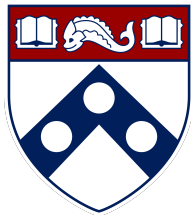


Randomized linear algebra



Yuxin Chen

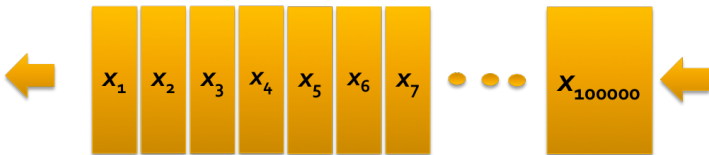
Wharton Statistics & Data Science, Spring 2022

Outline

- Approximate matrix multiplication
- Least squares approximation
- Low-rank matrix approximation

Main reference: "*Lecture notes on randomized linear algebra,*"
Michael W. Mahoney, 2016

Efficient large-scale data processing



When processing large-scale data (in particular, streaming data), we desire methods that can be performed with

- a few (e.g. one or two) passes of data
- limited memory (so impossible to store all data)
- low computational complexity

Key idea: dimension reduction via random sketching

- **random sampling:** randomly downsample data
 - often relies on the information of data
- **random projection:** rotates / projects data to lower dimensions
 - often data-agnostic

Approximate matrix multiplication

Matrix multiplication: a fundamental algebra task

Given $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times p}$, compute or approximate AB

Algorithm 6.1 Vanilla algorithm for matrix multiplication

```
1: for  $i = 1, \dots, m$  do  
2:   for  $k = 1, \dots, n$  do  
3:      $M_{i,k} = A_{i,:} B_{:,k}$   
4: return  $M$ 
```

Computational complexity: $O(mnp)$, or $O(n^3)$ if $m = n = p$

For simplicity, we shall assume $m = n = p$ unless otherwise noted.

Faster matrix multiplication?

- **Strassen algorithms:** exact matrix multiplication
 - Computational complexity $\approx O(n^{2.8})$
 - For various reasons, rarely used in practice
- Approximate solution?

A simple randomized algorithm

View AB as a sum of rank-one matrices (or outer products)

$$AB = \sum_{i=1}^n A_{:,i} B_{i,:}$$

Idea: randomly sample L rank-one components

Algorithm 6.2 Basic randomized algorithm for matrix multiplication

- 1: **for** $l = 1, \dots, L$ **do**
- 2: Pick $i_l \in \{1, \dots, n\}$ i.i.d. with prob. $\mathbb{P}\{i_l = k\} = p_k$
- 3: **return**

$$M = \sum_{l=1}^L \frac{1}{L p_{i_l}} A_{:,i_l} B_{i_l,:}$$

-
- $\{p_k\}$: importance sampling probabilities

A simple randomized algorithm

Rationale: M is an *unbiased* estimate of AB , i.e.

$$\begin{aligned}\mathbb{E}[M] &= \sum_{l=1}^L \sum_k \mathbb{P}\{i_l = k\} \frac{1}{Lp_k} \mathbf{A}_{:,k} \mathbf{B}_{k,:} \\ &= \sum_k \mathbf{A}_{:,k} \mathbf{B}_{k,:} = AB\end{aligned}$$

Clearly, the approximation error (e.g. $\|AB - M\|$) depends on $\{p_k\}$

Importance sampling probabilities

- **Uniform sampling** ($p_k \equiv \frac{1}{n}$): one can choose the sampling set before looking at data, so it's implementable via 1 pass over data

Intuitively, one may prefer biasing towards larger rank-1 components

- **Nonuniform sampling**

$$p_k = \frac{\|\mathbf{A}_{:,k}\|_2 \|\mathbf{B}_{k,:}\|_2}{\sum_l \|\mathbf{A}_{:,l}\|_2 \|\mathbf{B}_{l,:}\|_2}$$

- $\{p_k\}$ can be computed using one pass and $O(n)$ memory

Optimal sampling probabilities?

Let's measure the approximation error by $\mathbb{E} [\|M - AB\|_F^2]$

As it turns out, $\mathbb{E} [\|M - AB\|_F^2]$ is minimized by

$$p_k = \frac{\|A_{:,k}\|_2 \|B_{k,:}\|_2}{\sum_l \|A_{:,l}\|_2 \|B_{l,:}\|_2} \quad (6.1)$$

Thus, we call (6.1) the **optimal sampling probabilities**

Justification of the optimality of (6.1)

Since $\mathbb{E}[M] = AB$, one has

$$\begin{aligned}\mathbb{E} [\|M - AB\|_F^2] &= \mathbb{E} \left[\sum_{i,j} (M_{i,j} - A_{i,:} B_{:,j})^2 \right] = \sum_{i,j} \text{Var}[M_{i,j}] \\ &= \frac{1}{L} \sum_k \sum_{i,j} \frac{A_{i,k}^2 B_{k,j}^2}{p_k} - \frac{1}{L} \sum_{i,j} (A_{i,:} B_{:,j})^2 \quad (\text{check}) \\ &= \frac{1}{L} \sum_k \frac{1}{p_k} \|A_{:,k}\|_2^2 \|B_{k,:}\|_2^2 - \frac{1}{L} \|AB\|_F^2\end{aligned}\quad (6.2)$$

In addition, Cauchy-Schwarz yields $(\sum_k p_k) \left(\sum_k \frac{\alpha_k}{p_k}\right) \geq (\sum_k \sqrt{\alpha_k})^2$, with equality attained if $p_k \propto \sqrt{\alpha_k}$. This implies

$$\mathbb{E} [\|M - AB\|_F^2] \geq \frac{1}{L} \left(\sum_k \|A_{:,k}\|_2 \|B_{k,:}\|_2 \right)^2 - \frac{1}{L} \|AB\|_F^2,$$

where the lower bound is achieved when $p_k \propto \|A_{:,k}\|_2 \|B_{k,:}\|_2$

Error concentration

Practically, one often hopes that the approximation error is absolutely controlled most of the time. In other words, we desire an estimator which is sufficiently close to the truth **with high probability**

For approximate matrix multiplication, two error metrics are of particular interest

- Frobenius norm bound: $\|M - AB\|_F$
- spectral norm bound: $\|M - AB\|$

invoke **matrix concentration inequalities** to control these metrics

Recall: the matrix Bernstein inequality

Theorem 6.1 (Matrix Bernstein)

Let $\{\mathbf{X}_l \in \mathbb{R}^{d_1 \times d_2}\}$ be a sequence of independent zero-mean random matrices. Assume each random matrix satisfies $\|\mathbf{X}_l\| \leq R$. Define $V := \max \left\{ \left\| \mathbb{E} \left[\sum_{l=1}^L \mathbf{X}_l \mathbf{X}_l^\top \right] \right\|, \left\| \mathbb{E} \left[\sum_{l=1}^L \mathbf{X}_l^\top \mathbf{X}_l \right] \right\| \right\}$. Then,

$$\mathbb{P} \left\{ \left\| \sum_{l=1}^L \mathbf{X}_l \right\| \geq \tau \right\} \leq (d_1 + d_2) \exp \left(\frac{-\tau^2/2}{V + R\tau/3} \right)$$

Frobenius norm error of matrix multiplication

Theorem 6.2

Suppose $p_k \geq \frac{\beta \|\mathbf{A}_{:,k}\|_2 \|\mathbf{B}_{k,:}\|_2}{\sum_l \|\mathbf{A}_{:,l}\|_2 \|\mathbf{B}_{l,:}\|_2}$ for some quantity $0 < \beta \leq 1$. If $L \gtrsim \frac{\log n}{\beta}$, then with prob. exceeding $1 - O(n^{-10})$,

$$\|\mathbf{M} - \mathbf{AB}\|_{\text{F}} \lesssim \sqrt{\frac{\log n}{\beta L}} \|\mathbf{A}\|_{\text{F}} \|\mathbf{B}\|_{\text{F}}$$

Proof of Theorem 6.2

Clearly, $\text{vec}(\mathbf{M}) = \sum_{l=1}^L \mathbf{X}_l$, where

$\mathbf{X}_l = \sum_{k=1}^n \frac{1}{Lp_k} \mathbf{A}_{:,k} \otimes \mathbf{B}_{k,:}^\top \mathbb{1}\{i_l = k\}$. These matrices $\{\mathbf{X}_l\}$ obey

$$\|\mathbf{X}_l\|_2 \leq \max_k \frac{1}{Lp_k} \|\mathbf{A}_{:,k}\|_2 \|\mathbf{B}_{k,:}\|_2 \asymp \frac{1}{\beta L} \sum_{k=1}^n \|\mathbf{A}_{:,k}\|_2 \|\mathbf{B}_{k,:}\|_2 =: R$$

$$\mathbb{E} \left[\sum_{l=1}^L \|\mathbf{X}_l\|_2^2 \right] = L \sum_{k=1}^n \mathbb{P}\{i_l = k\} \frac{\|\mathbf{A}_{:,k}\|_2^2 \|\mathbf{B}_{k,:}\|_2^2}{L^2 p_k^2} \leq \underbrace{\frac{(\sum_{k=1}^n \|\mathbf{A}_{k,:}\|_2 \|\mathbf{B}_{k,:}\|_2)^2}{\beta L}}_{=: V}$$

Invoke matrix Bernstein to arrive at

$$\begin{aligned} \|\mathbf{M} - \mathbf{AB}\|_F &= \left\| \sum_{l=1}^L (\mathbf{X}_l - \mathbb{E}[\mathbf{X}_l]) \right\|_2 \lesssim \sqrt{V \log n} + R \log n \\ &\asymp \sqrt{\frac{\log n}{\beta L}} \left(\sum_{k=1}^n \|\mathbf{A}_{k,:}\|_2 \|\mathbf{B}_{k,:}\|_2 \right) \leq \sqrt{\frac{\log n}{\beta L}} \|\mathbf{A}\|_F \|\mathbf{B}\|_F \quad (\text{Cauchy-Schwarz}) \end{aligned}$$

Spectral norm error of matrix multiplication

Theorem 6.3

Suppose $p_k \geq \frac{\beta \|\mathbf{A}_{:,k}\|_2^2}{\|\mathbf{A}\|_F^2}$ for some quantity $0 < \beta \leq 1$, and

$L \gtrsim \frac{\|\mathbf{A}\|_F^2}{\beta \|\mathbf{A}\|^2 \log n}$. Then the estimate \mathbf{M} returned by Algorithm 6.2 obeys

$$\|\mathbf{M} - \mathbf{A}\mathbf{A}^\top\| \lesssim \sqrt{\frac{\log n}{\beta L}} \|\mathbf{A}\|_F \|\mathbf{A}\|$$

with prob. exceeding $1 - O(n^{-10})$

- If $L \gtrsim \underbrace{\frac{\|\mathbf{A}\|_F^2}{\|\mathbf{A}\|^2}}_{\text{stable rank}} \frac{\log n}{\varepsilon^2 \beta}$, then $\|\mathbf{M} - \mathbf{A}\mathbf{A}^\top\| \lesssim \varepsilon \|\mathbf{A}\|^2$

- can be generalized to approximate $\mathbf{A}\mathbf{B}$ (Magen, Zouzias '11)

Proof of Theorem 6.3

Write $\mathbf{M} = \sum_{l=1}^L \mathbf{Z}_l$, where $\mathbf{Z}_l = \sum_{k=1}^n \frac{1}{Lp_k} \mathbf{A}_{:,k} \mathbf{A}_{:,k}^\top \mathbb{1}\{i_l = k\}$. These matrices satisfy

$$\begin{aligned}\|\mathbf{Z}_l\|_2 &\leq \max_k \frac{\|\mathbf{A}_{:,k}\|_2^2}{Lp_k} \leq \frac{1}{\beta L} \|\mathbf{A}\|_F^2 =: R \\ \left\| \mathbb{E} \left[\sum_{l=1}^L \mathbf{Z}_l \mathbf{Z}_l^\top \right] \right\| &= \left\| L \sum_{k=1}^n \mathbb{P}\{i_l = k\} \frac{\|\mathbf{A}_{:,k}\|_2^2}{L^2 p_k^2} \mathbf{A}_{:,k} \mathbf{A}_{:,k}^\top \right\| \\ &= \frac{1}{\beta L} \|\mathbf{A}\|_F^2 \|\mathbf{A} \mathbf{A}^\top\| \\ &\leq \frac{1}{\beta L} \|\mathbf{A}\|_F^2 \|\mathbf{A}\|^2 =: V\end{aligned}$$

Invoke matrix Bernstein to conclude that with high prob.,

$$\begin{aligned}\|\mathbf{M} - \mathbf{A} \mathbf{A}^\top\| &= \left\| \sum_{l=1}^L (\mathbf{Z}_l - \mathbb{E}[\mathbf{Z}_l]) \right\| \lesssim \sqrt{V \log n} + B \log n \\ &\asymp \sqrt{\frac{\log n}{\beta L}} \|\mathbf{A}\|_F \|\mathbf{A}\|\end{aligned}$$

Matrix multiplication with one-sided information

What if we can only use the information about \mathbf{A} ?

For example, suppose $p_k \geq \frac{\beta \|\mathbf{A}_{:,k}\|_2^2}{\|\mathbf{A}\|_F^2}$. In this case, matrix Bernstein does NOT yield **sharp concentration**. But we can still use Markov's inequality to get some useful bound

Matrix multiplication with one-sided information

More precisely, when $p_k \geq \frac{\beta \|\mathbf{A}_{:,k}\|_2^2}{\|\mathbf{A}\|_F^2}$, it follows from (6.2) that

$$\begin{aligned}\mathbb{E} \left[\|\mathbf{M} - \mathbf{AB}\|_F^2 \right] &= \frac{1}{L} \sum_k \frac{1}{p_k} \|\mathbf{A}_{:,k}\|_2^2 \|\mathbf{B}_{k,:}\|_2^2 - \frac{1}{L} \|\mathbf{AB}\|_F^2 \\ &\leq \frac{1}{\beta L} \left(\sum_k \|\mathbf{B}_{k,:}\|_2^2 \right) \|\mathbf{A}\|_F^2 \\ &= \frac{\|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2}{\beta L}\end{aligned}$$

Hence, Markov's inequality yields that with prob. at least $1 - \frac{1}{\log n}$,

$$\|\mathbf{M} - \mathbf{AB}\|_F^2 \leq \frac{\|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2 \log n}{\beta L} \quad (6.3)$$

Least squares approximation

Least squares (LS) problems

Given $\mathbf{A} \in \mathbb{R}^{n \times d}$ ($n \gg d$) and $\mathbf{b} \in \mathbb{R}^n$, find the “best” vector s.t. $\mathbf{Ax} \approx \mathbf{b}$, i.e.

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^d} \quad \|\mathbf{Ax} - \mathbf{b}\|_2$$

If \mathbf{A} has full column rank, then

$$\mathbf{x}_{\text{ls}} = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} = \mathbf{V}_A \boldsymbol{\Sigma}_A^{-1} \mathbf{U}_A^\top \mathbf{b}$$

where $\mathbf{A} = \mathbf{U}_A \boldsymbol{\Sigma}_A \mathbf{V}_A^\top$ is the SVD of \mathbf{A} .

Methods for solving LS problems

Direct methods: computational complexity $O(nd^2)$

- *Cholesky decomposition:* compute upper triangular matrix \mathbf{R} s.t. $\mathbf{A}^\top \mathbf{A} = \mathbf{R}^\top \mathbf{R}$, and solve $\mathbf{R}^\top \mathbf{R} \mathbf{x} = \mathbf{A}^\top \mathbf{b}$
- *QR decomposition:* compute QR decomposition $\mathbf{A} = \mathbf{Q} \mathbf{R}$ (\mathbf{Q} : orthonormal; \mathbf{R} : upper triangular), and solve $\mathbf{R} \mathbf{x} = \mathbf{Q}^\top \mathbf{b}$

Iterative methods: computational complexity $O\left(\frac{\sigma_{\max}(\mathbf{A})}{\sigma_{\min}(\mathbf{A})} \log \frac{1}{\varepsilon}\right)$

- *conjugate gradient ...*

Randomized least squares approximation

Basic idea: generate a sketching / sampling matrix Φ (e.g. via random sampling, random projection), and solve instead

$$\tilde{\mathbf{x}}_{\text{ls}} = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \|\Phi(\mathbf{A}\mathbf{x} - \mathbf{b})\|_2$$

Goal: find Φ s.t.

$$\begin{aligned}\tilde{\mathbf{x}}_{\text{ls}} &\approx \mathbf{x}_{\text{ls}} \\ \|\mathbf{A}\tilde{\mathbf{x}}_{\text{ls}} - \mathbf{b}\|_2 &\approx \|\mathbf{A}\mathbf{x}_{\text{ls}} - \mathbf{b}\|_2\end{aligned}$$

Which sketching matrices enable good approximation?

We will start with two **deterministic** conditions that promise reasonably good approximations (Drineas et al '11)

Which sketching matrices enable good approximation?

Let $\mathbf{A} = \mathbf{U}_A \mathbf{\Sigma}_A \mathbf{V}_A^\top$ be the SVD of \mathbf{A} ...

- **Condition 1 (approximate isometry)**

$$\sigma_{\min}^2(\Phi \mathbf{U}_A) \geq \frac{1}{\sqrt{2}} \quad (6.4)$$

- says that $\Phi \mathbf{U}_A$ is an approximate isometry / rotation
- $1/\sqrt{2}$ can be replaced by other positive constants

Which sketching matrices enable good approximation?

Let $\mathbf{A} = \mathbf{U}_A \mathbf{\Sigma}_A \mathbf{V}_A^\top$ be the SVD of \mathbf{A} ...

- **Condition 2 (approximate orthogonality)**

$$\left\| \mathbf{U}_A^\top \mathbf{\Phi}^\top \mathbf{\Phi} (\mathbf{A} \mathbf{x}_{ls} - \mathbf{b}) \right\|_2^2 \leq \frac{\varepsilon}{2} \|\mathbf{A} \mathbf{x}_{ls} - \mathbf{b}\|_2^2 \quad (6.5)$$

- says that $\mathbf{\Phi} \mathbf{U}_A$ is roughly orthogonal to $\mathbf{\Phi} \underbrace{(\mathbf{A} \mathbf{x}_{ls} - \mathbf{b})}_{= (\mathbf{U}_A \mathbf{U}_A^\top - \mathbf{I}) \mathbf{b}}$
- even though this condition depends on \mathbf{b} , one can find $\mathbf{\Phi}$ satisfying this condition without using any information about \mathbf{b}

Can these conditions be satisfied?

Two extreme examples

1. $\Phi = I$, which satisfies

$$\begin{cases} \sigma_{\min}(\Phi U_A) & = \sigma_{\min}(U_A) = 1 \\ \left\| U_A^\top \Phi^\top \Phi (Ax_{ls} - b) \right\|_2 & = \left\| U_A^\top (I - U_A U_A^\top) b \right\|_2 = 0 \end{cases}$$

- easy to construct; hard to solve the subsampled LS problem

Can these conditions be satisfied?

Two extreme examples

2. $\Phi = U_A^\top$, which satisfies

$$\begin{cases} \sigma_{\min}(\Phi U_A) & = \sigma_{\min}(\mathbf{I}) = 1 \\ \left\| U_A^\top \Phi^\top \Phi (\mathbf{A}x_{\text{ls}} - \mathbf{b}) \right\|_2 & = \left\| U_A^\top (\mathbf{I} - U_A U_A^\top) \mathbf{b} \right\|_2 = 0 \end{cases}$$

- hard to construct (i.e. compute U_A); easy to solve subsampled LS problem

Quality of approximation

We'd like to assess the quality of approximation w.r.t. both fitting error and estimation error

Lemma 6.4

Under Conditions 1-2, the solution $\tilde{\mathbf{x}}_{\text{ls}}$ to the subsampled LS problem obeys

$$(i) \quad \|\mathbf{A}\tilde{\mathbf{x}}_{\text{ls}} - \mathbf{b}\|_2 \leq (1 + \varepsilon) \|\mathbf{A}\mathbf{x}_{\text{ls}} - \mathbf{b}\|_2$$

$$(ii) \quad \|\tilde{\mathbf{x}}_{\text{ls}} - \mathbf{x}_{\text{ls}}\|_2 \leq \frac{\sqrt{\varepsilon}}{\sigma_{\min}(\mathbf{A})} \|\mathbf{A}\mathbf{x}_{\text{ls}} - \mathbf{b}\|_2$$

Proof of Lemma 6.4(i)

The subsampled LS problem can be rewritten as

$$\begin{aligned}\min_{\mathbf{x} \in \mathbb{R}^d} \|\Phi \mathbf{b} - \Phi \mathbf{A} \mathbf{x}\|_2^2 &= \min_{\Delta \in \mathbb{R}^d} \|\Phi \mathbf{b} - \Phi \mathbf{A} (\mathbf{x}_{\text{ls}} + \Delta)\|_2^2 \\ &= \min_{\Delta \in \mathbb{R}^d} \|\Phi (\mathbf{b} - \mathbf{A} \mathbf{x}_{\text{ls}}) - \Phi \mathbf{A} \Delta\|_2^2 \\ &= \min_{\mathbf{z} \in \mathbb{R}^d} \left\| \Phi (\mathbf{b} - \mathbf{A} \mathbf{x}_{\text{ls}}) - \underbrace{\Phi U_A \mathbf{z}}_{= \mathbf{A}(\mathbf{x} - \mathbf{x}_{\text{ls}})} \right\|_2^2.\end{aligned}$$

Therefore, the optimal solution \mathbf{z}_{ls} obeys

$$\mathbf{z}_{\text{ls}} = (\mathbf{U}_A^\top \Phi^\top \Phi \mathbf{U}_A)^{-1} (\mathbf{U}_A^\top \Phi^\top) \Phi (\mathbf{b} - \mathbf{A} \mathbf{x}_{\text{ls}}).$$

Combine Conditions 1-2 to obtain

$$\|\mathbf{z}_{\text{ls}}\|_2^2 \leq \left\| (\mathbf{U}_A^\top \Phi^\top \Phi \mathbf{U}_A)^{-1} \right\|^2 \left\| \mathbf{U}_A^\top \Phi^\top \Phi (\mathbf{b} - \mathbf{A} \mathbf{x}_{\text{ls}}) \right\|_2^2 \leq 2\varepsilon \|\mathbf{b} - \mathbf{A} \mathbf{x}_{\text{ls}}\|_2^2$$

Proof of Lemma 6.4(i) (cont.)

Previous bounds further yield

$$\begin{aligned}\|\mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}_{ls}\|_2^2 &= \left\| \underbrace{\mathbf{b} - \mathbf{A}\mathbf{x}_{ls}}_{\perp U_A} + \underbrace{\mathbf{A}\mathbf{x}_{ls} - \mathbf{A}\tilde{\mathbf{x}}_{ls}}_{\in \text{range}(U_A)} \right\|_2^2 \\ &= \|\mathbf{b} - \mathbf{A}\mathbf{x}_{ls}\|_2^2 + \|\mathbf{A}\mathbf{x}_{ls} - \mathbf{A}\tilde{\mathbf{x}}_{ls}\|_2^2 \\ &= \|\mathbf{b} - \mathbf{A}\mathbf{x}_{ls}\|_2^2 + \|\mathbf{U}_A \mathbf{z}_{ls}\|_2^2 \\ &\leq \|\mathbf{b} - \mathbf{A}\mathbf{x}_{ls}\|_2^2 + \|\mathbf{z}_{ls}\|_2^2 \\ &\leq (1 + 2\varepsilon) \|\mathbf{b} - \mathbf{A}\mathbf{x}_{ls}\|_2^2\end{aligned}$$

Finally, we conclude the proof by recognizing that $\sqrt{1 + 2\varepsilon} \leq 1 + \varepsilon$.

Proof of Lemma 6.4(ii)

From the proof of Lemma 6.4(i), we know $\mathbf{A}\mathbf{x}_{\text{ls}} - \mathbf{A}\tilde{\mathbf{x}}_{\text{ls}} = \mathbf{U}_A\mathbf{z}_{\text{ls}}$ and $\|\mathbf{z}_{\text{ls}}\|_2^2 \leq \varepsilon\|\mathbf{b} - \mathbf{A}\mathbf{x}_{\text{ls}}\|_2^2$. These reveal that

$$\begin{aligned}\|\mathbf{x}_{\text{ls}} - \tilde{\mathbf{x}}_{\text{ls}}\|_2^2 &\leq \frac{\|\mathbf{A}(\mathbf{x}_{\text{ls}} - \tilde{\mathbf{x}}_{\text{ls}})\|_2^2}{\sigma_{\min}^2(\mathbf{A})} \\ &= \frac{\|\mathbf{U}_A\mathbf{z}_{\text{ls}}\|_2^2}{\sigma_{\min}^2(\mathbf{A})} \\ &\leq \frac{\|\mathbf{z}_{\text{ls}}\|_2^2}{\sigma_{\min}^2(\mathbf{A})} \\ &\leq \frac{\varepsilon\|\mathbf{b} - \mathbf{A}\mathbf{x}_{\text{ls}}\|_2^2}{\sigma_{\min}^2(\mathbf{A})}\end{aligned}$$

Quality of approximation (cont.)

By imposing further assumptions on \mathbf{b} , we can connect the error bound with $\|\mathbf{x}_{\text{ls}}\|_2$

Lemma 6.5

Suppose $\|\mathbf{U}_A \mathbf{U}_A^\top \mathbf{b}\|_2 \geq \gamma \|\mathbf{b}\|_2$ for some $0 < \gamma \leq 1$. Under Conditions 1-2, the solution $\tilde{\mathbf{x}}_{\text{ls}}$ to the subsampled LS problem obeys

$$\|\mathbf{x}_{\text{ls}} - \tilde{\mathbf{x}}_{\text{ls}}\|_2 \leq \sqrt{\varepsilon} \kappa(\mathbf{A}) \sqrt{\gamma^{-2} - 1} \|\mathbf{x}_{\text{ls}}\|_2$$

where $\kappa(\mathbf{A})$: condition number of \mathbf{A}

- $\|\mathbf{U}_A \mathbf{U}_A^\top \mathbf{b}\|_2 \geq \gamma \|\mathbf{b}\|_2$ says a nontrivial fraction of the energy of \mathbf{b} lies in $\text{range}(\mathbf{A})$

Proof of Lemma 6.5

Since $\mathbf{b} - \mathbf{Ax}_{\text{ls}} = (\mathbf{I} - \mathbf{U}_A \mathbf{U}_A^\top) \mathbf{b}$, one has

$$\begin{aligned} \|\mathbf{b} - \mathbf{Ax}_{\text{ls}}\|_2^2 &= \|(\mathbf{I} - \mathbf{U}_A \mathbf{U}_A^\top) \mathbf{b}\|_2^2 \\ &= \|\mathbf{b}\|_2^2 - \|\mathbf{U}_A \mathbf{U}_A^\top \mathbf{b}\|_2^2 \\ &\leq (\gamma^{-2} - 1) \|\mathbf{U}_A \mathbf{U}_A^\top \mathbf{b}\|_2^2 && \text{(since } \|\mathbf{U}_A \mathbf{U}_A^\top \mathbf{b}\|_2 \geq \gamma \|\mathbf{b}\|_2 \text{)} \\ &= (\gamma^{-2} - 1) \|\mathbf{Ax}_{\text{ls}}\|_2^2 && \text{(since } \mathbf{Ax}_{\text{ls}} = \mathbf{U}_A \mathbf{U}_A^\top \mathbf{b} \text{)} \\ &\leq (\gamma^{-2} - 1) \sigma_{\max}^2(\mathbf{A}) \|\mathbf{x}_{\text{ls}}\|_2^2 \end{aligned}$$

This combined with Lemma 6.4(ii) concludes the proof.

Connection with approximate matrix multiplication

Condition 1 can be guaranteed if

$$\left\| U_A^\top (\Phi^\top \Phi) U_A - \underbrace{U_A^\top U_A}_{=I} \right\| \leq 1 - \frac{1}{\sqrt{2}}$$

Condition 2 can be guaranteed if

$$\left\| U_A^\top (\Phi^\top \Phi) (\mathbf{A}x_{ls} - \mathbf{b}) - \underbrace{U_A^\top (\mathbf{A}x_{ls} - \mathbf{b})}_{=U_A^\top (\mathbf{I} - U_A U_A^\top) \mathbf{b} = \mathbf{0}} \right\|_2^2 \leq \frac{\epsilon}{2} \underbrace{\|U_A\|_2^2}_{=1} \|\mathbf{A}x_{ls} - \mathbf{b}\|_2^2$$

Both conditions can be viewed as approximate matrix multiplication
(by designing $\Phi\Phi^\top$)

A (slow) random projection strategy

Gaussian sampling: let $\Phi \in \mathbb{R}^{r \times n}$ be composed of i.i.d. Gaussian entries $\mathcal{N}(0, \frac{1}{r})$

- Conditions 1-2 are satisfied with high prob. if $r \gtrsim \frac{d \log d}{\epsilon}$ (exercise)
- implementing Gaussian sketching is expensive (computing ΦA takes time $\Omega(nrd) = \Omega(nd^2 \log d)$)

Another random subsampling strategy

Let's begin with Condition 1 and try Algorithm 6.2 with optimal sampling probabilities ...

Another random subsampling strategy

Leverage scores of \mathbf{A} are defined to be $\|(\mathbf{U}_A)_{:,i}\|_2$ ($1 \leq i \leq n$)

Nonuniform random subsampling: set $\Phi \in \mathbb{R}^{r \times n}$ to be a (weighted) random subsampling matrix s.t.

$$\mathbb{P} \left(\Phi_{i,:} = \frac{1}{\sqrt{r p_k}} e_k^\top \right) = p_k, \quad 1 \leq k \leq n$$

with $p_k \propto \|(\mathbf{U}_A)_{i,:}\|_2^2$

- still slow: needs to compute (exactly) leverage scores

Fast and data-agnostic sampling

Can we design **data-agnostic** sketching matrix Φ (i.e. independent of A, b) that allows **fast** computation while satisfying Conditions 1-2?

Subsampled randomized Hadamard transform (SRHT)

An SRHT matrix $\Phi \in \mathbb{R}^{r \times n}$ is

$$\Phi = RHD$$

- $D \in \mathbb{R}^{n \times n}$: diagonal matrix, whose entries are random $\{\pm 1\}$
- $H \in \mathbb{R}^{n \times n}$: Hadamard matrix (scaled by $1/\sqrt{n}$ so it's orthonormal)
- $R \in \mathbb{R}^{r \times n}$: uniform random subsampling

$$\mathbb{P}\left(\mathbf{R}_{i,:} = \sqrt{\frac{n}{r}} \mathbf{e}_k^\top\right) = \frac{1}{n}, \quad 1 \leq k \leq n$$

Subsampled randomized Hadamard transform

Key idea of SRHT:

- use HD to “uniformize” leverage scores (so that $\{\|(HDU_A)_{i,:}\|_2\}$ are more-or-less identical)
- subsample rank-one components uniformly at random

Uniformization of leverage scores

Lemma 6.6

For any fixed matrix $U \in \mathbb{R}^{n \times d}$, one has

$$\max_{1 \leq i \leq n} \|(\mathbf{H}DU)_{i,:}\|_2 \lesssim \frac{\log n}{\sqrt{n}} \|U\|_F$$

with prob. exceeding $1 - O(n^{-9})$

- HD preconditions U with high prob.; more precisely,

$$\frac{\|(\mathbf{H}DU)_{i,:}\|_2^2}{\sum_{l=1}^n \|(\mathbf{H}DU)_{l,:}\|_2^2} = \frac{\|(\mathbf{H}DU)_{i,:}\|_2^2}{\|U\|_F^2} \lesssim \frac{\log^2 n}{n} \quad (6.6)$$

Proof of Lemma 6.6

For any fixed matrix $U \in \mathbb{R}^{n \times d}$, one has

$$(\mathbf{H}DU)_{i,:} = \sum_{j=1}^n \underbrace{h_{i,j} D_{j,j}}_{\text{random on } \{\pm \frac{1}{\sqrt{n}}\}} U_{j,:},$$

which clearly satisfies $\mathbb{E}[(\mathbf{H}DU)_{i,:}] = \mathbf{0}$. In addition,

$$V := \mathbb{E} \left[\sum_{j=1}^n \|h_{i,j} D_{j,j} U_{j,:}\|_2^2 \right] = \frac{1}{n} \sum_{j=1}^n \|U_{j,:}\|_2^2 = \frac{1}{n} \|U\|_F^2$$

$$B := \max_j \|h_{i,j} D_{j,j} U_{j,:}\|_2 = \frac{1}{\sqrt{n}} \max_j \|U_{j,:}\|_2 \leq \frac{1}{\sqrt{n}} \|U\|_F$$

Invoke matrix Bernstein to demonstrate that with prob. $1 - O(n^{-10})$,

$$\|(\mathbf{H}DU)_{i,:}\|_2 \lesssim \sqrt{V \log n} + B \log n \lesssim \frac{\log n}{\sqrt{n}} \|U\|_F$$

Theoretical guarantees for SRHT

When uniform subsampling is adopted, one has $p_k = 1/n$. In view of Lemma 6.6,

$$p_k \geq \beta \frac{\|(\mathbf{H}\mathbf{D}\mathbf{U}_A)_{i,:}\|_2^2}{\sum_{l=1}^n \|(\mathbf{H}\mathbf{D}\mathbf{U}_A)_{l,:}\|_2^2}$$

with $\beta \asymp \log^{-2} n$. Apply Theorem 6.3 to yield

$$\begin{aligned} \left\| \mathbf{U}_A^\top \Phi^\top \Phi \mathbf{U}_A - \mathbf{I} \right\| &= \left\| \mathbf{U}_A^\top \Phi^\top \Phi \mathbf{U}_A - \mathbf{U}_A^\top \mathbf{U}_A \right\| \\ &= \left\| (\mathbf{U}_A^\top \mathbf{D}^\top \mathbf{H}^\top) \mathbf{R}^\top \mathbf{R} (\mathbf{H}\mathbf{D}\mathbf{U}_A) - (\mathbf{U}_A^\top \mathbf{D}^\top \mathbf{H}^\top) (\mathbf{H}\mathbf{D}\mathbf{U}_A) \right\| \\ &\leq 1/2 \end{aligned}$$

when $r \gtrsim \frac{\|\mathbf{H}\mathbf{D}\mathbf{U}_A\|_F^2 \log n}{\|\mathbf{H}\mathbf{D}\mathbf{U}_A\|_2^2 \beta} \asymp d \log^3 n$. This establishes Condition 1

Theoretical guarantees for SRHT

Similarly, Condition 2 is satisfied with high prob. if $r \gtrsim \frac{d \log^3 n}{\varepsilon}$
(exercise)

Back to least squares approximation

Preceding analysis suggests following algorithm

Algorithm 6.3 Randomized LS approximation (uniform sampling)

- 1: Pick $r \gtrsim \frac{d \log^3 n}{\varepsilon}$, and generate $\mathbf{R} \in \mathbb{R}^{r \times n}$, $\mathbf{H} \in \mathbb{R}^{n \times n}$ and $\mathbf{D} \in \mathbb{R}^{n \times n}$ (as described before)
 - 2: **return** $\tilde{\mathbf{x}} = (\mathbf{RHDA})^\dagger \mathbf{RHDb}$
-

- computational complexity:

$$O\left(\underbrace{nd \log \frac{n}{\varepsilon}}_{\text{compute } \mathbf{HDA}} + \underbrace{\frac{d^3 \log^3 n}{\varepsilon}}_{\text{solve subsampled LS } (rd^2)} \right)$$

An alternative approach: nonuniform sampling

Key idea of Algorithm 6.3 is to uniformize leverage scores followed by uniform sampling

Alternatively, one can also start by estimating leverage scores, and then apply **nonuniform sampling** accordingly

Fast approximation of leverage scores

Key idea: apply SRHT (or other fast Johnson-Lindenstrass transform) in appropriate places

$$\begin{aligned}\|U_{i,:}\|_2^2 &= \|e_i^\top U\|_2^2 = \|e_i^\top U U^\top\|_2^2 \\ &= \|e_i^\top \mathbf{A} \mathbf{A}^\dagger\|_2^2 \\ &\approx \|e_i^\top \mathbf{A} \mathbf{A}^\dagger \Phi_1^\top\|_2^2\end{aligned}$$

where $\Phi_1 \in \mathbb{R}^{r_1 \times n}$ is SRHT matrix

Issue: $\mathbf{A} \mathbf{A}^\dagger$ is expensive to compute; can we compute $\mathbf{A} \mathbf{A}^\dagger \Phi_1^\top$ in a fast manner?

Aside: pseudo inverse

A useful observation: $\mathbf{A}\mathbf{A}^\dagger\Phi^\top \approx \mathbf{A}(\Phi\mathbf{A})^\dagger$, where $\Phi \in \mathbb{R}^{r \times n}$ be SRHT matrix with sufficiently large $r \gg \frac{d \text{poly} \log n}{\varepsilon^2}$

It can be shown that (check Mahoney's lecture notes)

$$\|(\Phi\mathbf{U}_A)^\dagger - (\Phi\mathbf{U}_A)^\top\| \leq \varepsilon$$

$$\text{and } (\Phi\mathbf{A})^\dagger = \mathbf{V}_A\boldsymbol{\Sigma}_A^{-1}(\Phi\mathbf{U}_A)^\dagger$$

These mean

$$\begin{aligned}\mathbf{A}(\Phi\mathbf{A})^\dagger &= \mathbf{U}_A\boldsymbol{\Sigma}_A\mathbf{V}_A^\top\mathbf{V}_A\boldsymbol{\Sigma}_A^{-1}(\Phi\mathbf{U}_A)^\dagger \approx \mathbf{U}_A\boldsymbol{\Sigma}_A\mathbf{V}_A^\top\mathbf{V}_A\boldsymbol{\Sigma}_A^{-1}(\Phi\mathbf{U}_A)^\top \\ &= \mathbf{U}_A\mathbf{U}_A^\top\Phi^\top = \mathbf{A}\mathbf{A}^\dagger\Phi^\top\end{aligned}$$

Fast approximation of leverage scores

Continuing our key idea: apply SRHT (or other fast Johnson-Lindenstrass transform) in appropriate places

$$\begin{aligned}\|U_{i,:}\|_2^2 &\approx \|e_i^\top \mathbf{A}(\Phi_1 \mathbf{A})^\dagger\|_2^2 \\ &\approx \|e_i^\top \mathbf{A}(\Phi_1 \mathbf{A})^\dagger \Phi_2\|_2^2\end{aligned}$$

where $\Phi_1 \in \mathbb{R}^{r_1 \times n}$ and $\Phi_2 \in \mathbb{R}^{r_1 \times r_2}$ ($r_2 \asymp \text{poly log } n$) are both SRHT matrices

Fast approximation of leverage scores

Algorithm 6.4 Leverage scores approximation

- 1: Pick $r_1 \gtrsim \frac{d \log^3 n}{\varepsilon}$ and $r_2 \asymp \text{poly log } n$
 - 2: Compute $\Phi_1 \mathbf{A} \in \mathbb{R}^{r_1 \times d}$ and its QR decomposition, and let $\mathbf{R}_{\Phi_1 \mathbf{A}}$ be the “R” matrix from QR
 - 3: Construct $\Psi = \mathbf{A} \mathbf{R}_{\Phi_1 \mathbf{A}}^{-1} \Phi_2$
 - 4: **return** $\ell_i = \|\Psi_{i,:}\|_2$
-

- computational complexity: $O\left(\frac{nd \text{poly log } n}{\varepsilon^2} + \frac{d^3 \text{poly log } n}{\varepsilon^2}\right)$

Least squares approximation (nonuniform sampling)

Algorithm 6.5 Randomized LS approximation (nonuniform sampling)

- 1: Run Algorithm 6.4 to compute approximate leverage scores $\{\ell_k\}$, and set $p_k \propto \ell_k^2$
 - 2: Randomly sample $r \gtrsim \frac{d \text{poly} \log n}{\varepsilon}$ rows of \mathbf{A} and elements of \mathbf{b} using $\{p_k\}$ as sampling probabilities, rescaling each by $1/\sqrt{rp_k}$. Let $\Phi \mathbf{A}$ and $\Phi \mathbf{b}$ be the subsampled matrix and vector
 - 3: **return** $\tilde{\mathbf{x}}_{\text{ls}} = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \|\Phi \mathbf{A} \mathbf{x} - \Phi \mathbf{b}\|_2$
-

informally, Algorithm 6.5 returns a reasonably good solution with prob. $1 - O(1/\log n)$

Low-rank matrix approximation

Low-rank matrix approximation

Question: given a matrix $A \in \mathbb{R}^{n \times n}$, how to find a rank- k matrix that well approximates A

- One can compute SVD of $A = U\Sigma V^\top$, then return

$$A_k = U_k U_k^\top A$$

where U_k consists of top- k singular vectors

- In general, takes time $O(n^3)$, or $O(kn^2)$ (by power methods)
- Can we find faster algorithms if we only want “good approximation”?

Randomized low-rank matrix approximation

Strategy: find a matrix C (via, e.g., subsampling columns of A), and return

$$\underbrace{CC^\dagger A}_{\text{project } A \text{ onto column space of } C}$$

Question: how well can $CC^\dagger A$ approximate A ?

A simple paradigm

Algorithm 6.6

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

- As we will see, quality of approximation depends on size of

$$\underbrace{AA^T - CC^T}$$

connection with matrix multiplication

Quality of approximation (Frobenius norm)

One can also connect spectral-norm error with product of matrices

Lemma 6.7

The output of Algorithm 6.6 satisfies

$$\|\mathbf{A} - \mathbf{H}_k \mathbf{H}_k^\top \mathbf{A}\|_{\text{F}}^2 \leq \|\mathbf{A} - \mathbf{U}_k \mathbf{U}_k^\top \mathbf{A}\|_{\text{F}}^2 + 2\sqrt{k} \|\mathbf{A} \mathbf{A}^\top - \mathbf{C} \mathbf{C}^\top\|_{\text{F}}$$

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

- This holds for any \mathbf{C}
- Approximation error depends on the error in approximating product of two matrices

Proof of Lemma 6.7

To begin with, since \mathbf{H}_k is orthonormal, one has

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

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

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

Proof of Lemma 6.7

In addition,

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

Finally, one has $\|\mathbf{A} - \mathbf{U}_k \mathbf{U}_k^\top \mathbf{A}\|_{\text{F}}^2 = \|\mathbf{A}\|_{\text{F}}^2 - \sum_{i=1}^k \sigma_i^2(\mathbf{A})$.

Combining above results establishes the claim

Quality of approximation (spectral norm)

Lemma 6.8

The output of Algorithm 6.6 satisfies

$$\|A - H_k H_k^\top A\|^2 \leq \|A - U_k U_k^\top A\|^2 + 2\|AA^\top - CC^\top\|$$

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

Proof of Lemma 6.8

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 $\mathbf{x} \perp \mathbf{H}_k$,

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

This concludes the proof.

Back to low-rank matrix approximation

To ensure $\mathbf{A}\mathbf{A}^\top - \mathbf{C}\mathbf{C}^\top$ is small, we can do random subsampling / projection as before. For example:

Algorithm 6.7

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

Back to low-rank matrix approximation

Invoke Theorems 6.2 and 6.3 to see that with high prob.:

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

$$\|\mathbf{A} - \mathbf{H}_k \mathbf{H}_k^\top \mathbf{A}\|_F^2 \leq \|\mathbf{A} - \mathbf{U}_k \mathbf{U}_k^\top \mathbf{A}\|_F^2 + \varepsilon \|\mathbf{A}\|_F^2 \quad (6.7)$$

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

$$\|\mathbf{A} - \mathbf{H}_k \mathbf{H}_k^\top \mathbf{A}\|^2 \leq \|\mathbf{A} - \mathbf{U}_k \mathbf{U}_k^\top \mathbf{A}\|^2 + \varepsilon \|\mathbf{A}\|^2 \quad (6.8)$$

An improved multi-pass algorithm

Algorithm 6.8 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 $p_k \geq \frac{\beta \|(\mathbf{E}_l)_{:,k}\|_2^2}{\|\mathbf{E}_l\|_{\text{F}}^2}$, $1 \leq k \leq n$
 - 5: Randomly select r column indices with sampling prob. $\{p_k\}$ and append to \mathcal{S}
 - 6: **return** $\mathbf{C} = \mathbf{A}_{\mathcal{S}}$
-

An improved multi-pass algorithm

Theorem 6.9

Suppose $r \gtrsim \frac{k \log n}{\beta \varepsilon^2}$. With high prob.,

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

Proof of Theorem 6.9

We will prove it by induction. Clearly, the claim holds for $t = 1$ (according to (6.7)). Assume

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

and let \mathbf{Z} be the matrix of the columns of \mathbf{E}_t included in the sample. In view of (6.7),

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

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

$$\begin{aligned} \left\| \mathbf{E}_t - \mathbf{Z}\mathbf{Z}^\dagger \mathbf{E}_t \right\|_{\mathbb{F}}^2 &\leq \|\mathbf{E}_t - (\mathbf{E}_t)_k\|_{\mathbb{F}}^2 \\ &+ \frac{\varepsilon}{1-\varepsilon} \|\mathbf{A} - \mathbf{U}_k \mathbf{U}_k^\top \mathbf{A}\|_{\mathbb{F}}^2 + \varepsilon^t \|\mathbf{A}\|_{\mathbb{F}}^2 \end{aligned} \quad (6.9)$$

Proof of Theorem 6.9 (cont.)

If we can show that

$$\mathbf{E}_t - \mathbf{Z}\mathbf{Z}^\dagger \mathbf{E}_t = \mathbf{A} - \mathbf{C}^t(\mathbf{C}^t)^\dagger \mathbf{A} \quad (6.10)$$

$$\|\mathbf{E}_t - (\mathbf{E}_t)_k\|_{\mathbb{F}}^2 \leq \|\mathbf{A} - \mathbf{A}_k\|_{\mathbb{F}}^2 \quad (6.11)$$

then substitution into (6.9) yields

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

We can then use induction to finish proof

Proof of Theorem 6.9 (cont.)

It remains to justify (6.10) and (6.11).

To begin with, (6.10) follows from the definition of \mathbf{E}_t and the fact $\mathbf{Z}\mathbf{Z}^\dagger\mathbf{C}^{t-1}(\mathbf{C}^{t-1})^\dagger = \mathbf{0}$, which gives

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

Proof of Theorem 6.9 (cont.)

To show (6.11), note that $(\mathbf{E}_t)_k$ is best rank- k approximation of \mathbf{E}_t . This gives

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

where \mathbf{A}_k is best rank- k approximation of \mathbf{A} . Substitution into (6.9) establishes the claim for t

Multiplicative error bounds

So far, our results read

$$\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger \mathbf{A}\|_{\mathbf{F}}^2 \leq \|\mathbf{A} - \mathbf{A}_k\|_{\mathbf{F}}^2 + \text{additive error}$$

$$\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger \mathbf{A}\|^2 \leq \|\mathbf{A} - \mathbf{A}_k\|^2 + \text{additive error}$$

In some cases, one might prefer multiplicative error guarantees, e.g.

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

Two types of matrix decompositions

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

$$\hat{\mathbf{A}} = \mathbf{C}\mathbf{X}$$

for some matrix $\mathbf{X} \in \mathbb{R}^{r \times n}$

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

$$\hat{\mathbf{A}} = \mathbf{C}\mathbf{U}\mathbf{R}$$

for some matrix $\mathbf{U} \in \mathbb{R}^{r \times r}$

Generalized least squares problem

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

where \mathbf{X} is matrix (rather than vector)

- generalization of over-determined ℓ_2 regression
- optimal solution: $\mathbf{X}^{\text{ls}} = \mathbf{A}^\dagger \mathbf{B}$
- if $\text{rank}(\mathbf{A}) \leq k$, then $\mathbf{X}^{\text{ls}} = \mathbf{A}_k^\dagger \mathbf{B}$

Generalized least squares approximation

Randomized algorithm: construct an optimally weighted subsampling matrix $\Phi \in \mathbb{R}^{r \times n}$ with $r \gtrsim \frac{k^2}{\epsilon^2}$ and compute

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

Then informally, with high probability,

$$\begin{aligned} \|\mathbf{B} - \mathbf{A}\widetilde{\mathbf{X}}^{\text{ls}}\|_{\text{F}} &\leq (1 + \epsilon) \left\{ \min_{\mathbf{X}} \|\mathbf{B} - \mathbf{A}\mathbf{X}\|_{\text{F}} \right\} \\ \|\mathbf{X}^{\text{ls}} - \widetilde{\mathbf{X}}^{\text{ls}}\|_{\text{F}} &\leq \frac{\epsilon}{\sigma_{\min}(\mathbf{A}_k)} \left\{ \min_{\mathbf{X}} \|\mathbf{B} - \mathbf{A}\mathbf{X}\|_{\text{F}} \right\} \end{aligned}$$

Randomized algorithm for CX decomposition

Algorithm 6.9 Randomized algorithm for constructing CX matrix decompositions

- 1: Compute / approximate sampling probabilities $\{p_i\}_{i=1}^n$, where $p_i = \frac{1}{k} \|(\mathbf{U}_{A,k})_{:,i}\|_2^2$
 - 2: Use sampling probabilities $\{p_i\}$ to construct a rescaled random sampling matrix Φ
 - 3: Construct $\mathbf{C} = \mathbf{A}\Phi^\top$
-

Theoretical guarantees

Theorem 6.10

Suppose $r \gtrsim \frac{k \log k}{\varepsilon^2}$, then Algorithm 6.9 yields

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

Proof of Theorem 6.10

$$\begin{aligned} & \| \mathbf{A} - \underbrace{\mathbf{C}\mathbf{C}^\dagger}_{:= \mathbf{X}^{\text{ls}}} \mathbf{A} \|_{\text{F}} \\ &= \| \mathbf{A} - (\mathbf{A}\Phi^\top)(\mathbf{A}\Phi^\top)^\dagger \mathbf{A} \|_{\text{F}} \\ &\leq \| \mathbf{A} - (\mathbf{A}\Phi^\top)(\mathbf{P}_{A_k}\mathbf{A}\Phi^\top)^\dagger \mathbf{P}_{A_k}\mathbf{A} \|_{\text{F}} \quad (\mathbf{P}_{A_k} := \mathbf{U}_k\mathbf{U}_k^\top) \\ &\quad \text{since } \mathbf{X}^{\text{ls}} := \mathbf{C}^\dagger \mathbf{A} \text{ minimizes } \| \mathbf{A} - \mathbf{C}\mathbf{X} \|_{\text{F}} \\ &= \| \mathbf{A} - (\mathbf{A}\Phi^\top)(\mathbf{A}_k\Phi^\top)^\dagger \mathbf{A}_k \|_{\text{F}} \\ &\leq (1 + \varepsilon) \| \mathbf{A} - \mathbf{A}\mathbf{A}_k^\dagger \mathbf{A}_k \|_{\text{F}} \\ &= (1 + \varepsilon) \| \mathbf{A} - \mathbf{A}_k \|_{\text{F}} \end{aligned}$$

Reference

- "*Lecture notes on randomized linear algebra*," M. Mahoney, 2016.
- "*Randomized algorithms for matrices and data*," M. Mahoney, Foundations and Trends in Machine Learning, 2011.
- "*Low rank matrix-valued Chernoff bounds and approximate matrix multiplication*," A. Magen, Z. Anastasios, SODA, 2011.
- "*Faster least squares approximation*," P. Drineas, M. Mahoney, S. Muthukrishnan, T. Sarlos, Numerische mathematik, 2011.
- "*The fast Johnson-Lindenstrauss transform and approximate nearest neighbors*," N. Ailon, B. Chazelle, SIAM Journal on computing, 2009.
- "*Improved analysis of the subsampled randomized Hadamard transform*," J. Tropp, Advances in Adaptive Data Analysis, 2011

Reference

- "*Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions*," N. Halko, P. G. Martinsson, J. Tropp, SIAM review, 2011.