

## Lecture 25: Randomized Numerical Linear Algebra II

Lecturer: Yudong Chen

Scribe: Zhanpeng Zeng

In the last lecture, we discussed about the application of randomized numerical linear algebra (RandNLA) on approximate matrix multiplication and introduced the problem on least squares (LS) approximation. In this lecture, we continue on discussion on LS problem approximation.<sup>1,2</sup>

## 1 Least squares problem

Given  $\mathbf{A} \in \mathbb{R}^{n \times d}$  ( $n \gg d$ ) and  $\mathbf{b} \in \mathbb{R}^n$ , we would like to find the solution  $\mathbf{x}_{\text{ls}}$  to

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2. \quad (1)$$

Geometrically as shown in Figure 1, this problem is equivalent to projecting  $\mathbf{b}$  to  $\text{range}(\mathbf{A})$  (i.e., the column space of  $\mathbf{A}$ ), so the residual is orthogonal to the column space of  $\mathbf{A}$ :

$$(\mathbf{A}\mathbf{x}_{\text{ls}} - \mathbf{b}) \perp \text{range}(\mathbf{A}). \quad (2)$$

Later, we will see that if we want to approximately compute the least square solution, we want this geometric property to be approximately satisfied.

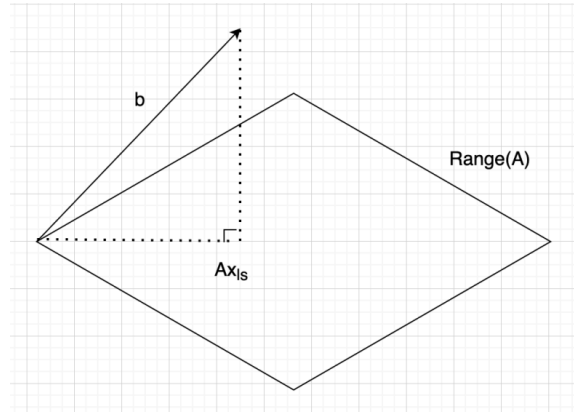


Figure 1: Geometric view of LS solution. Source: Lecture Note 24

Algebraically, the least squares solution  $\mathbf{x}_{\text{ls}}$  satisfies the *normal equation*  $\mathbf{A}^\top \mathbf{A} \mathbf{x}_{\text{ls}} = \mathbf{A}^\top \mathbf{b}$ . If  $\mathbf{A}$  has full column rank, then

$$\mathbf{x}_{\text{ls}} = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} = \underbrace{\mathbf{V}_A \boldsymbol{\Sigma}_A^{-1} \mathbf{U}_A^\top}_{=: \mathbf{A}^\dagger} \mathbf{b}, \quad (3)$$

where  $\mathbf{A} = \mathbf{U}_A \boldsymbol{\Sigma}_A \mathbf{V}_A^\top$  is the SVD of  $\mathbf{A}$  and  $\mathbf{A}^\dagger$  is called the *pseudo inverse* of  $\mathbf{A}$  (as  $\mathbf{A}^\dagger \mathbf{A} = \mathbf{I}_d$ ).

<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.

## 1.1 Methods

In practice, the closed form solution  $(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}$  is rarely used directly to compute  $\mathbf{x}_{\text{ls}}$ . Instead, LS problems are typically solved via one of the following methods:

1. **Direct methods:** these are variants of Gaussian elimination applied to the normal equation, which has computational complexity  $O(nd^2)$ .
  - (a) *Cholesky decomposition:* compute upper triangular matrix  $\mathbf{R}$  s.t.  $\mathbf{A}^\top \mathbf{A} = \mathbf{R}^\top \mathbf{R}$ , and solve  $\mathbf{R}^\top \mathbf{R} \mathbf{x} = \mathbf{A}^\top \mathbf{b}$
  - (b) *QR decomposition:* compute QR decomposition  $\mathbf{A} = \mathbf{Q} \mathbf{R}$  ( $\mathbf{Q}$ : orthonormal;  $\mathbf{R}$ : upper triangular), and solve  $\mathbf{R} \mathbf{x} = \mathbf{Q}^\top \mathbf{b}$

Direct methods yield high accuracy solutions. These methods are often used on dense  $\mathbf{A}$

2. **Iterative methods:** These are *Conjugate Gradient* and variants applied to the optimization problem and have computational complexity (roughly)  $O(nd \cdot \frac{\sigma_{\max}(\mathbf{A})}{\sigma_{\min}(\mathbf{A})} \log \frac{1}{\epsilon})$ . They are often used on large sparse problems. The complexity depends on the condition number of  $\mathbf{A}$ , so they may be slow on ill-conditioned problems.
3. **Randomized methods:** which is our focus. The general goal is to achieve a sample complexity in the form of  $\tilde{O}(nd/\epsilon)$ . In this case, randomized methods are linear in  $d$ , whereas direct methods scale with  $d^2$ , and they are fast on ill-conditioned problem compared to iterative methods.

## 2 Randomized least squares approximation

The **basic idea** is to generate a sketching / sampling matrix  $\Phi \in \mathbb{R}^{r \times n}$  (with  $r \ll n$ , e.g., via random sampling, random projection), and solve instead the smaller, subsampled LS problem

$$\tilde{\mathbf{x}}_{\text{ls}} = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \|\Phi(\mathbf{A} \mathbf{x} - \mathbf{b})\|_2, \quad (4)$$

where the problem size is reduced from  $n \times d$  to  $r \times d$ . Then, if the problem (4) is solved with direct methods, the time complexity is  $O(rd^2) \ll O(nd^2)$ .

Our **goal** is to find  $\Phi$  such that  $\tilde{\mathbf{x}}_{\text{ls}}$  is a good approximation to  $\mathbf{x}_{\text{ls}}$  in terms of both the estimation and fitting error, i.e.,

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

In the following, we will first state two *deterministic* conditions on  $\Phi$  that promise reasonably good approximations (Drineas et al., 2011). Later, we will find examples of  $\Phi$  that satisfy these conditions with high probability.

### 2.1 Deterministic conditions on $\Phi$

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

1. **Condition 1 (approximate isometry)**

$$\sigma_{\min}^2(\Phi \mathbf{U}_A) \geq \frac{1}{\sqrt{2}}. \quad (5)$$

The  $\frac{1}{\sqrt{2}}$  can be replaced by other positive constants. Note that  $\sigma_{\min}^2(\mathbf{U}_A) = 1$ . The condition (5) says that  $\sigma_{\min}^2(\Phi \mathbf{U}_A)$  is still bounded away from 0. Also note that  $\mathbf{U}_A$  is an isometry / rotation in the sense that

$$\|\mathbf{U}_A \mathbf{v}\|_2^2 = \mathbf{v}^\top \mathbf{U}_A^\top \mathbf{U}_A \mathbf{v} = \mathbf{v}^\top \mathbf{v} = \|\mathbf{v}\|_2^2$$

The condition (5) says that  $\Phi U_A$  is an approximate isometry / rotation in the sense that

$$\|\Phi U_A \mathbf{v}\|_2^2 \geq \frac{1}{\sqrt{2}} \|\mathbf{v}\|_2^2$$

One example of  $\Phi$  is a subsampling matrix that picks out a subset of the rows of  $\mathbf{A}$ . The condition (5) says after subsampling  $\mathbf{A}$ , the smallest singular value is approximately preserved.

## 2. Condition 2 (approximate orthogonality)

$$\left\| \mathbf{U}_A^\top \Phi^\top \Phi (\mathbf{A} \mathbf{x}_{\text{ls}} - \mathbf{b}) \right\|_2^2 \leq \frac{\varepsilon}{2} \|\mathbf{A} \mathbf{x}_{\text{ls}} - \mathbf{b}\|_2^2 \quad (6)$$

Recall the geometric interpretation of least squares: the residual of the least-squares solution  $\mathbf{x}_{\text{ls}}$  is orthogonal to  $\text{range}(\mathbf{U}_A) = \text{range}(\mathbf{A})$ :

$$\mathbf{U}_A^\top (\mathbf{A} \mathbf{x}_{\text{ls}} - \mathbf{b}) = 0$$

The condition (6) says that “subsampled residual”  $\Phi (\mathbf{A} \mathbf{x}_{\text{ls}} - \mathbf{b})$  is approximately orthogonal to the “subsampled range”  $\Phi \mathbf{U}_A$ . In other words, we want to  $\Phi$  to preserve the angle between the residual vector and the columns of  $\mathbf{A}$ . Even though this condition depends on  $\mathbf{b}$ , one can find  $\Phi$  satisfying this condition without using any information about  $\mathbf{b}$ , as we show later.

### 2.1.1 Examples

We provide two extreme examples that satisfy the two conditions above.

1.  $\Phi = \mathbf{I}$ , which satisfies

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

This  $\Phi$  is easy to construct, but does not actually reduce the computational complexity of the original LS problem.

2.  $\Phi = \mathbf{U}_A^\top$ , which satisfies

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

This  $\Phi$  is hard to construct: computing  $\mathbf{U}_A$  is at least as hard as solving the original problem. But using this  $\Phi$  leads to a subsampled LS problem that is easy to solve.

Later we will discuss better choices of  $\Phi$  that satisfy the two conditions, are easy to construct, and lead to a easily solvable subsampled LS problem.

## 2.2 Quality of approximation

When these two conditions are satisfied, the following lemma provide guarantee on the quality of approximation w.r.t. both fitting error and estimation error.

**Lemma 1.** *Under Conditions 1-2, the solution  $\tilde{\mathbf{x}}_{\text{ls}}$  to the subsampled LS problem (4) obeys*

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

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

**Proof**

(i) By a change of variable, the subsampled LS problem (4) can be rewritten as

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

where  $\Delta := \mathbf{x} - \mathbf{x}_{\text{ls}}$ . The minimization problem in the last line above is an LS problem over  $\mathbf{z}$ . The optimal solution  $\mathbf{z}_{\text{ls}}$  is given by

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

Recall that Condition 1 ensures that

$$\sigma_{\min}(\mathbf{U}_A^\top \Phi^\top \Phi \mathbf{U}_A) = \sigma_{\min}^2(\Phi \mathbf{U}_A) \geq \frac{1}{2},$$

and Condition 2 ensures that

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

Combining the last three equations gives

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

Previous bounds further yield

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

This completes the proof of part (i) in the lemma.

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

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

as claimed. □

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

**Lemma 2.** *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 (4) 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  $\mathbf{A}$ .

The condition  $\|\mathbf{U}_A \mathbf{U}_A^\top \mathbf{b}\|_2 \geq \gamma \|\mathbf{b}\|_2$  says that the quantity  $\mathbf{U}_A \mathbf{U}_A^\top \mathbf{b}$ , which is the projection of  $\mathbf{b}$  onto the column space of  $\mathbf{A}$ , is large. In other words, a nontrivial fraction of the energy of  $\mathbf{b}$  lies in  $\text{range}(\mathbf{A})$ .

**Proof**

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

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

This combined with Lemma 1(ii) concludes the proof. □

### 3 Randomized algorithms for least squares

Summarizing the last section, we have two deterministic conditions. If  $\Phi$  satisfies these two conditions, we are guaranteed to have a good solution both in terms of fitting error and estimation error.

**Condition 1** can be guaranteed if

$$\left\| \mathbf{U}_A^\top (\Phi^\top \Phi) \mathbf{U}_A - \underbrace{\mathbf{U}_A^\top \mathbf{U}_A}_{=\mathbf{I}} \right\| \leq 1 - \frac{1}{\sqrt{2}}. \quad (7)$$

This condition says that we want to approximate the matrix product  $\mathbf{U}_A^\top \mathbf{U}_A$  by constructing a good  $\Phi$ .

**Condition 2** can be guaranteed if

$$\left\| \underbrace{U_A^\top (\Phi^\top \Phi)}_{=U_A^\top (I - U_A U_A^\top)} (\mathbf{A} \mathbf{x}_{\text{ls}} - \mathbf{b}) - \underbrace{U_A^\top (\mathbf{A} \mathbf{x}_{\text{ls}} - \mathbf{b})}_{b=0} \right\|_2^2 \leq \frac{\varepsilon}{2} \underbrace{\|U_A\|^2}_{=1} \|\mathbf{A} \mathbf{x}_{\text{ls}} - \mathbf{b}\|_2^2. \quad (8)$$

This condition again says that we want to approximate the product of two matrices,  $U_A^\top$  and  $(\mathbf{A} \mathbf{x}_{\text{ls}} - \mathbf{b})$ , by constructing a good  $\Phi$ .

Both conditions can be viewed as approximate matrix multiplication via designing  $\Phi^\top \Phi$ . The key difference compared to approximate matrix multiplication is that the matrices  $U_A$  and  $\mathbf{A} \mathbf{x}_{\text{ls}} - \mathbf{b}$  are not directly given to us.

One data-agnostic choice of  $\Phi$  is **Gaussian sampling**. In particular, let  $\Phi \in \mathbb{R}^{r \times n}$  be composed of i.i.d. Gaussian entries  $\mathcal{N}(0, \frac{1}{r})$ . The expectation of  $\Phi^\top \Phi$  is proportional to the identity matrix. It is then easy to verify that Conditions 1-2 are satisfied with high prob. if  $r \gtrsim \frac{d \log d}{\varepsilon}$  (exercise). The problem is that implementing Gaussian sampling is expensive: computing  $\Phi \mathbf{A}$  takes time  $\Omega(nrd) = \Omega(nd^2 \log d)$ , which is the same as using direct methods to solve the original LS problem.

Can we construct a better  $\Phi$ ? Let us begin with Condition 1 and try the subsampled matrix multiplication approximation with the optimal sampling probabilities. That is, we try to choose a subsampling matrix  $\Phi$  such that

$$U_A^\top U_A \approx U_A^\top \Phi^\top \Phi U_A.$$

From our last lecture on approximate matrix multiplication, we know that if the subsampling probability is proportional to the norm of each row of  $U_A$ , then the approximation is good. These row norms play an important role in least squares problems and beyond, so we give them a name.

**Definition 1.** The *leverage scores* of  $\mathbf{A}$  are defined to be  $\|(U_A)_{k,:}\|_2$  ( $1 \leq k \leq n$ )

The above discussion suggests that one could perform **nonuniform random subsampling** using the leverage score of  $\mathbf{A}$ . In particular, one could set  $\Phi \in \mathbb{R}^{r \times n}$  to be a (weighted) random subsampling matrix such that

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

with  $p_k \propto \|(U_A)_{k,:}\|_2^2$ . This method, however, is still slow: it needs to compute (exactly) leverage scores, which takes  $\Omega(nd^2)$  time, same as direct methods.

Can we design a sketching matrix  $\Phi$  that allows *fast* computation while satisfying Conditions 1-2? Below we discuss two approaches. The first approach constructs a *data-agnostic*  $\Phi$  (i.e., independent of  $\mathbf{A}$ ,  $\mathbf{b}$ ) using Subsampled Randomized Hadamard Transform. The second approach is *data-dependent* and relies on approximate estimation of the leverage scores.

### 3.1 Subsampled Randomized Hadamard Transform (SRHT)

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

$$\Phi = \mathbf{R} \mathbf{H} \mathbf{D}, \quad (9)$$

where the three matrices  $\mathbf{R}$ ,  $\mathbf{H}$  and  $\mathbf{D}$  are given as follows.

- $\mathbf{D} \in \mathbb{R}^{n \times n}$  is a diagonal matrix, whose entries are random  $\{\pm 1\}$ ;
- $\mathbf{H} \in \mathbb{R}^{n \times n}$  is a Hadamard matrix (scaled by  $1/\sqrt{n}$  so it's orthonormal). A Hadamard matrix is a square matrix with  $\pm 1$  entries and orthogonal columns. Hadamard matrices can be constructed recursively by

$$\mathbf{H}_2 = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \quad \mathbf{H}_{2d} = \begin{bmatrix} \mathbf{H}_d & \mathbf{H}_d \\ \mathbf{H}_d & -\mathbf{H}_d \end{bmatrix}, d = 2, 3, \dots$$

- $\mathbf{R} \in \mathbb{R}^{r \times n}$  is uniform random subsampling matrix, where

$$\mathbb{P} \left( \mathbf{R}_{i,:} = \sqrt{\frac{n}{r}} \mathbf{e}_k^\top \right) = \frac{1}{n}, \quad 1 \leq k \leq n. \quad (10)$$

It is easy to verify that  $\mathbb{E} [\Phi^\top \Phi] = \mathbf{I}_n$ .

The **idea of SRHT** is using  $\mathbf{HD}$  to “uniformize” the leverage scores of  $\mathbf{A}$ , so that  $\{\|(\mathbf{HDU}_A)_{i,:}\|_2\}$  are more-or-less identical across  $i$ . Then we can subsample rank-one components *uniformly* at random, since the optimal sampling probability is proportional to the now-uniformized leverage scores. By using SRHT, we can avoid computing the leverage scores of  $\mathbf{U}_A$ .

Moreover, computing the product  $\Phi \mathbf{A} = (\mathbf{RHD})\mathbf{A}$  takes time  $\tilde{\Omega}(nd)$ , since the matrix-vector product  $\mathbf{H}\mathbf{u}$  can be computed in “FFT time”  $\Omega(n \log n)$  thanks to the structure of the Hadamard matrix  $\mathbf{H}$ .

The lemma below rigorously proves that  $\mathbf{HD}$  indeed approximately uniformizes the leverage scores.

**Lemma 3.** *For any fixed matrix  $\mathbf{U} \in \mathbb{R}^{n \times d}$ , one has*

$$\max_{1 \leq i \leq n} \|(\mathbf{HDU})_{i,:}\|_2 \lesssim \frac{\log n}{\sqrt{n}} \|\mathbf{U}\|_{\mathbb{F}}$$

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

In other words, the lemma says that with high probability  $\mathbf{HD}$  preconditions  $\mathbf{U}$ , in the sense that

$$\frac{\|(\mathbf{HDU})_{i,:}\|_2^2}{\sum_{l=1}^n \|(\mathbf{HDU})_{l,:}\|_2^2} = \frac{\|(\mathbf{HDU})_{i,:}\|_2^2}{\|\mathbf{U}\|_{\mathbb{F}}^2} \lesssim \frac{\log^2 n}{n}$$

**Proof**

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

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

It is clear that  $\mathbb{E}[(\mathbf{HDU})_{i,:}] = \mathbf{0}$ . In addition, we have

$$\begin{aligned} V &:= \mathbb{E} \left[ \sum_{j=1}^n \|h_{i,j} D_{j,j} \mathbf{U}_{j,:}\|_2^2 \right] = \frac{1}{n} \sum_{j=1}^n \|\mathbf{U}_{j,:}\|_2^2 = \frac{1}{n} \|\mathbf{U}\|_{\mathbb{F}}^2 \\ B &:= \max_j \|h_{i,j} D_{j,j} \mathbf{U}_{j,:}\|_2 = \frac{1}{\sqrt{n}} \max_j \|\mathbf{U}_{j,:}\|_2 \leq \frac{1}{\sqrt{n}} \|\mathbf{U}\|_{\mathbb{F}} \end{aligned}$$

Invoking the matrix Bernstein inequality, we obtain that with prob.  $1 - O(n^{-10})$ ,

$$\|(\mathbf{HDU})_{i,:}\|_2 \lesssim \sqrt{V \log n} + B \log n \lesssim \frac{\log n}{\sqrt{n}} \|\mathbf{U}\|_{\mathbb{F}}.$$

□

When uniform subsampling is adopted, one has  $p_k = 1/n$ . In view of Lemma 3, we have

$$p_k \geq \beta \frac{\|(\mathbf{HDU}_A)_{i,:}\|_2^2}{\sum_{l=1}^n \|(\mathbf{HDU}_A)_{l,:}\|_2^2} \quad (11)$$

with  $\beta \asymp \log^{-2} n$ . Then, applying Theorem 3 from last lecture (spectral norm error bound for approximate matrix multiplication), we get

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

provided that  $r \gtrsim \frac{\|\mathbf{H} \mathbf{D} \mathbf{U}_A\|_F^2 \log n}{\|\mathbf{H} \mathbf{D} \mathbf{U}_A\|_2^2 \beta} \asymp d \log^3 n$ . This shows that  $\Phi$ , which is given by SRHT, satisfies Condition 1.

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

Putting all together, we have the following randomized algorithm for approximate least squares.

---

**Algorithm 1** 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**  $\tilde{\mathbf{x}} = (\mathbf{R} \mathbf{H} \mathbf{D} \mathbf{A})^\dagger \mathbf{R} \mathbf{H} \mathbf{D} \mathbf{b}$
- 

The computational complexity of the above algorithm is roughly

$$O\left( \underbrace{nd \log \frac{n}{\varepsilon}}_{\text{compute } \mathbf{H} \mathbf{D} \mathbf{A}} + \underbrace{\frac{d^3 \log^3 n}{\varepsilon}}_{\text{solve subsampled LS } (rd^2)} \right). \quad (12)$$

When  $n$  is large, the first term is the dominant term. This archives our goal: linear in  $d$  and independent of the condition number.

### 3.2 Nonuniform sampling

The key idea of Algorithm 1 is to uniformize leverage scores followed by uniform sampling. There is an alternative approach, via nonuniform sampling. In particular, one can start by estimating leverage scores, and then apply *nonuniform sampling* accordingly.

The **key idea** is still applying SRHT (or other fast Johnson-Lindenstrass transform, e.g., Ailon and Chazelle 2009) but in several appropriate places. Observe that leverage scores can be approximated as

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

where  $\Phi_1 \in \mathbb{R}^{r_1 \times n}$  is SRHT matrix. But, the **issue** is that  $\mathbf{A} \mathbf{A}^\dagger$  is expensive to compute. Can we compute  $\mathbf{A} \mathbf{A}^\dagger \Phi_1^\top$  in a fast manner?

Let  $\Phi \in \mathbb{R}^{r \times n}$  be an SRHT matrix with sufficiently large  $r \gg \frac{d \text{poly} \log n}{\varepsilon^2}$ . With high probability, one has (Mahoney, 2016)

$$\|(\Phi \mathbf{U}_A)^\dagger - (\Phi \mathbf{U}_A)^\top\| \leq \varepsilon \quad \text{and} \quad (\Phi \mathbf{A})^\dagger = \mathbf{V}_A \Sigma_A^{-1} (\Phi \mathbf{U}_A)^\dagger.$$

This means that

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

Putting all together, we have the following approximation of the leverage scores:

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



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. Note that the pseudo inverse in the last line above is easy to compute, since  $\Phi_1 \mathbf{A}$  is a much smaller matrix than  $\mathbf{A}$ .

In Algorithm 2 we summarize the above procedure for approximating the leverage scores.

---

**Algorithm 2** Leverage scores approximation

---

- 1: Pick  $r_1 \gtrsim \frac{d \log^3 n}{\epsilon}$  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**  $\ell_i = \|\Psi_{i,:}\|_2$
- 

Once the leverage scores of  $\mathbf{A}$  are approximated, we can do nonuniform sampling with probability  $p_k$  proportional to the approximated leverage scores. This would allow us to satisfy Condition 1.

However, recall that we also need to approximate the product of  $\mathbf{U}_A$  and  $\mathbf{A} \mathbf{x}_{\text{ls}} - \mathbf{b}$  so as to satisfy Condition 2. The optimal sampling probability requires leverage scores for both matrices. We only have information for  $\mathbf{U}_A$  but no information for the residual  $\mathbf{A} \mathbf{x}_{\text{ls}} - \mathbf{b}$  (computing the residual is the same as solving the original LS problem). Fortunately, we can make use of the results for the one sided information setting of approximate matrix multiplication (Section 2.6 from last lecture). In particular, if we choose probability only with information from one matrix,  $\mathbf{U}_A$ , we can still have guarantee on the approximation error, albeit with a worst probability .

Putting all together, we have following randomized LS algorithm, which first estimates the leverage scores using Algorithm 2 and then perform nonuniform sampling.

---

**Algorithm 3** Randomized LS approximation (nonuniform sampling)

---

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

The computational complexity is roughly  $O\left(\frac{nd \text{ poly log } n}{\epsilon^2} + \frac{d^3 \text{ poly log } n}{\epsilon^2}\right)$ . Informally, Algorithm 3 returns a reasonably good solution with probability  $1 - O(1/\log n)$ .

## References

- Ailon, N. and Chazelle, B. (2009). The fast johnson–lindenstrauss transform and approximate nearest neighbors. *SIAM Journal on Computing*, 39(1):302–322.
- Drineas, P., Mahoney, M., Muthukrishnan, S., and Sarlos, T. (2011). Faster least squares approximation. *Numerische Mathematik*.
- Mahoney, M. W. (2016). Lecture notes on randomized linear algebra.