} - \mathbf{A}_k\|_F^2$ , and then substituting into the above inequality yields

$$\begin{aligned} \|\mathbf{E}'\|_F^2 &\leq \|\mathbf{A} - \mathbf{A}_k\|_F^2 + \varepsilon (\|\mathbf{A} - \mathbf{A}_k\|_F^2 + \varepsilon \|\mathbf{A}\|_F^2) \\ &= (1 + \varepsilon) \|\mathbf{A} - \mathbf{A}_k\|_F^2 + \varepsilon^2 \|\mathbf{A}\|_F^2 \\ &\quad \vdots \\ &\leq (1 + \varepsilon + \varepsilon^2 \cdots + \varepsilon^{t-1}) \|\mathbf{A} - \mathbf{A}_k\|_F^2 + \varepsilon^t \|\mathbf{A}\|_F^2 \end{aligned}$$

where to get the final inequality we recursively repeat this procedure  $t$  total times. □

**Proof of Theorem 3:** We will prove by induction. By (1), the theorem holds for  $t = 1$ .

Assume the theorem holds for  $t - 1$ :

$$\left\| \underbrace{\mathbf{A} - \mathbf{C}^{t-1} (\mathbf{C}^{t-1})^\dagger \mathbf{A}}_{:= \mathbf{E}_t} \right\|_F^2 \leq \frac{1}{1 - \varepsilon} \|\mathbf{A} - \mathbf{U}_k \mathbf{U}_k^\top \mathbf{A}\|_F^2 + \varepsilon^{t-1} \|\mathbf{A}\|_F^2,$$

and let  $\mathbf{Z}$  be the matrix of the columns of  $\mathbf{E}_t$  included in the sample. In view of (1),

$$\left\| \mathbf{E}_t - \mathbf{Z}\mathbf{Z}^\dagger \mathbf{E}_t \right\|_F^2 \leq \|\mathbf{E}_t - (\mathbf{E}_t)_k\|_F^2 + \varepsilon \|\mathbf{E}_t\|_F^2,$$

where  $(\mathbf{E}_t)_k$  is the best rank- $k$  approximation of  $\mathbf{E}_t$ . Combining the above two inequalities yields

$$\begin{aligned} \left\| \mathbf{E}_t - \mathbf{Z}\mathbf{Z}^\dagger \mathbf{E}_t \right\|_{\mathbb{F}}^2 &\leq \left\| \mathbf{E}_t - (\mathbf{E}_t)_k \right\|_{\mathbb{F}}^2 \\ &+ \frac{\varepsilon}{1-\varepsilon} \left\| \mathbf{A} - \mathbf{U}_k \mathbf{U}_k^\top \mathbf{A} \right\|_{\mathbb{F}}^2 + \varepsilon^t \left\| \mathbf{A} \right\|_{\mathbb{F}}^2. \end{aligned} \quad (3)$$

We claim (and will prove later) that

$$\mathbf{E}_t - \mathbf{Z}\mathbf{Z}^\dagger \mathbf{E}_t = \mathbf{A} - \mathbf{C}^t (\mathbf{C}^t)^\dagger \mathbf{A} \quad (4)$$

and

$$\left\| \mathbf{E}_t - (\mathbf{E}_t)_k \right\|_{\mathbb{F}}^2 \leq \left\| \mathbf{A} - \mathbf{A}_k \right\|_{\mathbb{F}}^2. \quad (5)$$

Substituting into (3) yields

$$\begin{aligned} \left\| \mathbf{A} - \mathbf{C}^t (\mathbf{C}^t)^\dagger \mathbf{A} \right\|_{\mathbb{F}}^2 &\leq \left\| \mathbf{A} - \mathbf{A}_k \right\|_{\mathbb{F}}^2 + \frac{\varepsilon}{1-\varepsilon} \left\| \mathbf{A} - \mathbf{A}_k \right\|_{\mathbb{F}}^2 + \varepsilon^t \left\| \mathbf{A} - \mathbf{A}_k \right\|_{\mathbb{F}}^2 \\ &= \frac{1}{1-\varepsilon} \left\| \mathbf{A} - \mathbf{A}_k \right\|_{\mathbb{F}}^2 + \varepsilon^t \left\| \mathbf{A} - \mathbf{A}_k \right\|_{\mathbb{F}}^2. \end{aligned}$$

This proves the theorem for  $t$ . It remains to justify (4) and (5).

To prove (4), note that  $\mathbf{Z}\mathbf{Z}^\dagger \mathbf{C}^{t-1} (\mathbf{C}^{t-1})^\dagger = \mathbf{0}$ . This gives

$$\mathbf{C}^t (\mathbf{C}^t)^\dagger = \mathbf{C}^{t-1} (\mathbf{C}^{t-1})^\dagger + \mathbf{Z}\mathbf{Z}^\dagger$$

Hence

$$\begin{aligned} \mathbf{A} - \mathbf{C}^t (\mathbf{C}^t)^\dagger \mathbf{A} &= \mathbf{A} - \mathbf{C}^{t-1} (\mathbf{C}^{t-1})^\dagger \mathbf{A} - \mathbf{Z}\mathbf{Z}^\dagger \mathbf{A} \\ &= \underbrace{\mathbf{A} - \mathbf{C}^{t-1} (\mathbf{C}^{t-1})^\dagger \mathbf{A}}_{:=\mathbf{E}_t} - \underbrace{\mathbf{Z}\mathbf{Z}^\dagger (\mathbf{A} - \mathbf{C}^{t-1} (\mathbf{C}^{t-1})^\dagger \mathbf{A})}_{:=\mathbf{E}_t} \\ &= \mathbf{E}_t - \mathbf{Z}\mathbf{Z}^\dagger \mathbf{E}_t \end{aligned}$$

To show (5), note that  $(\mathbf{E}_t)_k$  is best rank- $k$  approximation of  $\mathbf{E}_t$ . This gives

$$\begin{aligned} \left\| \mathbf{E}_t - (\mathbf{E}_t)_k \right\|_{\mathbb{F}}^2 &= \left\| (\mathbf{I} - \mathbf{C}^{t-1} (\mathbf{C}^{t-1})^\dagger) \mathbf{A} - ((\mathbf{I} - \mathbf{C}^{t-1} (\mathbf{C}^{t-1})^\dagger) \mathbf{A})_k \right\|_{\mathbb{F}}^2 \\ &\leq \left\| (\mathbf{I} - \mathbf{C}^{t-1} (\mathbf{C}^{t-1})^\dagger) \mathbf{A} - (\mathbf{I} - \mathbf{C}^{t-1} (\mathbf{C}^{t-1})^\dagger) \mathbf{A}_k \right\|_{\mathbb{F}}^2 \\ &\quad \text{since } (\mathbf{I} - \mathbf{C}^{t-1} (\mathbf{C}^{t-1})^\dagger) \mathbf{A}_k \text{ is rank-}k \\ &= \left\| (\mathbf{I} - \mathbf{C}^{t-1} (\mathbf{C}^{t-1})^\dagger) (\mathbf{A} - \mathbf{A}_k) \right\|_{\mathbb{F}}^2 \\ &\leq \left\| \mathbf{A} - \mathbf{A}_k \right\|_{\mathbb{F}}^2. \end{aligned}$$

□

## 2.4 Types of Matrix Decomposition

Having given some low-rank matrix approximation algorithms, we quickly take a high-level view of approaches to this problem.

- **CX decomposition:** let  $C \in \mathbb{R}^{n \times r}$  consist of  $r$  *actual* columns of  $A$ , and return

$$\hat{A} = CX$$

for some matrix  $X \in \mathbb{R}^{r \times n}$ . The multipass Algorithm 3 is an example of this type of decomposition.

- **CUR decomposition:** let  $C \in \mathbb{R}^{n \times r}$  consist of  $r$  *actual* columns of  $A$ , let  $R \in \mathbb{R}^{r \times n}$  consist of  $r$  *actual* rows of  $A$ , and return

$$\hat{A} = CUR$$

for some matrix  $U \in \mathbb{R}^{r \times r}$ . This will not be covered in this course but appears in the provided references.

## 2.5 An Example

The lecture audience asked what would go wrong if instead of using non-uniform sampling in Algorithm 2, we simply deterministically returned the columns of  $A$  with the largest norms. We discussed a few reasons.

1. Randomization allows us to create an unbiased estimator of  $A$ . If we instead fixed a deterministic strategy, then informally there will always exist inputs  $A$  on which we do badly.
2. Consider the following example:

$$A = \begin{bmatrix} 100 & 100 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 1 \end{bmatrix}.$$

A deterministic strategy which takes the largest-norm columns will choose the first two columns but miss (much of) the column space spanned by the remaining columns. This is undesirable because each of the first two columns span the same space, and so using both is redundant. The non-uniform sampling strategy, which samples columns proportionally to their squared norms, will also sample the first two columns too heavily while neglecting the remaining columns. This was still good enough to give us an *additive* error bound (the neglected columns have smaller norms), but as we will soon see, different sampling probabilities can achieve a superior *multiplicative* error bound.

## 3 Multiplicative Error Guarantees for CX Decomposition

So far we have developed CX decomposition results of the form

$$\begin{aligned} \|A - CC^\dagger A\|_F^2 &\leq \|A - A_k\|_F^2 + \text{additive error} \\ \|A - CC^\dagger A\|^2 &\leq \|A - A_k\|^2 + \text{additive error} \end{aligned}$$

via techniques from approximate matrix multiplication. A *stronger* guarantee would be multiplicative, such as

$$\|A - CC^\dagger A\|_F \leq (1 + \varepsilon) \|A - A_k\|_F.$$

To see why this might be desirable, consider the case that  $A$  is actually exactly low rank. A multiplicative guarantee would then give zero error (since  $\|A - A_k\|_F = 0$ ) whereas an additive guarantee would not. To achieve multiplicative error guarantees, we will use techniques from least squares instead of matrix multiplication.

### 3.1 Generalized Least Squares Problems

Generalized least squares problems have the form

$$\text{minimize}_{\mathbf{X}} \quad \|\mathbf{B} - \mathbf{A}\mathbf{X}\|_{\text{F}}^2$$

where  $\mathbf{X}$  is a matrix rather than a vector. The optimal solution to this problem generalizes the original least squares solution:  $\mathbf{X}^{\text{ls}} = \mathbf{A}^\dagger \mathbf{B}$ . This is just several independent least squares problems, one for each column of  $\mathbf{X}$  and  $\mathbf{B}$ . We can therefore still apply our randomized least squares techniques. We give a quick outline of how this works:

1. With  $r \gtrsim \frac{\text{rank}(\mathbf{A}) \cdot \log(\text{rank}(\mathbf{A}))}{\epsilon^2}$ , construct a optimally weighted subsampling matrix  $\Phi \in \mathbb{R}^{r \times n}$  (by approximating the leverage scores of  $\mathbf{A}$ ).
2. Compute

$$\widetilde{\mathbf{X}}^{\text{ls}} = (\Phi \mathbf{A})^\dagger \Phi \mathbf{B}.$$

Then informally, with high probability, we have

$$\|\mathbf{B} - \mathbf{A}\widetilde{\mathbf{X}}^{\text{ls}}\|_{\text{F}} \leq (1 + \epsilon) \left\{ \min_{\mathbf{X}} \|\mathbf{B} - \mathbf{A}\mathbf{X}\|_{\text{F}} \right\} \quad (6)$$

$$\|\mathbf{X}^{\text{ls}} - \widetilde{\mathbf{X}}^{\text{ls}}\|_{\text{F}} \leq \frac{\epsilon}{\sigma_{\min}(\mathbf{A}_k)} \left\{ \min_{\mathbf{X}} \|\mathbf{B} - \mathbf{A}\mathbf{X}\|_{\text{F}} \right\}. \quad (7)$$

### 3.2 Another Algorithm for CX Decomposition

Now we give our algorithm for CX decomposition which achieves a multiplicative error.

**Algorithm 4:** Randomized algorithm for CX decomposition

- 1: Compute sampling probabilities  $\{p_i\}_{i=1}^n$ , where  $p_i = \frac{1}{k} \|(\mathbf{U}_{\mathbf{A}_k^\top})_{i,:}\|_2^2$  (leverage scores of  $\mathbf{A}_k^\top$ )
- 2: Use sampling probabilities  $\{p_i\}$  to construct a rescaled random subsampling matrix  $\Phi \in \mathbb{R}^{r \times n}$
- 3: **return**  $\mathbf{C} = \mathbf{A}\Phi^\top$ , consisting of  $r$  columns of  $\mathbf{A}$

This algorithm only gives the  $\mathbf{C}$ , but now we can make use of generalized least squares to choose  $\mathbf{X} = \mathbf{C}^\dagger \mathbf{A} = \arg \min_{\mathbf{X}'} \|\mathbf{A} - \mathbf{C}\mathbf{X}'\|_{\text{F}}$ . Note that the most expensive part of this algorithm is computing  $\mathbf{A}_k$  and its leverage scores, for which again we need to use approximate algorithms.

We have the following guarantee for the above algorithm.

**Theorem 4.** *Suppose  $r \gtrsim \frac{k \log k}{\epsilon^2}$ , then Algorithm 4 yields*

$$\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger \mathbf{A}\|_{\text{F}} \leq (1 + \epsilon) \|\mathbf{A} - \mathbf{A}_k\|_{\text{F}}.$$

**Proof** of Theorem 4:

$$\begin{aligned} \|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger \mathbf{A}\|_{\text{F}} &\leq \|\mathbf{A} - \mathbf{C}(\mathbf{A}_k \Phi^\top)^\dagger \mathbf{A}_k\|_{\text{F}} \\ &\quad \text{since } \mathbf{X}^{\text{ls}} = \mathbf{C}^\dagger \mathbf{A} \text{ minimizes } \|\mathbf{A} - \mathbf{C}\mathbf{X}\|_{\text{F}} \\ &= \|\mathbf{A} - (\mathbf{A}\Phi^\top)(\mathbf{A}_k \Phi^\top)^\dagger \mathbf{A}_k\|_{\text{F}} \\ &\stackrel{(i)}{\leq} (1 + \epsilon) \|\mathbf{A} - \mathbf{A}\mathbf{A}_k^\dagger \mathbf{A}_k\|_{\text{F}} \\ &= (1 + \epsilon) \|\mathbf{A} - \mathbf{A}_k\|_{\text{F}} \end{aligned}$$



where the key inequality (i) follows from approximation guarantee (6) for the generalized least squares problem  $\min_{\mathbf{Y}} \|\mathbf{A} - \mathbf{Y}\mathbf{A}_k\|_F = \|\mathbf{A}_k^\top \mathbf{Y}^\top - \mathbf{A}^\top\|_F$ , as we will now detail. This problem has optimal solution  $(\mathbf{Y}^{\text{ls}})^\top = (\mathbf{A}_k^\top)^\dagger \mathbf{A}^\top$  which has objective value  $\|\mathbf{A}_k^\top (\mathbf{A}_k^\top)^\dagger \mathbf{A}^\top - \mathbf{A}^\top\|_F$ , while the subsampled problem has solution  $(\tilde{\mathbf{Y}}^{\text{ls}})^\top = (\Phi \mathbf{A}_k^\top)^\dagger \Phi \mathbf{A}^\top$ , so by the approximation guarantee its optimal value is within  $(1 + \varepsilon)$  that of the original problem:

$$\|\mathbf{A} - (\mathbf{A}\Phi^\top)(\mathbf{A}_k\Phi^\top)^\dagger \mathbf{A}_k\|_F \leq (1 + \varepsilon) \|\mathbf{A}_k^\top (\mathbf{A}_k^\top)^\dagger \mathbf{A}^\top - \mathbf{A}^\top\|_F.$$

Taking transpose gives the inequality (i). □

### 3.3 More Examples

The key difference between this strategy for choosing  $\mathbf{C}$  and the strategies from our previous Algorithms 2 and 3 is that we sample using leverage scores rather than column norms. As we will now see, this fixes the types of issues that we identified with our earlier example.

- Consider the matrix

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} 100 & 100 & 100 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 100\sqrt{3} & 0 \\ 0 & \sqrt{3} \end{bmatrix} \begin{bmatrix} 1/\sqrt{3} & 1/\sqrt{3} & 1/\sqrt{3} & 0 & 0 & 0 \\ 0 & 0 & 0 & 1/\sqrt{3} & 1/\sqrt{3} & 1/\sqrt{3} \end{bmatrix}. \end{aligned}$$

As we see from its SVD, all the leverage scores (of  $\mathbf{A}^\top$ ) are the same, causing us not to over-emphasize the high-norm columns.

- Consider the matrix

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \sqrt{5} \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1/\sqrt{5} & 1/\sqrt{5} & 1/\sqrt{5} & 1/\sqrt{5} & 1/\sqrt{5} \end{bmatrix}. \end{aligned}$$

The first column is essential to capturing the column space, whereas columns 2 through 6 are interchangeable. This is reflected in the leverage scores (1 for the first column versus  $1/\sqrt{5}$  for the remaining columns).

### 3.4 Final Remarks

- Our final Theorem 4 is a culmination of all of our results on randomized numerical linear algebra, making use of approximate matrix multiplication with nonuniform sampling proportional to leverage scores, in turn requiring the SRHT to approximate the leverage scores, and finally using guarantees from approximate least squares.
- Randomized numerical linear algebra is a rapidly-developing area, and more information is provided in the references.
- A long-term goal of the randomized linear algebra research program is to develop libraries for these problems which could rival LAPACK in not only speed, but also robustness and stability. Despite many theoretical advances, practical implementations are still works in progress.

## References

- [Ailon and Chazelle, 2009] Ailon, N. and Chazelle, B. (2009). The fast johnson–lindenstrauss transform and approximate nearest neighbors. *SIAM Journal on computing*, 39(1):302–322.
- [Drineas et al., 2011] Drineas, P., Mahoney, M. W., Muthukrishnan, S., and Sarlós, T. (2011). Faster least squares approximation. *Numerische mathematik*, 117(2):219–249.
- [Halko et al., 2011] Halko, N., Martinsson, P.-G., and Tropp, J. A. (2011). Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM review*, 53(2):217–288.
- [Magen and Zouzias, 2011] Magen, A. and Zouzias, A. (2011). Low rank matrix-valued chernoff bounds and approximate matrix multiplication. In *Proceedings of the twenty-second annual ACM-SIAM symposium on Discrete Algorithms*, pages 1422–1436. SIAM.
- [Mahoney, 2011] Mahoney, M. W. (2011). Randomized algorithms for matrices and data. *Foundations and Trends® in Machine Learning*, 3(2):123–224.
- [Mahoney, 2016] Mahoney, M. W. (2016). Lecture notes on randomized linear algebra.
- [Tropp, 2011] Tropp, J. A. (2011). Improved analysis of the subsampled randomized hadamard transform. *Advances in Adaptive Data Analysis*, 3(01n02):115–126.