

# Advanced Linear Algebra Week 10 Day 1

2018/11/19 – Jonathan Hayase, updated by Prof. Weiqing Gu

## 1 Application of Advanced Linear Algebra to Big Data

Consider the following:

$$\underbrace{L}_{\text{Linear}} \longleftrightarrow \underbrace{A}_{P^{-1}AP}$$

We are interested in invariants. Some examples of invariants are:  $\det(P^{-1}AP) = \det A$ ,  $\text{tr}(P^{-1}AP) = \text{tr} A$ , and  $\text{rank}(P^{-1}AP) = \text{rank} A$ . If these are the things that we care about then we have the freedom to choose  $P$  to fit our needs, and the quantities we are interested in do not change. This is a general technique that we would like to explore.

Suppose  $A$  (say, a data design matrix) is very big and can't fit into memory. Consider  $\Omega : \mathcal{M}_{n \times n}(\mathbb{K}) \rightarrow \text{BlockD}^n(\mathbb{K})$

$$A = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} \mapsto \begin{bmatrix} \boxed{h_1 \times h_1} & & 0 \\ & \ddots & \\ 0 & & \boxed{h_k \times h_k} \end{bmatrix} \leftarrow B$$

$$A\mathbf{x} = \mathbf{b} \quad A_i\mathbf{x}_i = \mathbf{b}_i \quad i = 1, 2, \dots, k$$

$$B = P^T A P$$

We wish to find  $P \in \text{perm}(n)$  which minimizes  $\|PAP^* - B_D\|$  with large probability.

### 1.1 Norms of Matrices

$A \rightarrow$  view this as an operator

We can talk about the norm that is induced by the vector norm like so:

$$\|A\| \triangleq \sup_{\|\mathbf{x}\|=1} \|A\mathbf{x}\| = \sup_{\mathbf{x} \in \mathbb{R}^n} \frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|}$$

However, there are other norms we might consider, for example:

$$\|A\| = \sqrt{\sum_{i,j} a_{ij}^2} = \sqrt{\text{tr}(A^T A)} = \|A\|_F$$

### 1.2 Johnson-Lindenstrauss lemma

Suppose we have  $\|A\mathbf{x}\| = \|\mathbf{x}\|$ , then  $A$  is orthogonal. Then the mapping from JL, if viewed as a matrix, could be thought of as “almost orthogonal”.

## 1.3 Sketching method as a tool for Linear Algebra

### 1.3.1 Problem 1

Approximating Leverage scores: A  $d$  dimensional subspace  $W$  contained in  $\mathbb{R}^n$  can be expressed written as  $\{x \mid \exists y \in \mathbb{R}^d, x = Uy\}$  for some  $U \in \mathbb{R}^{n \times d}$  with orthonormal columns. The squared Euclidean norms of rows of  $U$  are unique up to permutation. i.e. they depend only on  $A$  and are known as the Leverage scores of  $A$ .

Given  $A$ , we would like to output a list of its leverage scores up to  $1 \pm \epsilon$ .

e.g.  $A = [\mathbf{v}_1 \ \cdots \ \mathbf{v}_n]$ . Consider the column space of  $A$ .  $\implies$  can find an orthonormal basis of the columns space of  $A$ .

### 1.3.2 Problem 2

Least squares regression:

Given  $A \in \mathbb{R}^{n \times d}$  and  $\mathbf{b} \in \mathbb{R}^n$  want to compute  $\|A\tilde{\mathbf{x}} - \mathbf{b}\| \leq (1 + \epsilon) \min_{\mathbf{x} \in \mathbb{R}^d} \|A\mathbf{x} - \mathbf{b}\|$ .

### 1.3.3 Problem 3

Given  $A \in \mathbb{R}^{n \times d}$  and integer  $k > 0$ . Compute  $\tilde{A}_k \in \mathbb{R}^{n \times d}$  with  $\text{rank}(\tilde{A}) \leq k$  so that  $\|A - \tilde{A}_k\|_F \leq (1 + \epsilon) \min_{\text{rank}(A_k) < k} \|A - A_k\|$ .

Today, we will focus on Problem 2.

## 1.4 More on Least Squares Regression

Q: How to find an approximate solution  $\mathbf{x}$  to  $\min_{\mathbf{x}} \|A\mathbf{x} - \mathbf{b}\|$ . Goal: output  $\tilde{\mathbf{x}}$  for which  $\|A\tilde{\mathbf{x}} - \mathbf{b}\|_2 \leq (1 + \epsilon) \min \|A\mathbf{x} - \mathbf{b}\|$ .

Idea: Draw  $S$  from a  $k \times n$  random family of matrices for value  $k \ll n$ : Compute  $SA$  and  $S\mathbf{b}$  and output the solution  $\tilde{\mathbf{x}}$  to  $\min \|(SA)\mathbf{x} - S\mathbf{b}\|$ . e.g.  $S$  is  $d^2/\epsilon \times n$  matrix of i.i.d. normal random variables. E.g.

$$S = [\pm e_{i_1} \ \pm e_{i_2} \ \cdots \ \pm e_{i_k}]$$
$$[e_2 \ -e_1 \ e_1 \ \cdots] = \begin{bmatrix} 0 & -1 & 1 & \cdots & 0 & 0 \\ 1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 0 & \cdots & -1 & 0 \\ 0 & 0 & 0 & \cdots & 0 & 1 \end{bmatrix} \begin{bmatrix} RA_1 \\ \vdots \\ RA_n \end{bmatrix}$$

Original approach/analysis was very long and complicated, then Nelson and Nguyen<sup>1</sup> used Advanced Linear Algebra + Probability to give a much simpler proof.

Key idea: Consider  $[A \mid \mathbf{b}] \triangleq B$  and consider the column space of  $B$ . Let  $U$  be a matrix with columns which form an orthonormal basis of the column space of  $B$ . Claim: It suffices to show  $\|IU\mathbf{x}\|_2 = (1 \pm \epsilon)\|\mathbf{x}\|_2$ . This will imply  $\|S(A\mathbf{x} - \mathbf{b})\|_2 = (1 \pm \epsilon)\|A\mathbf{x} - \mathbf{b}\|_2$  for all  $\mathbf{x}$ .

---

<sup>1</sup>is this right...

## 1.5 Basics of Least Squares

$$y^{(i)} = \theta_0 + \theta_1 x_1^{(i)} + \cdots + \theta_n x_n^{(i)}$$

$$\begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(i)} \\ \vdots \\ y^{(N)} \end{bmatrix} = \underbrace{\begin{bmatrix} \vdots & & & \vdots \\ 1 & x_1^{(i)} & \cdots & x_n^{(i)} \\ \vdots & & & \vdots \end{bmatrix}}_X \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

or

$$\mathbf{y} = X\boldsymbol{\theta} = \underbrace{\theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_m x_m}_{\text{column space of } X}$$

where we wish to find  $\boldsymbol{\theta}$ . Then we may write

$$X^T \mathbf{y} = X^T X \boldsymbol{\theta} \implies \boldsymbol{\theta} = (X^T X)^{-1} (X^T \mathbf{y})$$

## 1.6 An aside on the product of a matrix and its adjoint

$$\mathbf{x}^* A^* A \mathbf{x} = (A\mathbf{x})^* (A\mathbf{x}) = \|A\mathbf{x}\|^2 \geq 0 \implies \lambda_i \geq 0$$

so

$$\|A\|_2 = \sqrt{\lambda_{\max}(A^* A)} = \sigma_{\max}(A)$$

then

$$P^{-1} A P = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$$

$$\text{tr}(P^{-1} A P) = \lambda_1 + \cdots + \lambda_n$$

and

$$\|A\|_F = \sqrt{\text{tr}(A^* A P P^{-1})} = \sqrt{\lambda_1 + \cdots + \lambda_n} = \sqrt{\sum \sigma_i^2} = \left( \sum \sigma_i \right)^{1/2}$$

## 1.7 Why orthonormal bases

Suppose we have an orthonormal basis  $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ . Then any  $\mathbf{x}$  can be written as  $\mathbf{x} = x_1 \mathbf{u}_1 + \cdots + x_n \mathbf{u}_n$  and  $A\mathbf{x} = x_1 A\mathbf{u}_1 + \cdots + x_n A\mathbf{u}_n = x_1 \lambda_1 \mathbf{u}_1 + \cdots + x_n \lambda_n \mathbf{u}_n$ .

$$\begin{aligned} \|A\mathbf{x}\| &= \sqrt{(\lambda_1 x_1)^2 + \cdots + (\lambda_n x_n)^2} \geq \sqrt{(\lambda_n x_1)^2 + \cdots + (\lambda_n x_n)^2} = \sqrt{\lambda_n^2 (x_1^2 + \cdots + x_n^2)} \\ &= \sqrt{\lambda_n^2} \|\mathbf{x}\| = \sqrt{\lambda_n^2} \end{aligned}$$

If  $A$  is hermitian then  $\|A\mathbf{x}\| \geq \lambda_{\min}$ . Similarly,

$$\|A\mathbf{x}\| \leq (\lambda_1 x_1)^2 + \cdots + (\lambda_n x_n)^2 = \lambda_{\max} = (1 \pm \epsilon) \|A\mathbf{x} - b\|_2$$

for all  $\mathbf{x}$ .

## 1.8 Picking up from before...

Claim:  $SU$  is  $\frac{(d+1)^2}{\epsilon^2} \times (d+1)$  matrix.

Claim:  $\|(SU)^T SU - I\|_2 \leq \|U^T S^T SU - I\|_F \leq \epsilon$ .

Definition: An oblivious subspace embedding (OSE) is a distribution  $D$  over matrices  $\Pi \in \mathbb{R}^{n \times m}$  given some parameters  $\epsilon, \delta$  such that for any linear subspace  $W \subseteq \mathbb{R}^n$  with  $\dim W = d$ , the following holds:

$$\mathcal{P}_{\Pi \sim D} (\forall x \in W, \|\Pi x\|_2 \in (1 \pm \epsilon)\|x\|_2) > \frac{2}{3}$$

N N showed that an OSE exists with  $m = O(d^2/\epsilon^2)$  and where  $\Pi \in \text{supp}(\mathcal{O})$  has exactly  $s = 1$  nonzero entries per column. (This improves Woodruff's result.)

Goal: Obtain a fast randomized Algorithm for several numerical linear algebra problems. We focus on Least Squares problem  $\arg\min_{x \in \mathbb{R}^d} \|Ax - b\|$ .

Key idea: Use sketch as application to  $\ell_2$ -estimation in data streams only require  $h$  to be pairwise independent and  $\sigma$  4-wise independent. Claim: A matrix  $\Pi$  preserving the Euclidean norm of all vectors  $x \in W$  up to  $1 \pm \epsilon$  is equivalent to

$$\Pi U y = (\pm \epsilon) \|y\|$$

simultaneously for  $y \in \mathbb{R}^d$  if and only if all singular values lie in the interval  $[1 - \epsilon, 1 + \epsilon]$ .

Claim/Recall: Eigenvalues of orthonormal matrices have norm 1. E.g.

$$A = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

has no eigenvalues geometrically, but with complex eigenvalues  $\lambda_i, |\lambda_i| = 1$ .

Proof: Let  $\lambda, x$  be an eigenvalue/vector pair, so  $Ax = \lambda x$ .

$$\begin{aligned} (Ax)^*(Ax) &= (\overline{Ax})^T (Ax) \\ &= \overline{x}^T \overline{A}^T Ax \\ &= \overline{x}^T x = \|x\|^2 \end{aligned}$$

But

$$(Ax)^*(Ax) = \|Ax\|^2 = \|\lambda x\|^2 = |\lambda|^2 \|x\|^2$$

so  $|\lambda| = 1$ .

Write  $S = (\Pi U)^*(\Pi U)$  (note,  $S$  here is not the same as before) so that we want to show that all of the eigenvalues of  $S$  lie in  $[(1 - \epsilon)^2, (1 + \epsilon)^2]$ .

Trick:  $S = I + (S - I)$ . Use Weyl's inequality (described in the next section).

Let  $M = S$  with eigenvalues  $\mu_i$ ,  $H = I$  with eigenvalues 1, and  $P = S - I$  with eigenvalues  $\rho_i$ .

$$-\|S - I\|_2 \leq \rho_n \leq \mu_i - 1 \leq \rho_1 \leq \|S - I\|_2$$

Because  $S = I + (S - I)$ , we can show all eigenvalues of  $S$  are  $1 \pm \|S - I\|$ . And we want to bound  $\|S - I\|$ . We ultimately are showing that  $S$  is the JL transformation.

### 1.8.1 Weyl's inequality

Let  $M, H, P$  be  $n \times n$  hermitian matrices with only real eigenvalues. where  $M$  has eigenvalues  $\mu_1 \geq \dots \geq \mu_n$ ,  $H$  had eigenvalues  $\gamma_1 \geq \dots \geq \gamma_n$ , and  $P$  has eigenvalues  $\rho_1 \geq \dots \geq \rho_n$ . Then for  $1 \leq i \leq n$ , it holds that  $\rho_n \leq \mu_i - \gamma_i \leq \rho_1$ . Proof [Tao 12].

## 1.9 Continuing from before

We want to show that  $\|S - I\|$  is small with good probability by Markov's inequality.

$$\mathbb{P}(\|S - I\| \geq t) = \mathbb{P}(\|S - I\|^2 \geq t^2) \leq \frac{\|S - I\|^2}{t^2} \leq \frac{\|S - I\|_F}{t^2}$$

and then we proceed to bound.