

# Random linear maps

Daniel Hsu

COMS 4772

1

JL lemma

2

## JL lemma

**Johnson and Lindenstrauss (1984) theorem.** There is a constant  $C > 0$  such that the following holds. For any  $\varepsilon \in (0, 1/2)$ , point set  $S \subset \mathbb{R}^d$  of cardinality  $|S| = n$ , and  $k \in \mathbb{N}$  such that  $k \geq \frac{C \log n}{\varepsilon^2}$ , there exists a linear map  $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$  such that

$$(1-\varepsilon)\|\mathbf{x}-\mathbf{y}\|_2^2 \leq \|f(\mathbf{x})-f(\mathbf{y})\|_2^2 \leq (1+\varepsilon)\|\mathbf{x}-\mathbf{y}\|_2^2 \quad \text{for all } \mathbf{x}, \mathbf{y} \in S.$$

- ▶ There is a randomized procedure to efficiently construct  $f$ .
- ▶ Target dimension  $k$  need not depend on original dimension  $d$ .
- ▶ Any data analysis based on Euclidean distances among  $n$  points can be approximately carried out in dimension  $O(\log n)$ .
  - ▶ E.g., nearest-neighbor computations, many clustering procedures

3

## Proofs of JL lemma

Many ways to (randomly) construct  $f$  that proves the lemma.

1. Original construction:

$$f(\mathbf{x}) = \sqrt{\frac{d}{k}} \mathbf{A} \mathbf{x}$$

where rows of  $\mathbf{A}$  are orthonormal basis (ONB) for  $k$ -dimensional subspace chosen uniformly at random.

2. Simpler construction (Indyk & Motwani, 1998):

$$f(\mathbf{x}) = \frac{1}{\sqrt{k}} \mathbf{A} \mathbf{x}$$

where  $\mathbf{A}$  is a random matrix whose entries are iid  $N(0, 1)$ .

- ▶ Can replace  $N(0, 1)$  with any subgaussian distribution with mean zero and unit variance.

4

## Uniformly random unit vector

Pick  $Z_1, Z_2, \dots, Z_d$  iid  $N(0, 1)$ , and set

$$\mathbf{U} := \frac{(Z_1, Z_2, \dots, Z_d)}{\sqrt{Z_1^2 + Z_2^2 + \dots + Z_d^2}}.$$

Aside: if  $\mathbf{U}$  and  $W_d \sim \chi^2(d)$  are independent, then

$$\sqrt{W_d} \mathbf{U} \sim N(\mathbf{0}, I).$$

5

## ONB for uniformly random $k$ -dimensional subspace

- ▶ Pick  $\mathbf{U}_1$  uniformly at random from  $S^{d-1}$ .
  - ▶ Let columns of  $\mathbf{V}_1$  be ONB for subspace orthogonal to  $\text{span}\{\mathbf{U}_1\}$ .
- ▶ Pick  $\mathbf{U}_2$  uniformly at random from  $\mathbf{V}_1 S^{d-2}$ .
  - ▶ Let columns of  $\mathbf{V}_2$  be ONB for subspace orthogonal to  $\text{span}\{\mathbf{U}_1, \mathbf{U}_2\}$ .
- ▶ Pick  $\mathbf{U}_3$  uniformly at random from  $\mathbf{V}_2 S^{d-3}$ .
  - ▶ Let columns of  $\mathbf{V}_3$  be ONB for subspace orthogonal to  $\text{span}\{\mathbf{U}_1, \mathbf{U}_2, \mathbf{U}_3\}$ .
- ▶ ...
- ▶ Mapping is

$$f(\mathbf{x}) = \sqrt{\frac{d}{k}} \begin{bmatrix} \langle \mathbf{U}_1, \mathbf{x} \rangle \\ \langle \mathbf{U}_2, \mathbf{x} \rangle \\ \vdots \\ \langle \mathbf{U}_k, \mathbf{x} \rangle \end{bmatrix}.$$

6

## ONB for uniformly random $k$ -dimensional subspace

### Easier method:

- ▶ Pick  $k \times d$  random matrix  $\mathbf{A}$  with all entries iid  $N(0, 1)$ .
- ▶ Run *Gram-Schmidt orthogonalization* on the rows.

7

## Requirements of the randomized construction

- ▶  $f$  is a linear map, so  $f(\mathbf{x}) - f(\mathbf{y}) = f(\mathbf{x} - \mathbf{y})$ .
- ▶  $f$  “works” for all  $\binom{n}{2}$  squared lengths  $\|f(\mathbf{x} - \mathbf{y})\|_2^2$ :

$$(1 - \varepsilon)\|\mathbf{x} - \mathbf{y}\|_2^2 \leq \|f(\mathbf{x} - \mathbf{y})\|_2^2 \leq (1 + \varepsilon)\|\mathbf{x} - \mathbf{y}\|_2^2.$$

- ▶ Equivalently, ensure for each of  $\binom{n}{2}$  unit vectors  $\mathbf{v} := \frac{\mathbf{x} - \mathbf{y}}{\|\mathbf{x} - \mathbf{y}\|_2}$ ,

$$1 - \varepsilon \leq \|f(\mathbf{v})\|_2^2 \leq 1 + \varepsilon.$$

- ▶ **Proof strategy:** prove that, for any such unit vector  $\mathbf{v}$ ,

$$\mathbb{P}\left(\|f(\mathbf{v})\|_2^2 \notin [1 - \varepsilon, 1 + \varepsilon]\right) \leq \frac{2}{n^2}.$$

- ▶ By a union bound over all  $\binom{n}{2}$  choices of  $\mathbf{v}$ , we achieve the required properties with probability at least  $1/n$ .

8

## Key lemma

**Key lemma:** for any fixed  $\mathbf{v} \in S^{d-1}$ ,

$$\mathbb{P}\left(\|f(\mathbf{v})\|_2^2 \notin [1 - \varepsilon, 1 + \varepsilon]\right) \leq \frac{2}{n^2}.$$

- ▶ **Simple construction:**  $f(\mathbf{v}) = \frac{1}{\sqrt{k}} \mathbf{A} \mathbf{v}$ , where  $\mathbf{A}$  is  $k \times d$  random matrix with iid  $N(0, 1)$  entries.
- ▶ Each entry of  $\mathbf{A} \mathbf{v}$  is a linear combination of iid  $N(0, 1)$  random variables: for  $Z \sim N(0, 1)$ ,

$$\sum_{j=1}^d A_{i,j} v_j \stackrel{\text{dist}}{=} \left( \sum_{j=1}^d v_j^2 \right)^{1/2} Z = Z.$$

- ▶ So distribution of  $\|\mathbf{A} \mathbf{v}\|_2^2$  is same as that of  $\sum_{i=1}^k Z_i^2$ , where  $Z_1, Z_2, \dots, Z_k$  are iid  $N(0, 1)$ .
- ▶ I.e.,  $Y := \|\mathbf{A} \mathbf{v}\|_2^2 \sim \chi^2(k)$ .

9

## Proof of key lemma

**To prove:** for  $Y \sim \chi^2(k)$ ,

$$\mathbb{P}(Y \notin k[1 - \varepsilon, 1 + \varepsilon]) \leq \frac{2}{n^2}.$$

- ▶ Recall:  $Y$  is  $(4k, 4)$ -subexponential, so

$$\mathbb{P}(Y \geq k + t) \leq \exp\left(-\min\left\{t^2/k, t\right\}/8\right).$$

- ▶ Also can show that  $-Y$  is  $2k$ -subgaussian, so

$$\mathbb{P}(Y \leq k - t) = \mathbb{P}(-Y \geq -k + t) \leq \exp\left(-t^2/(4k)\right).$$

- ▶ For  $t := k\varepsilon$ , each bound is at most  $\exp(-k\varepsilon^2/8)$ .
- ▶ Proof follows by using assumption  $k \geq \frac{16 \ln(n)}{\varepsilon^2}$ . □

10

## Finishing the proof of JL lemma

- ▶ For any pair of distinct points  $\mathbf{x}, \mathbf{y} \in S$ ,

$$\mathbb{P}\left(\frac{\|f(\mathbf{x}) - f(\mathbf{y})\|_2^2}{\|\mathbf{x} - \mathbf{y}\|_2^2} \notin [1 - \varepsilon, 1 + \varepsilon]\right) \leq 2 \exp(-k\varepsilon^2/8) \leq \frac{2}{n^2}.$$

- ▶ Union bound over all  $\binom{n}{2}$  pairs:

$$\mathbb{P}\left(\exists \mathbf{x}, \mathbf{y} \in S. \frac{\|f(\mathbf{x}) - f(\mathbf{y})\|_2^2}{\|\mathbf{x} - \mathbf{y}\|_2^2} \notin [1 - \varepsilon, 1 + \varepsilon]\right) \leq \binom{n}{2} \frac{2}{n^2}.$$

- ▶ Therefore, with probability at least  $1/n$ ,

$$\frac{\|f(\mathbf{x}) - f(\mathbf{y})\|_2^2}{\|\mathbf{x} - \mathbf{y}\|_2^2} \in [1 - \varepsilon, 1 + \varepsilon] \quad \text{for all } \mathbf{x}, \mathbf{y} \in S. \quad \square$$

- ▶ *Note:* success probability is  $1 - \delta$  if  $k \geq \frac{16 \ln(n) + 8 \ln(1/\delta)}{\varepsilon^2}$ .

11

## Original construction

**Original construction:**

$$f(\mathbf{x}) = \sqrt{\frac{d}{k}} \mathbf{A} \mathbf{x}$$

where rows of  $\mathbf{A}$  are ONB for  $k$ -dimensional subspace chosen uniformly at random.

- ▶ Elementary proof by Dasgupta and Gupta (2002) also reduces to similar key lemma: for any *fixed*  $\mathbf{v} \in S^{d-1}$ ,

$$\mathbb{P}\left(\|f(\mathbf{v})\|_2^2 \notin [1 - \varepsilon, 1 + \varepsilon]\right) \leq 2 \exp(-\Omega(k\varepsilon^2)).$$

- ▶ *Key insight:* Distribution of  $\|\mathbf{A}\mathbf{v}\|_2^2$  is the same as  $\|\mathbf{R}\mathbf{U}\|_2^2$ , where  $\mathbf{R}$ 's rows are ONB for *fixed*  $k$ -dimensional subspace, and  $\mathbf{U}$  is a uniformly random unit vector in  $S^{d-1}$ .

12

## Fast JL transform

13

## Computational issues

- ▶  $d$  = original dimension;  $k$  = target dimension.
- ▶ Time to apply  $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$  is  $O(kd)$ .
  - ▶ Due to matrix-vector multiplication.
  - ▶ Not obvious how to speed-up this up because matrix is mostly unstructured.

14

## Using a structured random matrix

- ▶ **Simple idea:** suppose  $\mathbf{M}$  is *sparse*, i.e.,  $\text{nnz}(\mathbf{M}) \ll kd$ .
  - ▶ Can multiply vector by  $\mathbf{M}$  in time  $O(\text{nnz}(\mathbf{M}))$ .
  - ▶ Still want  $\mathbf{M}$  to satisfy “JL property”: for any *fixed*  $\mathbf{x} \in S^{d-1}$ ,

$$\mathbb{P}\left(\|\mathbf{M}\mathbf{x}\|_2^2 \notin [1 - \varepsilon, 1 + \varepsilon]\right) \leq 2 \exp\left(-\Omega(k\varepsilon^2)\right).$$

15

## Sparse random matrix

Define  $\mathbf{M}$  to be  $k \times d$  random matrix with iid entries

$$M_{i,j} := \frac{1}{\sqrt{\theta k}} A_{i,j} B_{i,j},$$

where  $A_{i,j} \sim N(0, 1)$  and  $B_{i,j} \sim \text{Bern}(\theta)$ , which are also independent of each other.

- ▶ Write as  $\mathbf{M} = \frac{1}{\sqrt{\theta k}} (\mathbf{A} \odot \mathbf{B})$ .
- ▶ Scaling ensures  $\mathbb{E} \|\mathbf{M}\mathbf{x}\|_2^2 = 1$  for every  $\mathbf{x} \in S^{d-1}$ .
- ▶  $\mathbb{E}(\text{nnz}(\mathbf{M})) = \theta kd$ .
- ▶ Great if we can use  $\theta = O(1/d + 1/k)$ , which would give  $\mathbb{E}(\text{nnz}(\mathbf{M})) = O(k + d)$ .
- ▶ But does it satisfy JL property?
  - ▶ Depends on  $\mathbf{x} \dots$

16



## JL property for sparse random matrix

$$\|\mathbf{M}\mathbf{x}\|_2^2 = \sum_{i=1}^k \left( \sum_{j=1}^d \frac{1}{\sqrt{\theta k}} A_{i,j} B_{i,j} x_j \right)^2 \stackrel{\text{dist}}{=} \frac{1}{\theta k} \sum_{i=1}^k \left( \sum_{j=1}^d B_{i,j} x_j^2 \right) Z_i^2$$

where  $Z_1, Z_2, \dots, Z_k$  are iid  $N(0, 1)$ .

- Suppose  $\mathbf{x} = (1, 0, \dots, 0)$ .

- $\|\mathbf{M}\mathbf{x}\|_2^2$  depends only on first column of  $\mathbf{M}$ :

$$\|\mathbf{M}\mathbf{x}\|_2^2 \stackrel{\text{dist}}{=} \frac{1}{\theta k} \sum_{i=1}^k B_{i,1} Z_i^2.$$

- Variance is  $\approx 3/(\theta k)$ , which is  $O(\varepsilon^2)$  only if  $\theta = \Omega(1/(k\varepsilon^2))$ .

17

## JL property for sparse random matrix

$$\|\mathbf{M}\mathbf{x}\|_2^2 = \sum_{i=1}^k \left( \sum_{j=1}^d \frac{1}{\sqrt{\theta k}} A_{i,j} B_{i,j} x_j \right)^2 \stackrel{\text{dist}}{=} \frac{1}{\theta k} \sum_{i=1}^k \left( \sum_{j=1}^d B_{i,j} x_j^2 \right) Z_i^2$$

where  $Z_1, Z_2, \dots, Z_k$  are iid  $N(0, 1)$ .

- Suppose instead  $\mathbf{x} = (d^{-1/2}, d^{-1/2}, \dots, d^{-1/2})$ .

- Averaging effect: with high probability,

$$\sum_{j=1}^d B_{i,j} x_j^2 = \frac{1}{d} \sum_{j=1}^d B_{i,j} = \theta \pm O\left(\sqrt{\frac{\theta}{d}} + \frac{1}{d}\right).$$

- Just need  $\theta = \Omega(1/d)$ . In general, just need  $\theta = \Omega(\|\mathbf{x}\|_\infty^2)$ .

18

## Densification

- ▶ Sparse random matrix not great for *sparse unit vectors*, but great for *dense unit vectors*, which have

$$\|\mathbf{x}\|_\infty^2 = \max_{i \in [d]} x_i^2 \approx \frac{1}{d}.$$

- ▶ **Idea:** compose two linear maps.

1. “Densifying” orthogonal transformation:

(maybe sparse)  $\mathbf{x} \mapsto \mathbf{Q}\mathbf{x}$  (likely dense).

2. Sparse linear map:

$$\mathbf{Q}\mathbf{x} \mapsto \frac{1}{\sqrt{\theta k}}(\mathbf{A} \odot \mathbf{B})(\mathbf{Q}\mathbf{x}).$$

19

## Simple densification (picture)

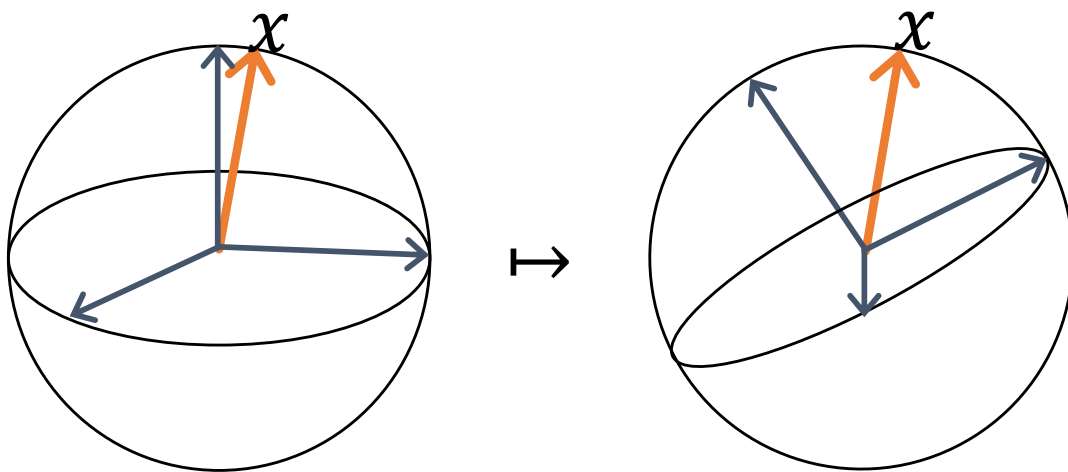


Figure 1: Densifying orthogonal transformation

20

## Simple densification

- ▶ Let  $\mathbf{Q}$  be uniformly random  $d \times d$  orthogonal matrix.
  - ▶  $i$ -th row  $\mathbf{Q}_i^\top$  of  $\mathbf{Q}$  is a uniformly random unit vector.
  - ▶  $i$ -th entry of  $\mathbf{Q}\mathbf{x}$  is  $\langle \mathbf{Q}_i, \mathbf{x} \rangle$ .
- ▶ Can show that

$$\mathbb{P}(|\langle \mathbf{Q}_i, \mathbf{x} \rangle| \geq \varepsilon) \leq 2e^{-\varepsilon^2(d-1)/2}.$$

- ▶ Union bound  $\Rightarrow$  with high probability,

$$\langle \mathbf{Q}_i, \mathbf{x} \rangle^2 \leq O\left(\frac{\log d}{d}\right) \text{ for all } i = 1, 2, \dots, d.$$

21

## Faster densification

- ▶ Unfortunately, uniformly random orthogonal matrix also mostly unstructured; time to apply is  $O(d^2)$ .
- ▶ **Insight of (Ailon and Chazelle, 2006):** can use highly structured “densifying” orthogonal matrix:

$$\mathbf{x} \mapsto \frac{1}{\sqrt{d}} \mathbf{H} \mathbf{D} \mathbf{x}.$$

- ▶  $\mathbf{H} = \mathbf{H}_d$  is the  $d \times d$  Hadamard matrix (not random).
- ▶  $\mathbf{D}$  is random diagonal matrix where diagonal entries are iid Rademacher.

22

## Hadamard matrices

- ▶ Recursive definition (for  $d$  a power of two):

$$\mathbf{H}_1 := +1, \quad \mathbf{H}_d := \begin{bmatrix} +\mathbf{H}_{d/2} & +\mathbf{H}_{d/2} \\ +\mathbf{H}_{d/2} & -\mathbf{H}_{d/2} \end{bmatrix}.$$

- ▶ Example:  $d = 4$

$$\mathbf{H}_4 = \begin{bmatrix} +1 & +1 & +1 & +1 \\ +1 & -1 & +1 & -1 \\ +1 & +1 & -1 & -1 \\ +1 & -1 & -1 & +1 \end{bmatrix}.$$

- ▶ **Fact 1:**  $\frac{1}{\sqrt{d}}\mathbf{H}_d$  is orthogonal, and so is  $\frac{1}{\sqrt{d}}\mathbf{H}_d\mathbf{D}$ .
- ▶ **Fact 2:** Multiplication by  $\mathbf{D}$  requires  $O(d)$  time.
- ▶ **Fact 3:** Multiplication by  $\mathbf{H}_d$  requires  $O(d \log d)$  time!

23

## Hadamard transform via divide-and-conquer

- ▶ To compute  $\mathbf{H}_d\mathbf{x}$ :
  - ▶ Partition  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$ , so  $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^{d/2}$ .
  - ▶ Recursively compute  $\mathbf{H}_{d/2}\mathbf{x}_1$  and  $\mathbf{H}_{d/2}\mathbf{x}_2$ .
  - ▶ Compute  $\mathbf{H}_{d/2}\mathbf{x}_1 + \mathbf{H}_{d/2}\mathbf{x}_2$  and  $\mathbf{H}_{d/2}\mathbf{x}_1 - \mathbf{H}_{d/2}\mathbf{x}_2$ .
  - ▶ Return  $\mathbf{H}_d\mathbf{x} = \begin{bmatrix} \mathbf{H}_{d/2}\mathbf{x}_1 + \mathbf{H}_{d/2}\mathbf{x}_2 \\ \mathbf{H}_{d/2}\mathbf{x}_1 - \mathbf{H}_{d/2}\mathbf{x}_2 \end{bmatrix}$ .
- ▶ Total time:  $O(d \log d)$ .

24

## Analysis of randomized Hadamard transform

- ▶ Let  $\mathbf{Y} := \frac{1}{\sqrt{d}} \mathbf{H} \mathbf{D} \mathbf{x}$  for *fixed* unit vector  $\mathbf{x} \in S^{d-1}$ .
- ▶ Want to show that  $\|\mathbf{Y}\|_\infty^2 = O\left(\frac{\log d}{d}\right)$  with high probability.
- ▶ For each  $i = 1, 2, \dots, d$ ,

$$Y_i = \frac{1}{\sqrt{d}} \sum_{j=1}^d H_{i,j} \sigma_j x_j \stackrel{\text{dist}}{=} \frac{1}{\sqrt{d}} \sum_{j=1}^d x_j \sigma_j,$$

where  $\sigma_1, \sigma_2, \dots, \sigma_d$  are iid Rademacher.

- ▶ Each  $Y_i$  has mean zero and is 1-subgaussian, so with high probability,

$$Y_i^2 \leq O\left(\frac{\log d}{d}\right) \text{ for all } i = 1, 2, \dots, d.$$

25

## Overall random linear map (picture)

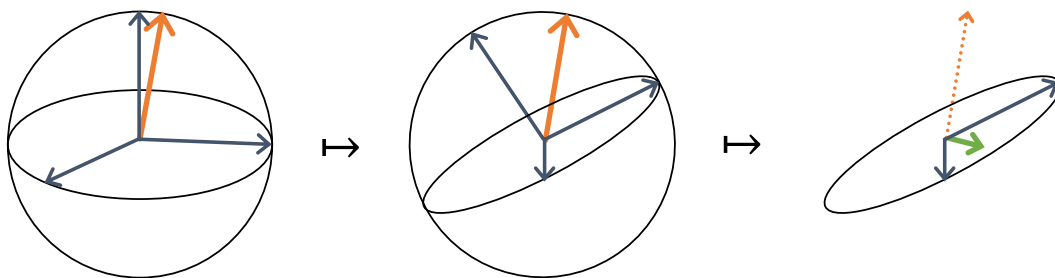


Figure 2: Randomized Hadamard transform + sparse random linear map

26

## Overall random linear map

- ▶ Overall linear map from  $\mathbb{R}^d$  to  $\mathbb{R}^k$ :

1. Densification (randomized Hadamard transform):

$$\mathbf{x} \mapsto \mathbf{y} := \frac{1}{\sqrt{d}} \mathbf{H} \mathbf{D} \mathbf{x}.$$

2. Dimension reduction (sparse random linear map):

$$\mathbf{y} \mapsto \frac{1}{\sqrt{\theta k}} (\mathbf{A} \odot \mathbf{B}) \mathbf{y}.$$

- ▶ Overall running time:  $O(d \log d + \theta k d)$ .
- ▶ Can use  $\theta \approx \frac{\log d}{d}$ , so running time is  $O((d + k) \log d)$ .