

# CS-GY 9223 I: Lecture 12

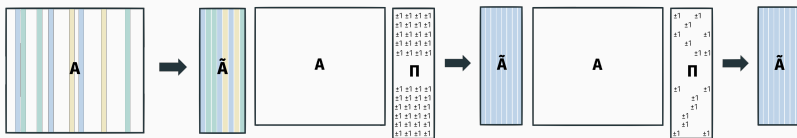
## Randomized numerical linear algebra, fast Johnson-Lindenstrauss Transform

---

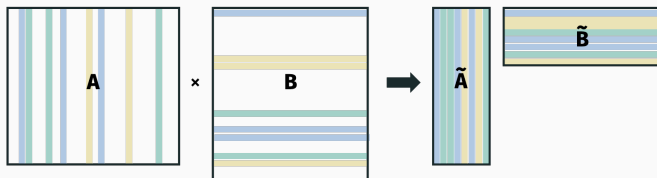
NYU Tandon School of Engineering, Prof. Christopher Musco

**Main idea:** If you want to compute singular vectors or eigenvectors, multiply two matrices, solve a regression problem, etc.:

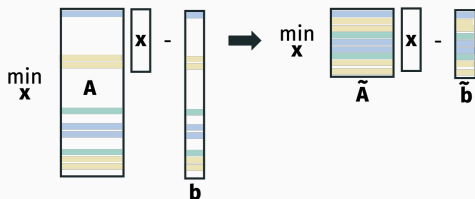
1. Compress your matrices using a randomized method.
2. Solve the problem on the smaller or sparser matrix.
  - $\tilde{A}$  called a “sketch” or “coreset” for  $A$ .



Approximate matrix multiplication:

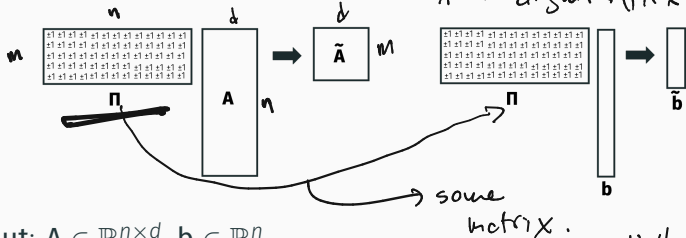


Approximate regression:



# SKETCHED REGRESSION

Randomized approximate regression using a Johnson-Lindenstrauss Matrix:



Input:  $A \in \mathbb{R}^{n \times d}$ ,  $b \in \mathbb{R}^n$ .

Algorithm: Let  $\tilde{x}^* = \arg \min_x \|\tilde{\Pi}Ax - \tilde{\Pi}b\|_2$ .

Goal: Want  $\|\underline{\tilde{A}}\tilde{x}^* - b\|_2^2 \leq \underline{(1 + \epsilon)} \min_x \|\underline{Ax} - b\|_2^2$

If  $\Pi \in \mathbb{R}^{m \times n}$ , how large does  $m$  need to be? Is it even clear this should work as  $m \rightarrow \infty$ ?

## Theorem (Randomized Linear Regression)

Let  $\mathbf{\Pi}$  be a properly scaled JL matrix (random Gaussian, sign, sparse random, etc.) with  $m = O\left(\frac{d \log(d/\epsilon) + \log(1/\delta)}{\epsilon^2}\right)$  rows.

Then with probability  $(1 - \delta)$ , for any  $\mathbf{A} \in \mathbb{R}^{n \times d}$  and  $\mathbf{b} \in \mathbb{R}^n$ ,

$$\|\mathbf{A}\tilde{\mathbf{x}}^* - \mathbf{b}\|_2^2 \leq (1 + \epsilon) \min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$$

where  $\tilde{\mathbf{x}}^* = \arg \min_{\mathbf{x}} \|\mathbf{\Pi}\mathbf{A}\mathbf{x} - \mathbf{\Pi}\mathbf{b}\|_2^2$ .

$$m = O\left(\frac{d}{\epsilon^2}\right)$$

when  $\epsilon = O(1)$ ,

$$m = O(d)$$

# SKETCHED REGRESSION

Claim: Suffices to prove that for all  $x \in \mathbb{R}^d$ ,

$$(1 - \epsilon) \|Ax - b\|_2^2 \leq \|\Pi Ax - \Pi b\|_2^2 \leq (1 + \epsilon) \|Ax - b\|_2^2$$

$$x^* = \underset{x}{\operatorname{argmin}} \|Ax - b\|_2^2.$$

Want to prove:  $\|A\tilde{x}^* - b\|_2^2 \leq (1 + \epsilon) \|Ax^* - b\|_2^2.$

$$\|A\tilde{x}^* - b\|_2^2 \leq \frac{1}{1 - \epsilon} \|\Pi A\tilde{x}^* - \Pi b\|_2^2 \leq \frac{1}{1 - \epsilon} \|\Pi Ax^* - \Pi b\|_2^2$$

by fact that  $\tilde{x}^*$  is opt for  $\Pi A, \Pi b$

For small  $\epsilon$ ,  $\frac{1 + \epsilon}{1 - \epsilon} = 1 + O(\epsilon)$

$$\leq \frac{(1 + \epsilon)}{(1 - \epsilon)} \|Ax^* - b\|_2^2.$$

$\epsilon^1 = \frac{\epsilon}{\text{constant}} \rightarrow$  gives result 6

## Lemma (Distributional JL)

If  $\Pi$  is chosen to a properly scaled random Gaussian matrix, sign matrix, sparse random matrix, etc., with  $O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$  rows then for any fixed  $\mathbf{y}$ ,

$$(1 - \epsilon)\|\mathbf{y}\|_2^2 \leq \|\Pi\mathbf{y}\|_2^2 \leq (1 + \epsilon)\|\mathbf{y}\|_2^2$$

with probability  $(1 - \delta)$ .

**Corollary:** For any fixed  $\mathbf{x}$ , with probability  $(1 - \delta)$ ,

$$\left[ (1 - \epsilon)\|\mathbf{Ax} - \mathbf{b}\|_2^2 \leq \|\Pi\mathbf{Ax} - \Pi\mathbf{b}\|_2^2 \leq (1 + \epsilon)\|\mathbf{Ax} - \mathbf{b}\|_2^2 \right]$$

$$\mathbf{y} = \mathbf{Ax} - \mathbf{b} \quad \downarrow \quad \|\mathbf{y}\|_2^2 \quad \downarrow \quad \|\Pi\mathbf{y}\|_2^2$$

How do we go from "for any fixed  $\mathbf{x}$ " to "for all  $\mathbf{x} \in \mathbb{R}^d$ ".

This statement requires establishing a Johnson-Lindenstrauss type bound for an infinity of possible vectors  $(\mathbf{Ax} - \mathbf{b})$ , which obviously can't be tackled with a union bound argument.

Union bound: prove a statement with  $\delta/\#\text{ events}$   $\rightarrow$  false



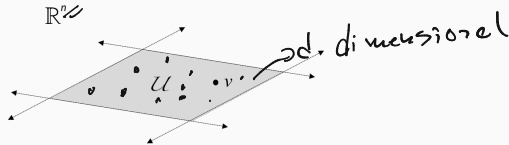
## SUBSPACE EMBEDDINGS

### Theorem (Subspace Embedding from JL)

Let  $\mathcal{U} \subset \mathbb{R}^n$  be a  $d$ -dimensional linear subspace in  $\mathbb{R}^n$ . If  $\Pi \in \mathbb{R}^{m \times d}$  is chosen from any distribution  $\mathcal{D}$  satisfying the Distributional JL Lemma, then with probability  $1 - \delta$ ,

$$(1 - \epsilon)\|\mathbf{v}\|_2^2 \leq \|\Pi\mathbf{v}\|_2^2 \leq (1 + \epsilon)\|\mathbf{v}\|_2^2$$

for all  $\mathbf{v} \in \mathcal{U}$ , as long as  $m = O\left(\frac{d + \log(1/\delta)}{\epsilon^2}\right)$ <sup>1</sup>.  $\approx O\left(\frac{d}{\epsilon^2}\right)$



<sup>1</sup>It's possible to obtain a slightly tighter bound of  $O\left(\frac{d + \log(1/\delta)}{\epsilon^2}\right)$ . It's a nice challenge to try proving this.

## SUBSPACE EMBEDDING TO APPROXIMATE REGRESSION

**Corollary:** If we choose  $\Pi$  and properly scale, then with  $O(d/\epsilon^2)$  rows,

$$\|Ax - b\|_2^2 \leq \|\Pi Ax - \Pi b\|_2^2 \leq (1 + \epsilon) \|Ax - b\|_2^2$$

for all x and thus

$$m = O\left(\frac{d+1}{\epsilon^2}\right)$$

$$\|A\tilde{x}^* - b\|_2^2 \leq (1 + \epsilon) \min_x \|Ax - b\|_2^2.$$

**I.e., our main theorem is proven.**

**Proof:** Apply Subspace Embedding Thm. to the  $(d+1)$  dimensional subspace spanned by  $A$ 's  $d$  columns and  $b$ . Every vector  $Ax - b$  lies in this subspace.

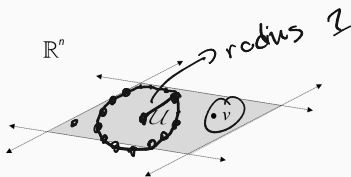
$\hookrightarrow$   $n$  dimension vector, but always in  $d+1$  dimensional subspace of  $\mathbb{R}^n$

## Theorem (Subspace Embedding from JL)

Let  $\mathcal{U} \subset \mathbb{R}^n$  be a  $d$ -dimensional linear subspace in  $\mathbb{R}^n$ . If  $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$  is chosen from any distribution  $\mathcal{D}$  satisfying the Distributional JL Lemma, then with probability  $1 - \delta$ ,

$$(1 - \epsilon)\|\mathbf{v}\|_2^2 \leq \|\mathbf{\Pi}\mathbf{v}\|_2^2 \leq (1 + \epsilon)\|\mathbf{v}\|_2^2 \quad (1)$$

for all  $\mathbf{v} \in \mathcal{U}$ , as long as  $m = O\left(\frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon^2}\right)$



## SUBSPACE EMBEDDING PROOF

**Observation:** The theorem holds as long as (1) holds for all  $\mathbf{w}$  on the unit sphere in  $\mathcal{U}$ . Denote the sphere  $S_{\mathcal{U}}$ :

$$S_{\mathcal{U}} = \{\mathbf{w} \mid \mathbf{w} \in \mathcal{U} \text{ and } \|\mathbf{w}\|_2 = 1\}.$$

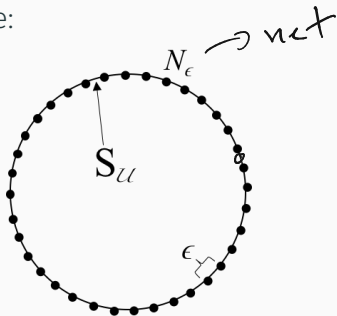
**Follows from linearity:** Any point  $\mathbf{v} \in \mathcal{U}$  can be written as  $c\mathbf{w}$  for some scalar  $c$  and some point  $\mathbf{w} \in S_{\mathcal{U}}$ .

- If  $(1 - \epsilon)\|\mathbf{w}\|_2 \leq \|\Pi\mathbf{w}\|_2 \leq (1 + \epsilon)\|\mathbf{w}\|_2$ .
- then  $c(1 - \epsilon)\|\mathbf{w}\|_2 \leq c\|\Pi\mathbf{w}\|_2 \leq c(1 + \epsilon)\|\mathbf{w}\|_2$ ,
- and thus  $(1 - \epsilon)\|c\mathbf{w}\|_2 \leq \|\Pi c\mathbf{w}\|_2 \leq (1 + \epsilon)\|c\mathbf{w}\|_2$ .

$$\begin{array}{ccc} \downarrow & & \downarrow \\ \|\mathbf{v}\|_2 & & \|\Pi\mathbf{v}\|_2 \end{array}$$

## SUBSPACE EMBEDDING PROOF

**Intuition:** There are not too many “different” points on a  $d$ -dimensional sphere:



$N_\epsilon$  is called an “ $\epsilon$ ”-net.

If we can prove

$$\|w\|_{\infty}(1 - \epsilon) \leq \|\Pi w\|_2 \leq (1 + \epsilon) \|w\|_{\infty}$$

for all points  $w \in N_\epsilon$ , we can hopefully extend to all of  $S_U$ .

## Lemma ( $\epsilon$ -net for the sphere)

For any  $\epsilon \leq 1$ , there exists a set  $N_\epsilon \subset S_{\mathcal{U}}$  with  $|N_\epsilon| = \left(\frac{4}{\epsilon}\right)^d$  such that  $\forall \mathbf{v} \in S_{\mathcal{U}}$ ,

$$\min_{\mathbf{w} \in N_\epsilon} \|\mathbf{v} - \mathbf{w}\| \leq \epsilon.$$

1. Preserving norms of all points in net  $N_\epsilon$ . $\rightarrow 1/|N_\epsilon|$ 

Set  $\delta' = \left(\frac{\epsilon}{4}\right)^d \cdot \delta$ . By a union bound, with probability  $1 - \delta$ , for all  $w \in N_\epsilon$ ,

$$\|w\|_2(1 - \epsilon) \leq \|\Pi w\|_2 \leq (1 + \epsilon)\|w\|_2 \quad \delta$$

 $1 - \delta' \cdot |N_\epsilon|$ 

as long as  $\Pi$  has  $O\left(\frac{\log(1/\delta')}{\epsilon^2}\right) = O\left(\frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon^2}\right)$  rows.

$$\begin{aligned} \log(1/\delta') &= \log\left(\frac{1}{(\epsilon/4)^d \cdot \delta}\right) = \log(1/\delta) + \log\left(\frac{4^d}{\epsilon^d}\right) \\ &= \log(1/\delta) + d \log 4/\epsilon \end{aligned}$$

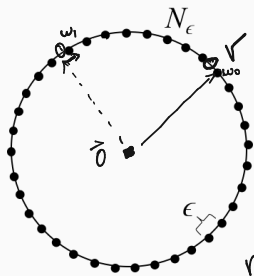
# SUBSPACE EMBEDDING PROOF

## 2. Writing any point in sphere as linear comb. of points in $N_\epsilon$ .

For some  $w_0, w_1, w_2 \dots \in N_\epsilon$ , any  $v \in S_U$  can be written:

$$v = \underbrace{w_0 + c_1 w_1 + c_2 w_2 + \dots}_{\text{remainder}}$$

for constants  $c_1, c_2, \dots$  where  $|c_i| \leq \epsilon^i$



$$v = w_0 + r_0 \rightarrow \text{remainder}$$

where  $\|r_0\|_2 \leq \epsilon$

$$\frac{r_0}{\|r_0\|_2} = w_1 + r_1$$

$$r_0 = \underbrace{\|r_0\|_2}_{\leq \epsilon} w_1 + \underbrace{\|r_0\|_2 r_1}_{\text{norm} \leq \epsilon^2}$$



# SUBSPACE EMBEDDING PROOF

$$(1-\epsilon) \leq \|\Pi w_0\|_2 \leq (1+\epsilon)$$

## 3. Preserving norm of v.

for all  $w_0 \in N_\epsilon$

Applying triangle inequality, we have

$$\begin{aligned} \|\Pi v\|_2 &= \|\Pi w_0 + c_1 \Pi w_1 + c_2 \Pi w_2 + \dots\| \\ &\leq \|\Pi w_0\| + \epsilon \|\Pi w_1\| + \epsilon^2 \|\Pi w_2\| + \dots \\ &\leq (1+\epsilon) + \underbrace{\epsilon(1+\epsilon)}_{\leq 2\epsilon} + \epsilon^2 \underbrace{(1+\epsilon)}_{\leq 2\epsilon^2} + \dots \\ &\leq \underline{1+O(\epsilon)}. \end{aligned}$$

$$\begin{aligned} &\|\Pi w_0 + \Pi r\|_2 \\ &\leq \|\Pi w_0\|_2 \\ &\quad + \underbrace{\|\Pi r\|_2}_{\text{huge? even if } \|r\|_2 \leq \epsilon} \end{aligned}$$

$$\|\Pi v\|_2 \leq (1+O(\epsilon)) \|v\|_2$$

3. Preserving norm of  $v$ .

Similarly,

$$\begin{aligned}
 \|\Pi v\|_2 &= \|\Pi w_0 + c_1 \Pi w_1 + c_2 \Pi w_2 + \dots\| \\
 &\geq \|\Pi w_0\| - \epsilon \|\Pi w_1\| - \epsilon^2 \|\Pi w_2\| - \dots \\
 &\geq (1 - \epsilon) - \epsilon(1 + \epsilon) - \epsilon^2(1 + \epsilon) - \dots \\
 &\geq 1 - O(\epsilon).
 \end{aligned}$$

So we have proven

$$1 - O(\epsilon) \leq \|\Pi \mathbf{v}\|_2 \leq 1 + O(\epsilon)$$

for all  $\mathbf{v} \in S_{\mathcal{U}}$ , which in turn implies for small  $\epsilon$ ,

$$1 - O(\epsilon) \leq \|\Pi \mathbf{v}\|_2^2 \leq 1 + O(\epsilon)$$

Adjusting  $\epsilon$  proves the Subspace Embedding theorem.

**Theorem (Subspace Embedding from JL)**

Let  $\mathcal{U} \subset \mathbb{R}^n$  be a  $d$ -dimensional linear subspace in  $\mathbb{R}^n$ . If  $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$  is chosen from any distribution  $\mathcal{D}$  satisfying the Distributional JL Lemma, then with probability  $1 - \delta$ ,

$$(1 - \epsilon)\|\mathbf{v}\|_2 \leq \|\mathbf{\Pi}\mathbf{v}\|_2 \leq (1 + \epsilon)\|\mathbf{v}\|_2 \quad (2)$$

for all  $\mathbf{v} \in \mathcal{U}$ , as long as  $m = O\left(\frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon^2}\right)$

### Theorem (Randomized Linear Regression)

Let  $\mathbf{\Pi}$  be a properly scaled JL matrix (random Gaussian, sign, sparse random, etc.) with  $m = O\left(\frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon^2}\right)$  rows.

Then with probability  $(1 - \delta)$ , for any  $\mathbf{A} \in \mathbb{R}^{n \times d}$  and  $\mathbf{b} \in \mathbb{R}^n$ ,

$$\|\mathbf{A}\tilde{\mathbf{x}}^* - \mathbf{b}\|_2^2 \leq (1 + \epsilon) \min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$$

where  $\tilde{\mathbf{x}}^* = \arg \min_{\mathbf{x}} \|\mathbf{\Pi A x} - \mathbf{\Pi b}\|_2^2$ .

**Subspace embeddings have many other applications!**

For example, if  $m = O(k/\epsilon)$ ,  $\mathbf{\Pi A}$  can be used to compute an approximate partial SVD, which leads to a  $(1 + \epsilon)$  approximate low-rank approximation for  $\mathbf{A}$ .

## Lemma ( $\epsilon$ -net for the sphere)

For any  $\epsilon \leq 1$ , there exists a set  $N_\epsilon \subset \underline{S_{\mathcal{U}}}$  with  $|N_\epsilon| = \left(\frac{4}{\epsilon}\right)^d$  such that  $\forall \mathbf{v} \in S_{\mathcal{U}}$ ,

$$\min_{\mathbf{w} \in N_\epsilon} \|\mathbf{v} - \mathbf{w}\| \leq \epsilon.$$

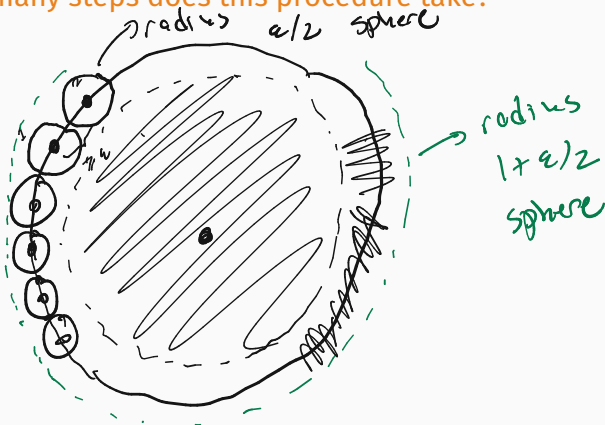
Imaginary algorithm for constructing  $N_\epsilon$ :

- Set  $N_\epsilon = \{\}$
- While such a point exists, choose an arbitrary point  $\mathbf{v} \in S_{\mathcal{U}}$  where  $\nexists \mathbf{w} \in N_\epsilon$  with  $\|\mathbf{v} - \mathbf{w}\| \leq \epsilon$ . Set  $N_\epsilon = N_\epsilon \cup \{\mathbf{w}\}$ .



After running this procedure, we have  $N_\epsilon = \{\mathbf{w}_1, \dots, \mathbf{w}_{|N_\epsilon|}\}$  and  $\min_{\mathbf{w} \in N_\epsilon} \|\mathbf{v} - \mathbf{w}\| \leq \epsilon$  for all  $\mathbf{v} \in S_{\mathcal{U}}$  as desired.

How many steps does this procedure take?



Can place a ball of radius  $\epsilon/2$  around each  $w_i$  without intersecting any other balls. All of these balls live in a ball of radius  $1 + \epsilon/2$ .

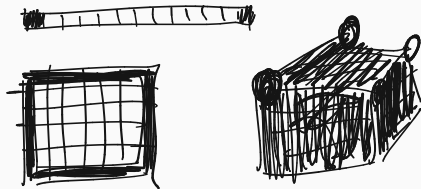
# $\epsilon$ -NET FOR THE SPHERE

Volume of  $d$  dimensional ball of radius  $r$  is

$$\text{vol}(d, r) = c r^d,$$

where  $c$  is a constant that depends on  $d$ , but not  $r$ . From previous slide we have:

$$\begin{aligned} \text{vol}(d, \epsilon/2) \cdot |N_\epsilon| &\leq \text{vol}(d, 1 + \epsilon/2) \\ |N_\epsilon| &\leq \frac{\text{vol}(d, 1 + \epsilon/2)}{\text{vol}(d, \epsilon/2)} = \frac{c \cdot (1 + \epsilon/2)^d}{c \cdot (\epsilon/2)^d} \leq 2 \\ &\leq \left(\frac{4}{\epsilon}\right)^d \leq \left(\frac{4}{\epsilon}\right)^d \end{aligned}$$





## RUNTIME CONSIDERATION

For  $\epsilon, \delta = O(1)$ , we need  $\mathbf{\Pi}$  to have  $m = O(d)$  rows.

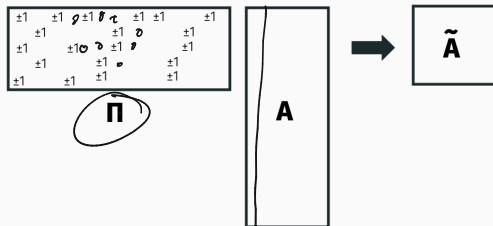
- Cost to solve  $\|\mathbf{Ax} - \mathbf{b}\|_2^2$ :
  - $O(nd^2)$  time for direct method. Need to compute  $(\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{b}$ .  $\mathbf{A}^T\mathbf{A} \rightarrow O(nd^2)$
  - $O(nd) \cdot (\# \text{ of iterations})$  time for iterative method (GD, AGD, conjugate gradient method).  $2\mathbf{A}^T\mathbf{A}\mathbf{x} - 2\mathbf{A}^T\mathbf{b} \rightarrow O(nd)$  time
- Cost to solve  $\|\mathbf{\Pi Ax} - \mathbf{\Pi b}\|_2^2$ :
  - $O(d^3)$  time for direct method.
  - $O(d^2) \cdot (\# \text{ of iterations})$  time for iterative method.

## RUNTIME CONSIDERATION

But time to compute  $\Pi A$  is an  $(m \times n) \times (n \times d)$  matrix multiply:  $O(mnd) = O(nd^2)$  time.

to compute  $\Pi A$ .

Goal: Develop faster Johnson-Lindenstrauss projections.



Typically using sparse and structured matrices.

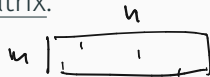
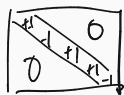
# THE FAST JOHNSON-LINDENSTRAUSS TRANSFORM

Subsampled Randomized Hadamard Transform (SHRT)  
(Ailon-Chazelle, 2006):

Construct  $\mathbf{\Pi} \in \mathbb{R}^{m \times n}$  as follows:

$$\mathbf{\Pi} = \sqrt{\frac{n}{m}} \cdot \mathbf{S} \mathbf{H} \mathbf{D}, \text{ where}$$

- $\mathbf{S} \in \mathbb{R}^{m \times n}$  is a row subsampling matrix. Each row has a single 1 in a random column, all other entries 0.
- $\mathbf{D} \in n \times n$  is a diagonal matrix with each entry uniform  $\pm 1$ .
- $\mathbf{H} \in [n \times n]$  is a Hadamard matrix.



## HADAMARD MATRICES

Assume for now that  $n$  is a power of 2. For  $i = 0, 1, \dots$ ,  $H_i$  is a Hadamard matrix with dimension  $2^i \times 2^i$ .

$$H_0 = 1 \quad H_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

$$H_2 = \frac{1}{\sqrt{4}} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{bmatrix}$$

$$H_k = \frac{1}{\sqrt{2}} \begin{bmatrix} H_{k-1} & H_{k-1} \\ H_{k-1} & -H_{k-1} \end{bmatrix}$$

How long does it take to compute  $Hx$  for a vector  $x \in \mathbb{R}^n$ ?

$\Pi y = \underline{\text{SMD}} y$

# HADAMARD MATRICES

Property 1: Can compute  $\Pi x = \text{SHD}x$  in  $O(n \log n)$  time.

$n \log n$  time  $O(n \log n)$  time  
 $\frac{n}{2}$   
 $Dx = y$

$H_y \rightarrow H_n$  where  $k = \log_2(n)$

$$\underline{H_k} y = \begin{bmatrix} H_{k-1} & H_{k-1} \\ H_{k-1} & -H_{k-1} \end{bmatrix} \begin{bmatrix} y_1 \\ \dots \\ y_r \end{bmatrix} = \begin{bmatrix} H_{k-1} y_1 + H_{k-1} y_2 \\ \dots \\ H_{k-1} y_1 - H_{k-1} y_2 \end{bmatrix}$$

compute:  $H_{k-1} y_1, H_{k-1} y_2$

$$T_k = 2T_{k-1} + n$$

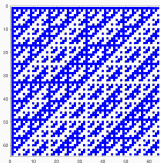
$$T_n = c n \log_2 n$$

Compare to  $O(nm)$  time for random Gaussian or  $\pm 1 \Pi \in \mathbb{R}^{m \times n}$ .

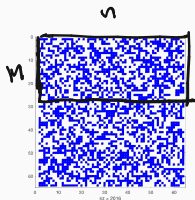
$$T_{k-1} = \frac{n}{2} \log_2(n/2)$$

$$T_k = n \log_2(n/2) + n \rightarrow n \log_2(n)$$

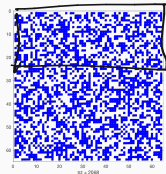
# RANDOMIZED HADAMARD TRANSFORM



Deterministic  
Hadamard matrix.



Randomized  
Hadamard PHD.



Fully random sign  
matrix.

# JOHNSON-LINDENSTRAUSS WITH SHRTS

## Theorem (JL from SRHT)

Let  $\Pi \in \mathbb{R}^{m \times n}$  be a subsampled randomized Hadamard transform with  $m = O\left(\frac{\log(n/\delta)^2 \log(1/\delta)}{\epsilon^2}\right)$  rows. Then for any fixed  $y$ ,

↓  
small cost in embedding dimension

$$\left( (1 - \epsilon) \|y\|_2^2 \leq \|\Pi y\|_2^2 \leq (1 + \epsilon) \|y\|_2^2 \right)$$

with probability  $(1 - \delta)$ .

$$m = O\left(\frac{\log 1/\delta}{\epsilon^2}\right)$$

$$\boxed{\text{Rand Generator}} \Bigg|_m = \boxed{\text{SRHT}} \Bigg|_n$$

# HADAMARD MATRICES ARE ORTHOGONAL

Property 2: For any  $k = 0, 1, \dots$ , we have  $H_k^T H_k = I$ .

$$\begin{aligned} & \frac{1}{\sqrt{2}} \begin{bmatrix} H_{n-1}^T & H_{n-1}^T \\ H_{n-1}^T & -H_{n-1}^T \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} H_{n-1} & H_{n-1} \\ H_{n-1} & -H_{n-1} \end{bmatrix} \\ &= \frac{1}{2} \begin{bmatrix} H_{n-1}^T H_{n-1} + H_{n-1}^T H_{n-1} & H_{n-1}^T H_{n-1} - H_{n-1}^T H_{n-1} \\ H_{n-1}^T H_{n-1} - H_{n-1}^T H_{n-1} & H_{n-1}^T H_{n-1} + H_{n-1}^T H_{n-1} \end{bmatrix} \\ & \quad \downarrow \quad \downarrow \\ & \quad \quad \quad 0 \quad \quad \quad 0 \\ &= \frac{1}{2} \begin{bmatrix} 2I & 0 \\ 0 & 2I \end{bmatrix} = \begin{bmatrix} I \\ I \end{bmatrix} \end{aligned}$$



# RANDOMIZED HADAMARD ANALYSIS

We want to show that  $\|\sqrt{\frac{n}{m}} \mathbf{S} \mathbf{H} \mathbf{D} \mathbf{y}\|_2^2 \approx \|\mathbf{y}\|_2^2$ .

Let  $\underline{z} \in \mathbb{R}^n = \mathbf{H} \mathbf{D} \mathbf{y}$ .

• Claim:  $\|\mathbf{z}\|_2^2 = \|\mathbf{y}\|_2^2$ , exactly.

•  $\|\mathbf{S} \mathbf{H} \mathbf{D} \mathbf{y}\|_2^2 = \frac{n}{m} \|\mathbf{S} \mathbf{z}\|_2^2 =$  subsample of  $\mathbf{z}$ .

•  $\mathbb{E} \left[ \frac{n}{m} \|\mathbf{S} \mathbf{z}\|_2^2 \right] = \|\mathbf{z}\|_2^2$ .

$$\|\mathbf{D} \mathbf{y}\|_2^2 = \mathbf{y}^T \mathbf{D} \mathbf{D} \mathbf{y} = \|\mathbf{y}\|_2^2$$

$$\|\mathbf{H} \mathbf{D} \mathbf{y}\|_2^2$$

$$\mathbf{y}^T \underbrace{\mathbf{D} \mathbf{H}^T \mathbf{H} \mathbf{D}}_{\mathbf{I}} \mathbf{y} = \|\mathbf{D} \mathbf{y}\|_2^2 = \|\mathbf{y}\|_2^2$$

What would  $\mathbf{z}$  have to look like for  $\|\mathbf{S} \mathbf{z}\|_2^2$  to look very different from  $\|\mathbf{z}\|_2^2$  with high probability? I.e. when does subsampling fail. When does subsampling work?



$$\frac{n}{m} \|\mathbf{S} \mathbf{z}\|_2^2 \approx \|\mathbf{z}\|_2^2 = \|\mathbf{D} \mathbf{y}\|_2^2$$

# RANDOMIZED HADAMARD ANALYSIS

## Lemma (SHRT mixing lemma)

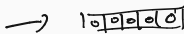
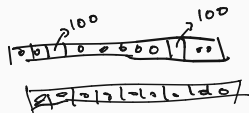
Let  $\underline{H}$  be an  $(n \times n)$  Hadamard matrix and  $\underline{D}$  a random  $\pm 1$  diagonal matrix. Let  $\underline{z} = \underline{H}\underline{D}\mathbf{y}$  for some  $\mathbf{y} \in \mathbb{R}^n$ . With probability  $1 - \delta$ ,

$$|z_i| \leq c \cdot \sqrt{\frac{\log(n/\delta)}{n}} \|\mathbf{y}\|_2$$

for some fixed constant  $c$ .  $\|\mathbf{z}\|_2^2 = \|\delta\|_2^2$

If all entries in  $\mathbf{z}$  were uniform magnitude, we would have

$|z_i| = \frac{1}{\sqrt{n}} \|\mathbf{y}\|_2$ . So we are very close to uniform with high probability.



$$z_i^2 = \frac{\|\mathbf{z}\|_2^2}{n} = \frac{\|\delta\|_2^2}{n}$$

$$\sum z_i^2 = \|\mathbf{z}\|_2^2 = \|\delta\|_2^2$$

$$|z_i| = \frac{1}{\sqrt{n}} \|\delta\|_2$$

# RANDOMIZED HADAMARD ANALYSIS

SHRT mixing lemma proof:  $z = \underbrace{\begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}}_H \begin{bmatrix} D \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$

Let  $\mathbf{h}_i^T$  be the  $i^{\text{th}}$  row of  $\mathbf{H}$ .  $\underline{z}_i = \underline{\mathbf{h}_i^T \mathbf{D} \mathbf{y}}$  where:

$$\mathbf{h}_i^T \mathbf{D} = \frac{1}{\sqrt{n}} \begin{bmatrix} 1 & 1 & -1 & -1 \end{bmatrix} \begin{bmatrix} R_1 & & & \\ & R_2 & & \\ & & R_3 & \\ & & & R_4 \end{bmatrix}$$

where  $R_1, \dots, R_n$  are random  $\pm 1$ 's.  $\rightarrow$  "Rademacher random variable"

This is equivalent to

$$\mathbf{h}_i^T \mathbf{D} = \frac{1}{\sqrt{n}} \begin{bmatrix} R_1 & R_2 & R_3 & R_4 \end{bmatrix} \cdot \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix}$$

# RANDOMIZED HADAMARD ANALYSIS

SHRT mixing lemma proof:

$$\begin{aligned} \mathbb{E}[\sqrt{n} z_i] &= \mathbb{E}[\sum R_i y_i] \\ &= \sum \underbrace{\mathbb{E}[R_i]}_0 y_i \end{aligned}$$

So we have, for all  $i$ ,

$$\underbrace{(\sqrt{n} z_i)} = \underline{h_i^T D y} = \frac{1}{\sqrt{4n}} \sum_{i=1}^n R_i y_i \quad \text{Var}[\sqrt{n} z_i] = \sum \text{Var}[R_i y_i] = \sum y_i^2 \text{Var}[R_i]$$

- $\sqrt{n} \cdot z_i$  is a random variable with mean 0 and variance  $= \sum y_i^2 = \|\mathbf{y}\|_2^2$ , which is a sum of independent random variables.
- By Central Limit Theorem, we expect that:

$$\Pr[|\sqrt{n} \cdot z_i| \geq t \|\mathbf{y}\|_2] \leq e^{-O(t^2)}$$

$$e^{-\sqrt{\log(4/\delta)}^2} = \frac{\delta}{4}$$

- Setting  $t$  gives  $\Pr[|z_i| \geq O\left(\sqrt{\frac{\log(n/\delta)}{n}} \|\mathbf{y}\|_2\right)] \leq \frac{\delta}{n}$
- Applying a union bound to all  $n$  entries of  $\mathbf{z}$  gives the SHRT mixing lemma.

holds for all  $i$  w/ prob.  $1 - \delta$

## RADEMACHER CONCENTRATION

Formally, need to use Bernstein type concentration inequality to prove the bound:

### Lemma (Rademacher Concentration)

Let  $R_1, \dots, R_n$  be Rademacher random variables (i.e. uniform  $\pm 1$ 's). Then for any vector  $\mathbf{a} \in \mathbb{R}^n$ ,

$$\Pr \left[ \sum_{i=1}^n R_i a_i \geq t \|\mathbf{a}\|_2 \right] \leq e^{-t^2/2}.$$

↙ from previous page

## FINISHING UP

With probability  $1 - \delta$ , we have that all  $z_i \leq O\left(\sqrt{\frac{\log(n/\delta)}{n}} \|y\|_2\right)$ .

We want to analyze:  $L \approx \|y\|_2^2 = \|z\|_2^2$

$$L = \left\| \sqrt{\frac{n}{m}} \text{SHD}_y \right\|_2^2 = \frac{1}{m} \left\| \sqrt{n} \text{Sz} \right\|_2^2 = \frac{1}{m} \sum_{i=1}^m (\sqrt{n} z_{j_i})^2$$

where  $j_i$  is a random index in  $1, \dots, n$ .

We have that  $\mathbb{E}L = \|z\|_2^2 = \|y\|_2^2$  and  $L$  is a sum of random variables, each bounded by  $O(\log(n/\delta))$ , which means they have bounded variance.

$$|z_{j_i}| \leq \sqrt{\log(n/\delta)} \cdot \|y\|_2$$

Apply a Chernoff/Hoeffding bound to get that

$|L - \|y\|_2^2| \leq \epsilon \|y\|_2^2$  with probability  $1 - \delta$  as long as:

$\downarrow$   
 $\mathbb{E}[L]$

$$m \geq O\left(\frac{\log^2(n/\delta) \log(1/\delta)}{\epsilon^2}\right).$$



# JOHNSON-LINDENSTRAUSS WITH SHRTS

## Theorem (JL from SRHT)

Let  $\Pi \in \mathbb{R}^{m \times n}$  be a subsampled randomized Hadamard transform with  $m = O\left(\frac{\log(n/\delta) \log(1/\delta)}{\epsilon^2}\right)$  rows. Then for any fixed  $y$ ,

$$(1 - \epsilon)\|y\|_2^2 \leq \|\Pi y\|_2^2 \leq (1 + \epsilon)\|y\|_2^2$$

with probability  $(1 - \delta)$ .

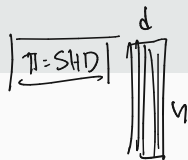
Can be improved to  $m = O\left(\frac{\log(n/\delta) \log(1/\delta)}{\epsilon^2}\right)$ .

**Upshot for regression:** Compute  $\Pi A$  in  $O(nd \log n)$  time instead of  $O(nd^2)$  time. Compress problem down to  $\tilde{A}$  with  $O(d^2)$  dimensions.

$$s = O(\epsilon^d)$$

$$\log(n/\epsilon^d) \cdot \log(1/\epsilon^d)$$

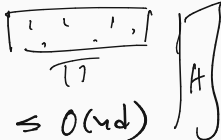
$$\textcircled{d} \log(n/\epsilon) \textcircled{d} \log(1/\epsilon)$$



## BRIEF COMMENT ON OTHER METHODS

$O(nd \log n)$  is nearly linear in the size of  $\mathbf{A}$  when  $\mathbf{A}$  is dense.

Clarkson-Woodruff 2013, STOC Best Paper: Possible to compute  $\tilde{\mathbf{A}}$  with  $\text{poly}(d)$  rows in:

$$\underbrace{O(\text{nnz}(\mathbf{A}))}_{\leq O(nd)} \text{ time.}$$


The diagram shows a matrix  $\mathbf{A}$  with a single row highlighted by a horizontal line. Below the matrix, a vertical brace indicates the number of rows is  $\text{poly}(d)$ . The text  $O(\text{nnz}(\mathbf{A}))$  is underlined, and an arrow points from this underlined term to the expression  $\leq O(nd)$ .

$\Pi$  is chosen to be an ultra-sparse random matrix. Uses totally different techniques (you can't do JL +  $\epsilon$ -net).

Lead to a whole class of matrix algorithms (for regression, SVD, etc.) which run in time:

$$O(\text{nnz}(\mathbf{A})) + \text{poly}(d, \epsilon).$$

$$\text{nnz}(\mathbf{A}) \circ (\# \text{ iterations}) \rightarrow f(k, \epsilon)$$



## WHAT WERE AILON AND CHAZELLE THINKING?



Simple, inspired algorithm that has been used for accelerating:

- Vector dimensionality reduction
- Linear algebra
- Locality sensitive hashing (SimHash)
- Randomized kernel learning methods (we will discuss after Thanksgiving)

```
m = 20|;  
c1 = (2*randi(2,1,n)-3).*y;  
c2 = sqrt(n)*fwht(dy);  
c3 = c2(randperm(n));  
z = sqrt(n/m)*c3(1:m);
```

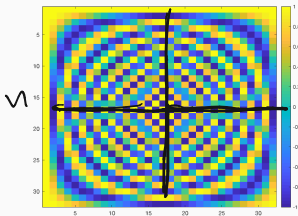
## WHAT WERE AILON AND CHAZELLE THINKING?

The Hadamard Transform is closely related to the Discrete Fourier Transform.

$$F_{j,k} = e^{-2\pi i \frac{j \cdot k}{n}}$$

$$i = \sqrt{-1}$$

$$F^*F = I.$$



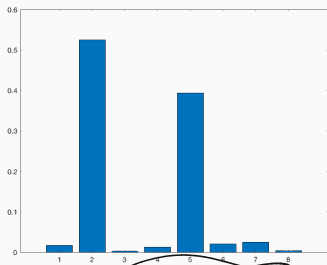
$$\begin{aligned} |e^{-2\pi i \frac{j \cdot k}{n}}| \\ = 1 \end{aligned}$$

Real part of  $F_{j,k}$ .

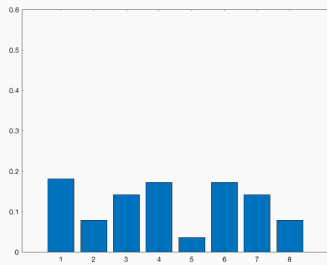
Fy computes the Fourier-transform of the vector y. Can be computed in  $O(n \log n)$  time using a divide and conquer algorithm (the Fast Fourier Transform).

# THE UNCERTAINTY PRINCIPAL

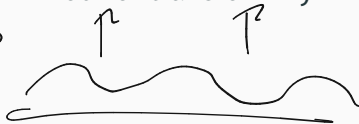
The Uncertainty Principal (informal): A function and its Fourier transform cannot both be concentrated.



Vector  $y$ .



Fourier transform  $Fy$ .



## THE UNCERTAINTY PRINCIPAL

$$S F y$$
$$\rightarrow [S H z] \rightarrow S H D z$$

Sampling does not preserve norms, i.e.  $\|S y\|_2 \neq \|y\|_2$  when  $y$  has a few large entries.

Taking a Fourier transform exactly eliminates this hard case, without changing  $y$ 's norm.