

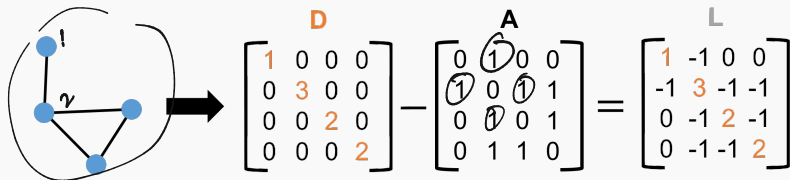
CS-GY 6763: Lecture 11

Randomized numerical linear algebra, ϵ -net arguments.

NYU Tandon School of Engineering, Prof. Christopher Musco

LAST CLASS

Represent undirected graph as symmetric matrix: $n \times n$ adjacency matrix A and graph Laplacian $L = D - A$ where D is the diagonal degree matrix.

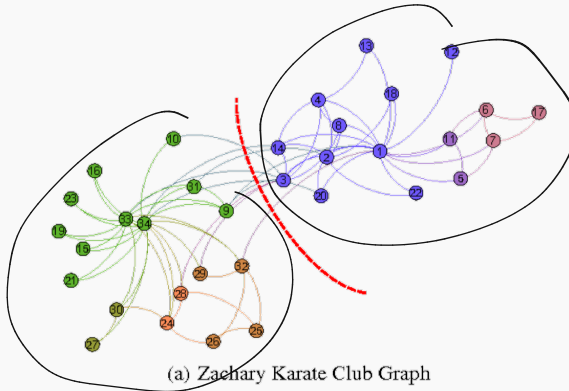


$L = B^T B$ where B is the “edge-vertex incidence” matrix.

of edges $\leftarrow B = \begin{bmatrix} 1 & -1 & 0 & 0 \\ 0 & 1 & -1 & 0 \\ 0 & 1 & 0 & -1 \\ 0 & 0 & 1 & -1 \end{bmatrix}$

Balanced Cut: Partition nodes along a cut that:

- Has few crossing edges: $|\{(u, v) \in E : u \in B, v \in C\}|$ is small.
- Separates large partitions: $|B|, |C|$ are not too small.



RELAX AND ROUND

$$\underline{B}, \underline{C} \quad B \cup C = \{1, \dots, n\}$$

We observed that $\mathbf{x}^T \mathbf{L} \mathbf{x} = \sum_{(i,j) \in E} (\mathbf{x}(i) - \mathbf{x}(j))^2$. If \mathbf{c} is a “cut indicator vector” for a cut between node set B and C – i.e. $\mathbf{c}[i] = 1$ for all $i \in B$ and -1 elsewhere, then it followed that:

$$\mathbf{c}^T \mathbf{L} \mathbf{c} = 4 \cdot \text{cut}(B, C) \quad | \quad \underline{-1 \quad -1 \quad -1 \quad -1}$$

We used this basic fact to argue heuristically that the smallest eigenvectors of \mathbf{L} can be used to find balanced cuts in a graph.

Note: \mathbf{c} often denote by $\chi_{B,C}$.

“Relax and round” algorithm:

- Relax problem $\min \mathbf{c}^T \mathbf{L} \mathbf{c}$ by not requiring \mathbf{c} to be a binary cut-indicator vector.
- Showed that second smallest eigenvector \mathbf{v}_{n-1} of \mathbf{L} solved the relaxed problem.
- Round this vector to be a cut indicator vector: all negative entries rounded to -1 , all positive entries rounded to 1 .

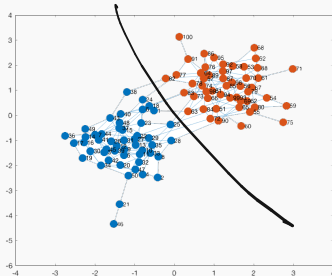
Main theoretical result: This approach is hard to analyze in general, but can be proven to work well on random graphs drawn from the stochastic block model!.

STOCHASTIC BLOCK MODEL

Stochastic Block Model (Planted Partition Model):

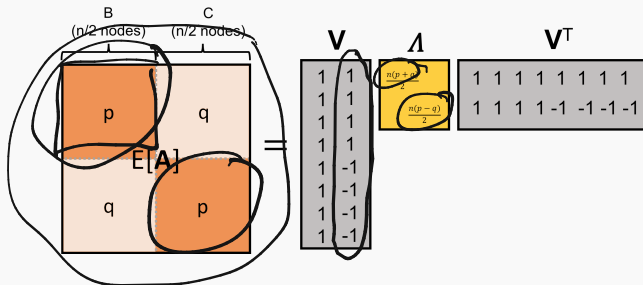
Let $G_n(p, q)$ be a distribution over graphs on n nodes, split equally into two groups B and C , each with $n/2$ nodes.

- Any two nodes in the **same group** are connected with probability p (including self-loops).
- Any two nodes in **different groups** are connected with prob. $q < p$.



EXPECTED ADJACENCY SPECTRUM

$\mathbb{E}[A] = p \cdot I - \mathbb{E}[L]$, so smallest eigenvectors of $\mathbb{E}[L]$ are equal to largest of $\mathbb{E}[A]$.



- $\mathbf{v}_1 = \mathbf{1}$ with eigenvalue $\lambda_1 = \frac{(p+q)n}{2}$.
- $\mathbf{v}_2 = \chi_{B,C}$ with eigenvalue $\lambda_2 = \frac{(p-q)n}{2}$.
- $\chi_{B,C}(i) = 1$ if $i \in B$ and $\chi_{B,C}(i) = -1$ for $i \in C$.

If we compute \mathbf{v}_2 then we recover the communities B and C .

Upshot: The second small eigenvector of $\mathbb{E}[\mathbf{L}]$ (i.e. the second largest of $\mathbb{E}[\mathbf{A}]$) is $\chi_{B,C}$ – the indicator vector for the cut between the communities.

- If the random graph G (equivilantly \mathbf{A} and \mathbf{L}) were exactly equal to its expectation, partitioning using this eigenvector would exactly recover communities B and C .

How do we show that a matrix (e.g., \mathbf{A}) is close to its expectation? **Matrix concentration inequalities.**

Matrix Concentration Inequality: If $p \geq O\left(\frac{\log^4 n}{n}\right)$, then with high probability

$$\|A - \mathbb{E}[A]\|_2 \leq O(\sqrt{pn}).$$

where $\|\cdot\|_2$ is the matrix **spectral** norm (operator norm).

For $X \in \mathbb{R}^{n \times d}$, $\|X\|_2 = \max_{z \in \mathbb{R}^d: \|z\|_2=1} \|Xz\|_2 = \sigma_1(X)$.

For the stochastic block model application, we want to show that the second eigenvectors of \mathbf{A} and $\mathbb{E}[\mathbf{A}]$ are close. How does this relate to their difference in spectral norm?

Davis-Kahan Eigenvector Perturbation Theorem: Suppose $\underline{\mathbf{A}}, \underline{\bar{\mathbf{A}}} \in \mathbb{R}^{d \times d}$ are symmetric with $\|\underline{\mathbf{A}} - \underline{\bar{\mathbf{A}}}\|_2 < \epsilon$ and eigenvectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ and $\bar{\mathbf{v}}_1, \bar{\mathbf{v}}_2, \dots, \bar{\mathbf{v}}_n$. Letting $\theta(\mathbf{v}_i, \bar{\mathbf{v}}_i)$ denote the angle between \mathbf{v}_i and $\bar{\mathbf{v}}_i$, for all i :

$$\sin[\theta(\mathbf{v}_i, \bar{\mathbf{v}}_i)] \leq \frac{\epsilon}{\min_{j \neq i} |\lambda_i - \lambda_j|}$$

where $\lambda_1, \dots, \lambda_n$ are the eigenvalues of $\bar{\mathbf{A}}$.

APPLICATION TO STOCHASTIC BLOCK MODEL

Claim 1 (Matrix Concentration): For $p \geq O\left(\frac{\log^4 n}{n}\right)$,

$$\|A - \mathbb{E}[A]\|_2 \leq O(\sqrt{pn}).$$

Claim 2 (Davis-Kahan): For $p \geq O\left(\frac{\log^4 n}{n}\right)$,

$$\sin \theta(\underline{v}_2, \bar{v}_2) \leq \frac{O(\sqrt{pn})}{\min_{j \neq 2} |\lambda_j - \lambda_2|} \leq \frac{O(\sqrt{pn})}{(p-q)n/2} = O\left(\frac{\sqrt{p}}{(p-q)\sqrt{n}}\right)$$

Recall: $\mathbb{E}[A]$, has eigenvalues $\lambda_1 = \frac{(p+q)n}{2}$, $\lambda_2 = \frac{(p-q)n}{2}$, $\lambda_i = 0$ for $i \geq 3$.

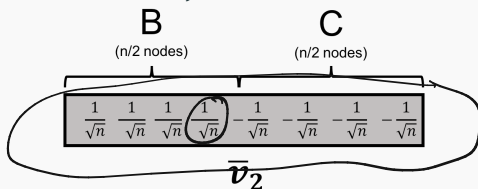
$$\min_{j \neq 2} |\lambda_2 - \lambda_j| = \min\left(\underline{qn}, \frac{(p-q)n}{2}\right).$$

Assume $\frac{(p-q)n}{2}$ will be the minimum of these two gaps.

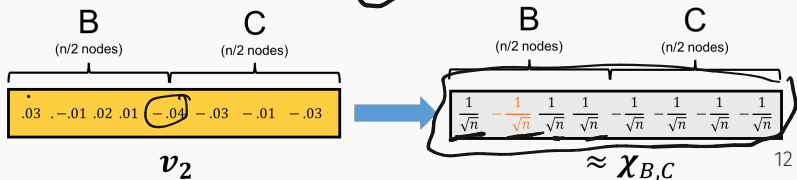
APPLICATION TO STOCHASTIC BLOCK MODEL

So far: $\sin \theta(\mathbf{v}_2, \bar{\mathbf{v}}_2) \leq O\left(\frac{\sqrt{p}}{(p-q)\sqrt{n}}\right)$. What does this give us?

- Can show that this implies $\|\mathbf{v}_2 - \bar{\mathbf{v}}_2\|_2^2 \leq O\left(\frac{p}{(p-q)^2 n}\right)$ (exercise).
- $\bar{\mathbf{v}}_2$ is $\frac{1}{\sqrt{n}}\chi_{B,C}$: the community indicator vector.



- To understand how well rounding recovers $\bar{\mathbf{v}}_2$, need to understand how many locations \mathbf{v}_2 and $\bar{\mathbf{v}}_2$ can differ in sign.



γ_{in}

Main argument:

- Every i where $v_2(i), \bar{v}_2(i)$ differ in sign contributes $\geq \frac{1}{n}$ to $\|v_2 - \bar{v}_2\|_2^2$
- We know that $\|v_2 - \bar{v}_2\|_2^2 \leq O\left(\frac{p}{(p-q)^2 n}\right)$.
- So v_2 and \bar{v}_2 differ in sign in at most $O\left(\frac{p}{(p-q)^2}\right)$ positions.

↓
 $\frac{1}{5}$

Upshot: If G is a stochastic block model graph with adjacency matrix \mathbf{A} , if we compute its second large eigenvector v_2 and assign nodes to communities according to the sign pattern of this vector, we will correctly assign all but $O\left(\frac{p}{(p-q)^2}\right)$ nodes.

- **Hard case:** $p = c/n$ for some factor c . Even when $p - q = O(1/n)$, assign all but an $O(n)$ fraction of nodes correctly. E.g., assign 99% of nodes

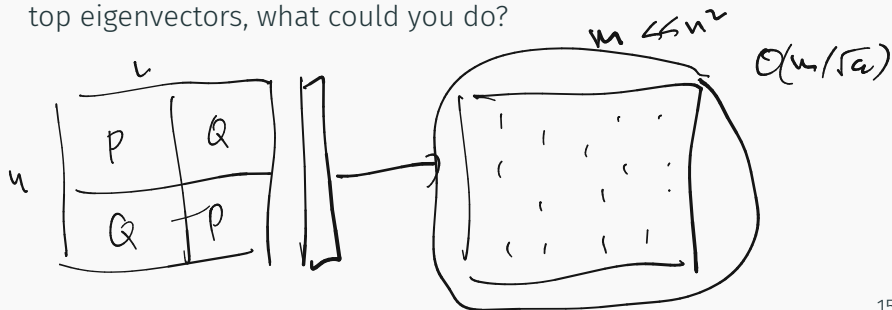
↓
c.n for $c \ll 1$

RANDOMIZED NUMERICAL LINEAR ALGEBRA

Forget about the previous problem, but still consider the matrix $\mathbf{M} = \mathbb{E}[\mathbf{A}]$.

- Dense $n \times n$ matrix.
- Computing top eigenvectors takes $\approx \underline{\underline{O(n^2/\sqrt{\epsilon})}}$ time.

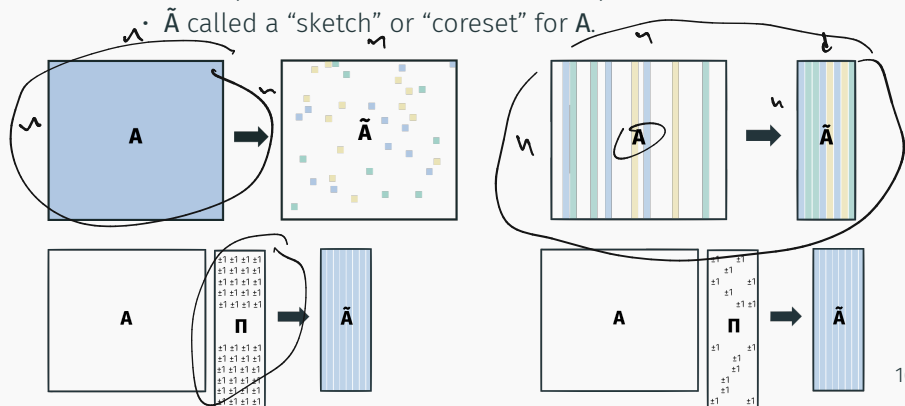
If someone asked you to speed this up and return approximate top eigenvectors, what could you do?



RANDOMIZED NUMERICAL LINEAR ALGEBRA

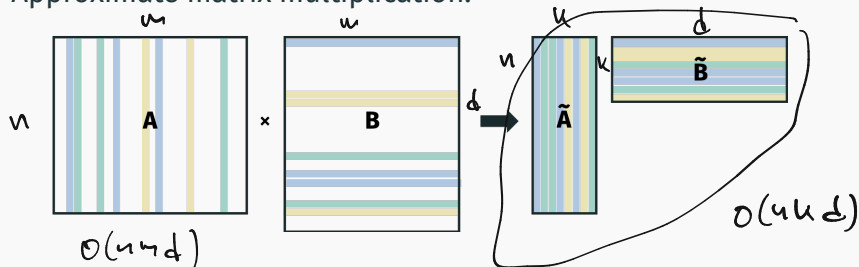
Main idea: If you want to compute singular vectors, multiply two matrices, solve a regression problem, etc.:

1. Compress your matrices using a randomized method (e.g. subsampling).
2. Solve the problem on the smaller or sparser matrix.

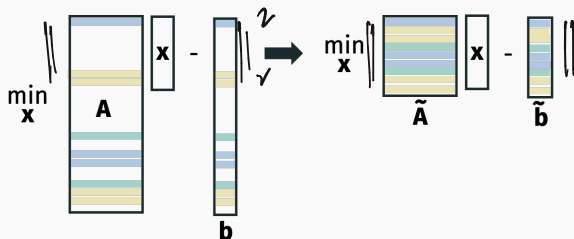


RANDOMIZED NUMERICAL LINEAR ALGEBRA

Approximate matrix multiplication:



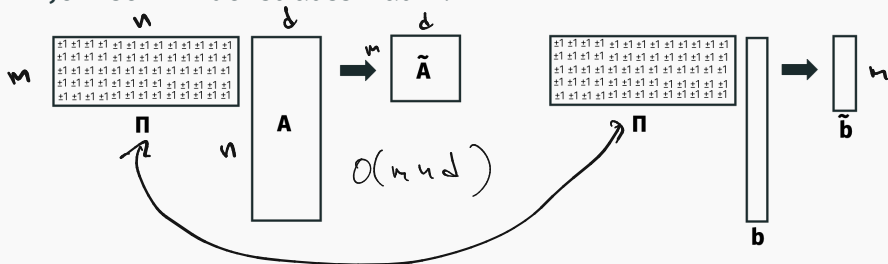
Approximate regression:



SKETCHED REGRESSION

$O(nd^2)$ $O(nd^2)$
 Randomized approximate regression using a
 Johnson-Lindenstrauss Matrix:

$O(nd \cdot \# \text{ iterations})$
 $O(nd \cdot \# \text{ iterations})$



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

Goal: Let $x^* = \arg \min_x \|Ax - b\|_2^2$. Let $\tilde{x} = \arg \min_x \|\tilde{A}x - \tilde{b}\|_2^2$

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

If $P \in \mathbb{R}^{m \times n}$, 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 $m = O\left(\frac{d}{\epsilon^2}\right)$ rows¹. Then with probability 9/10, for any $\mathbf{A} \in \mathbb{R}^{n \times d}$ and $\mathbf{b} \in \mathbb{R}^n$,

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

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

¹This can be improved to $O(d/\epsilon)$ with a tighter analysis

- Prove this theorem using an ϵ -net argument, which is a popular technique for applying our standard concentration inequality + union bound argument to an infinite number of events.
- These sort of arguments appear all the time in theoretical algorithms and ML research, so this lecture is as much about the technique as the final result.
- You will use and ϵ -net argument to prove a matrix concentration inequality on your problem set.

SKETCHED REGRESSION

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

$$x^* = \arg \min \|Ax - b\|_2$$

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

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

$$\tilde{x} = \arg \min \|\Pi Ax - \Pi b\|_2^2$$

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

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

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

$$(1 + \epsilon) \leq \frac{1}{1 - \epsilon} \leq (1 + 2\epsilon)$$

$$(1 + c\epsilon)(1 + \epsilon) = 1 + c\epsilon + \epsilon + c\epsilon^2 \leq 1 + (2c + 1)\epsilon$$

Lemma (Distributional JL)

If Π 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 y ,

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

with probability $(1 - \delta)$.



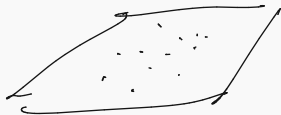
Corollary: For any fixed x , with probability $(1 - \delta)$,

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

$$\underline{y = Ax - b} \quad \Pi y = \Pi(Ax - b) = \Pi Ax - \Pi b$$

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

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



Note that all vectors of the form $(Ax - b)$ lie in a low dimensional subspace: spanned by $d + 1$ vectors, where d is the width of A . **So even though the set is infinite, it is “simple” in some way. Parameterized by just $d + 1$ numbers.**

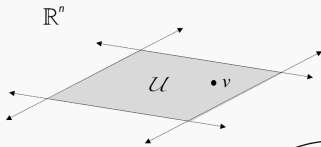
SUBSPACE EMBEDDINGS

Theorem (Subspace Embedding from JL)

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

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

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



²It'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 Π and properly scale, then with $O(d/\epsilon^2)$ rows,

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

for all \mathbf{x} and thus

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

i.e., our main theorem is proven.

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

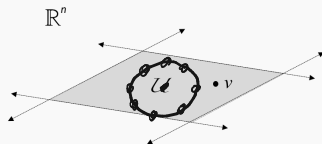
SUBSPACE EMBEDDINGS

Theorem (Subspace Embedding from JL)

Let $\mathcal{U} \subset \mathbb{R}^n$ be a d -dimensional linear subspace in \mathbb{R}^n . If $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$ 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 embeddings have tons of other applications!

SUBSPACE EMBEDDING PROOF

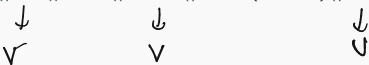
$$(1 - \epsilon)\|\mathbf{v}\|_2^2 \leq \|\Pi\mathbf{v}\|_2^2 \leq (1 + \epsilon)\|\mathbf{v}\|_2^2 \quad (2)$$

First Observation: The theorem holds as long as (2) holds for all \mathbf{w} on the unit sphere in \mathcal{U} . Denote the sphere $S_{\mathcal{U}}$:

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

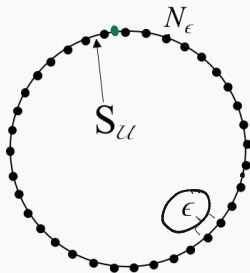
Follows from linearity: Any point $\mathbf{v} \in \mathcal{U}$ can be written as $c\mathbf{w}$ for some scalar c and some point $\mathbf{w} \in S_{\mathcal{U}}$. $c = \|\mathbf{v}\|_2$

- 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 $(1 - \epsilon)\|c\mathbf{w}\|_2 \leq \|\Pi c\mathbf{w}\|_2 \leq (1 + \epsilon)\|c\mathbf{w}\|_2$.



SUBSPACE EMBEDDING PROOF

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



N_ϵ is called an “ ϵ ”-net.

If we can prove

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

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

Lemma (ϵ -net for the sphere)

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

$$\min_{w \in N_\epsilon} \|\underline{v} - \underline{w}\|_2 \leq \epsilon.$$

Take this claim to be true for now: we will prove later.

SUBSPACE EMBEDDING PROOF

Failure prob
over all events

of events

$$n \text{ is } O\left(\frac{d}{\epsilon^2}\right)$$

1. Preserving norms of all points in net N_ϵ .

$$\delta = \left(\frac{4}{a}\right)^d \cdot \delta'$$

failure prob for one event

Set $\delta = \frac{\epsilon^2}{4}$. By a union bound, with probability $1 - \delta$, for all $w \in N_\epsilon$,

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

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

$$\delta' = \left(\frac{a}{4}\right)^d \cdot \delta$$

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

SUBSPACE EMBEDDING PROOF

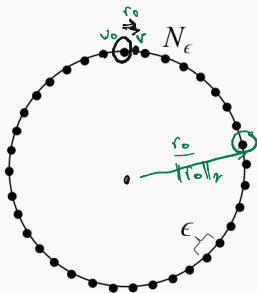
2. Writing any point in sphere as linear comb. of points in N_ϵ .

For some $w_0, w_1, w_2 \dots \in N_\epsilon$, any $v \in S_{\mathcal{U}}$ can be written:

$$v = w_0 + c_1 w_1 + c_2 w_2 + \dots$$

$\leq \epsilon$ $\leq \epsilon^2$

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



$$w_0 = \arg \min_{w \in N_\epsilon} \|w - v\|_2 \quad r_0 = v - w_0$$

$$w_1 = \arg \min_{w \in N_\epsilon} \|w - r_0 / \|r_0\|_2\|_2 \quad r_1 = \frac{r_0}{\|r_0\|_2} - w_1$$

$$v = \underbrace{w_0}_{\leq \epsilon} + \underbrace{\|r_0\|_2}_{\leq \epsilon^2} \underbrace{w_1}_{\leq \epsilon} + \underbrace{\|r_0\|_2 \|r_1\|_2}_{\leq \epsilon^3} \underbrace{w_2}_{\leq \epsilon} + \dots$$

SUBSPACE EMBEDDING PROOF

$$\Pi v = \Pi (w_0 + c_1 w_1 + c_2 w_2 + \dots)$$

3. Preserving norm of v.

$$\begin{aligned} c_1 &\leq \epsilon \\ c_2 &\leq \epsilon^2 \\ c_3 &\leq \epsilon^3 \end{aligned}$$

Applying triangle inequality, we have

$$\begin{aligned} \|\Pi v\|_2 &= \|\Pi w_0 + c_1 \Pi w_1 + c_2 \Pi w_2 + \dots\| \leq \|\Pi w_0\| + \|c_1 \Pi w_1\| \\ &\leq \|\Pi w_0\| + \epsilon \|\Pi w_1\| + \epsilon^2 \|\Pi w_2\| + \dots + \|c_2 \Pi w_2\| + \dots \\ &\leq (1 + \epsilon) + \epsilon(1 + \epsilon) + \epsilon^2(1 + \epsilon) + \dots \\ &\leq \underbrace{(1 + O(\epsilon))}_{= (1 + \epsilon)} \|v\|_2 = (1 + \epsilon) (1 + \epsilon + \epsilon^2 + \epsilon^3 + \dots) \end{aligned}$$

$$\|\Pi w_i\| \leq (1 + \epsilon) \|w_i\| = (1 + \epsilon) \cdot 1$$

$$\leq 1 + 2\epsilon$$

$$\leq (1 + \epsilon)(1 + 2\epsilon)$$

$$\leq (1 + 5\epsilon)$$

3. Preserving norm of v .

Similarly,

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

$$\Pi w_0 = (\Pi v - \Pi c_1 w_1 - \Pi c_2 w_2 - \dots)$$

$$\|\Pi w_0\| \leq \|\Pi v\|_2 + \|\Pi c_1 w_1\| + \dots$$

So we have proven

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

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

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

Adjusting ϵ proves the Subspace Embedding theorem.

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 $1 - \delta$,

$$(1 - \epsilon) \|\mathbf{v}\|_2^2 \leq \|\Pi \mathbf{v}\|_2^2 \leq (1 + \epsilon) \|\mathbf{v}\|_2^2 \quad (3)$$

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 embeddings have many other applications!

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

Lemma (ϵ -net for the sphere)

For any $\epsilon \leq 1$, there exists a set $N_\epsilon \subset S_{\mathcal{U}}$ with $|N_\epsilon| = \frac{4}{\epsilon^2}$ such that $\forall \mathbf{v} \in S_{\mathcal{U}}$,

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

$$\log\left(\frac{4}{\epsilon^2}\right)$$

$$\log(4/\epsilon)$$

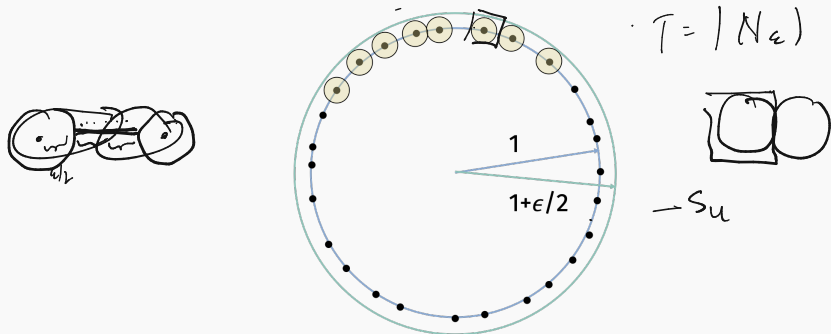
Imaginary algorithm for constructing N_ϵ :

- 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_j without intersecting any other balls. All of these balls live in a ball of radius $1 + \epsilon/2$.

ϵ -NET FOR THE SPHERE

Volume of d dimensional ball of radius r is

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

where c is a constant that depends on d , but not r . From

previous slide we have:

$$1 + \epsilon/2 \leq 2$$

$$\text{vol}(d, \epsilon/2) \cdot |N_\epsilon| \leq \text{vol}(d, 1 + \epsilon/2)$$
$$|N_\epsilon| \leq \frac{\text{vol}(d, 1 + \epsilon/2)}{\text{vol}(d, \epsilon/2)}$$

sum of volumes
of small balls

$$\leq \left(\frac{1 + \epsilon/2}{\epsilon/2} \right)^d \leq \left(\frac{2}{\epsilon} \right)^d$$

You can actually show that $m = O\left(\frac{d + \log(1/\delta)}{\epsilon}\right)$ suffices to be a d dimensional subspace embedding, instead of the bound we proved of $m = O\left(\frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon}\right)$.

The trick is to show that a constant factor net is actually all that you need instead of an ϵ factor.

RUNTIME CONSIDERATION

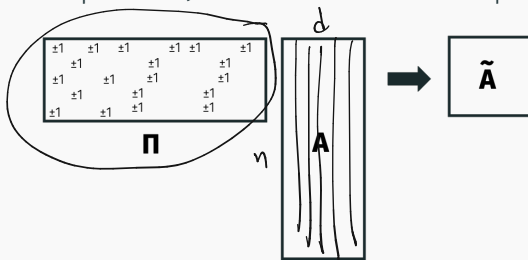
For $\epsilon, \delta = O(1)$, we need $\mathbf{\Pi}$ to have $m = O(d)$ rows.

- Cost to solve $\|\mathbf{Ax} - \mathbf{b}\|_2^2$:
 - $O(nd^2)$ time for direct method. Need to compute $(\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{b}$.
 - $O(nd) \cdot (\# \text{ of iterations})$ time for iterative method (GD, AGD, conjugate gradient method).
- 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 ΠA is an $(m \times n) \times (n \times d)$ matrix multiply: $O(mnd) = O(nd^2)$ time!

Goal: Develop faster Johnson-Lindenstrauss projections.



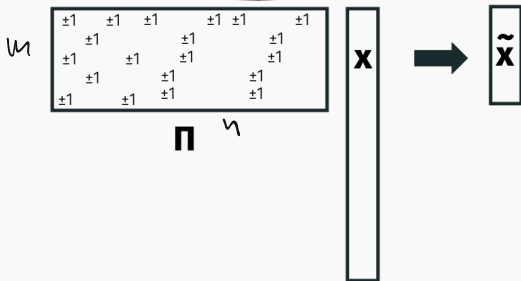
Typically using sparse and structured matrices.

We will describe a construction where ΠA can be computed in $O(nd \log n)$ time. $O(nd)$

RETURN TO SINGLE VECTOR PROBLEM

Goal: Develop methods that reduce a vector $\mathbf{x} \in \mathbb{R}^n$ down to $m \approx \frac{\log(1/\delta)}{\epsilon^2}$ dimensions in $O(mn)$ time and guarantee:

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



$O(mn)$

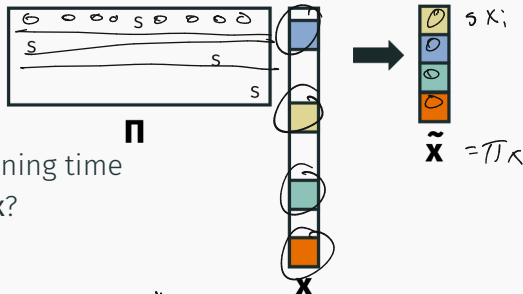
We will learn about a truly brilliant method that runs in $O(n \log n)$ time. **Preview:** Will involve Fast Fourier Transform in disguise.

FIRST ATTEMPT

$$\mathbb{E}\{z_i^2\} = \frac{1}{n} x_1^2 + \frac{1}{n} x_2^2 + \dots + \frac{1}{n} x_n^2$$

Let Π be a **random sampling matrix**. Every row contains a value of $s = \sqrt{n/m}$ in a single location, and is zero elsewhere.

subsampling matrix



What's the running time to compute Πx ?

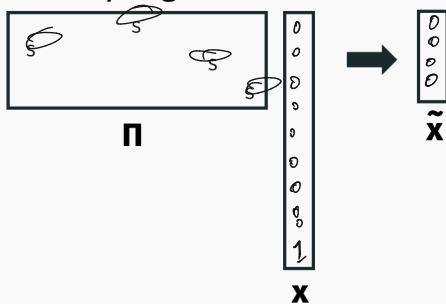
$$\|\Pi x\|_2^2 = \sum_{i=1}^m (s z_i)^2 = s^2 \sum_{i=1}^m z_i^2 = \frac{n}{m} \sum_{i=1}^m z_i^2 \quad \text{where } z_i \sim \text{Unif}(x_1, \dots, x_n)$$

$$\mathbb{E}[\|\Pi x\|_2^2] = \frac{n}{m} \sum_{i=1}^m \mathbb{E}\{z_i^2\} = \frac{n}{m} \sum_{i=1}^m \frac{1}{n} \sum_{j=1}^n x_j^2 = \|x\|_2^2$$

FIRST ATTEMPT

So $\mathbb{E}\|\Pi\mathbf{x}\|_2^2 = \|\mathbf{x}\|_2^2$ in expectation. To show it is close with high probability we would need to apply a concentration inequality. How do you think this will work out?

subsampling matrix



VARIANCE ANALYSIS

$$\|\Pi \mathbf{x}\|_2^2 = \frac{1}{n} \sum_{i=1}^n z_i^2 \quad \text{where } z_i \sim \text{Unif}(x_1, \dots, x_n)$$

$$\sigma^2 = \text{Var}[\|\Pi \mathbf{x}\|_2^2] = \frac{n^2}{n^2} \sum_{i=1}^n \text{Var}[z_i^2] = \frac{1}{n} \cdot \|\mathbf{x}\|_4^4$$

$$\text{Var}[z_i^2] = \mathbb{E}[(z_i^2)^2] - \mathbb{E}[z_i^2]^2 = \mathbb{E}[z_i^4] = \frac{1}{n} \sum_{j=1}^n x_j^4 = \frac{1}{n} \|\mathbf{x}\|_4^4$$

Recall Chebyshev's Inequality:

$$\Pr[|\|\Pi \mathbf{x}\|_2^2 - \|\mathbf{x}\|_2^2| \leq 10 \cdot \sigma] \leq \frac{1}{100}$$

We want additive error $|\|\Pi \mathbf{x}\|_2^2 - \|\mathbf{x}\|_2^2| \leq \epsilon \|\mathbf{x}\|_2^2$

VARIANCE ANALYSIS

We need to choose m so that:

$$10 \sqrt{\frac{n}{m}} \|x\|_4^2 \leq \epsilon \|x\|_2^2.$$

$$\|x\|_4 = \left(\sum_{i=1}^n x_i^4 \right)^{1/4}$$

$$10 \frac{\sqrt{4}}{\sqrt{n}} \sqrt{n} \in \epsilon n$$

How do these two norms compare?

$$\|x\|_4^2 = \left(\sum_{i=1}^n x_i^4 \right)^{1/2}$$

$$\|x\|_2^2 = \sum_{i=1}^n x_i^2$$

Consider 2 extreme cases:

$$\|x\|_4^2 = 1$$

$$\|x\|_2^2 = 1$$

$$x = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$x = \begin{bmatrix} 1 \\ -1 \\ \vdots \\ +1 \end{bmatrix}$$
$$\|x\|_4^2 = \sqrt{n}$$
$$\|x\|_2^2 = n$$

VARIANCE FOR SMOOTH FUNCTIONS

We need to choose m so that:

$$\frac{1}{10} \sqrt{\frac{n}{m}} \|x\|_4^2 \leq \epsilon \|x\|_2^2.$$

Suppose x is very evenly distributed. I.e., for all $i \in 1, \dots, n$,

$$\underline{x_i^2} \leq \frac{c}{n} \sum_{i=1}^n x_i^2 = \frac{c}{n} \|x\|_2^2 \quad \underline{x_i^2} \leq \frac{c}{n} \|x\|_2^2$$

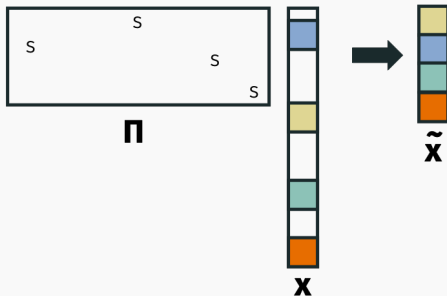
Claim: $\|x\|_4^2 \leq \frac{c}{\sqrt{n}} \|x\|_2^2$. So $m = O(c/\epsilon^2)$ samples suffices.³

³Using the right Bernstein bound we can prove $m = O(c \log(1/\delta)/\epsilon^2)$ suffices for failure probability δ .

VECTOR SAMPLING

So sampling does work to preserve the norm of \mathbf{x} , but only when the vector is relatively “smooth” (not concentrated). Do we expect to see such vectors in the wild?

subsampling matrix



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

Key idea: First multiply \mathbf{x} by a “mixing matrix” \mathbf{M} which ensures it cannot be too concentrated in one place.

\mathbf{M} should have the property that $\|\mathbf{M}\mathbf{x}\|_2^2 = \|\mathbf{x}\|_2^2$ exactly, or is very close. Then we will multiply by a subsampling matrix \mathbf{S} to do the actual dimensionality reduction:

$$\underline{\mathbf{P}}\mathbf{x} = \underline{\mathbf{S}}\mathbf{M}\mathbf{x}$$

Oh... and \mathbf{M} needs to be fast to multiply by!

$$O(n \log n)$$

THE FAST JOHNSON-LINDENSTRAUSS TRANSFORM

Good mixing matrices should look random:

$$\begin{array}{|c|} \hline +1 & -1 & +1 & +1 & +1 & -1 & +1 & -1 \\ \hline -1 & -1 & -1 & +1 & +1 & +1 & -1 & -1 \\ \hline +1 & -1 & +1 & +1 & +1 & -1 & -1 & -1 \\ \hline +1 & +1 & +1 & +1 & -1 & +1 & -1 & +1 \\ \hline -1 & -1 & +1 & +1 & -1 & +1 & +1 & -1 \\ \hline -1 & +1 & -1 & -1 & -1 & +1 & -1 & -1 \\ \hline -1 & +1 & -1 & +1 & -1 & -1 & -1 & +1 \\ \hline \end{array} \quad \begin{array}{|c|} \hline \\ \hline \\ \hline \\ \hline \\ \hline \\ \hline \\ \hline \\ \hline \end{array}$$

M **x**

For this approach to work, we need to be able to compute $\mathbf{M}\mathbf{x}$ very quickly. So we will use a **pseudorandom** matrix instead.

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

$\Pi = SM$ where $M = HD$:

- $D \in n \times n$ is a diagonal matrix with each entry uniform ± 1 .
- $H \in n \times n$ is a Hadamard matrix.

The Hadamard matrix is an orthogonal matrix closely related to the discrete Fourier matrix. It has two critical properties:

1. $\|Hv\|_2^2 = \|v\|_2^2$ exactly. Thus $\|HDx\|_2^2 = \|x\|_2^2$
2. $\|Hv\|_2^2$ can be computed in $O(n \log n)$ time.

HADAMARD MATRICES RECURSIVE DEFINITION

Assume that n is a power of 2. For $k = 0, 1, \dots$, the k^{th} Hadamard matrix \mathbf{H}_k is a $2^k \times 2^k$ matrix defined by:

$$H_0 = 1 \quad H_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \quad H_2 = \frac{1}{\sqrt{4}} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{bmatrix}$$

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

The $n \times n$ Hadamard matrix has all entries as $\pm \frac{1}{\sqrt{n}}$.

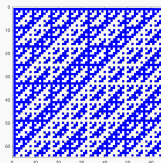
HADAMARD MATRICES ARE ORTHOGONAL

Property 1: For any $k = 0, 1, \dots$, we have $\|\mathbf{H}_k \mathbf{v}\|_2^2 = \|\mathbf{v}\|_2^2$ for all \mathbf{v} .
I.e., \mathbf{H}_k is orthogonal.

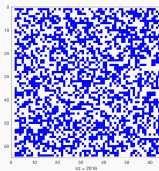
Property 2: Can compute $\mathbf{P}\mathbf{x} = \mathbf{S}\mathbf{H}\mathbf{D}\mathbf{x}$ in $O(n \log n)$ time.

RANDOMIZED HADAMARD TRANSFORM

Property 3: The randomized Hadamard matrix is a good “mixing matrix” for smoothing out vectors.



Deterministic
Hadamard matrix.



Randomized
Hadamard **PHD**.



Fully random sign
matrix.

Blue squares are $1/\sqrt{n}$'s, white squares are $-1/\sqrt{n}$'s.

Lemma (SHRT mixing lemma)

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

$$(z_i)^2 \leq \frac{c \log(n/\delta)}{n} \|\mathbf{z}\|_2^2$$

for some fixed constant c .

The vector is very close to uniform with high probability. As we saw earlier, we can thus argue that $\|\mathbf{S}\mathbf{z}\|_2^2 \approx \|\mathbf{z}\|_2^2$. I.e. that:

$$\|\mathbf{I}\mathbf{x}\|_2^2 = \|\mathbf{S}\mathbf{H}\mathbf{D}\mathbf{x}\|_2^2 \approx \|\mathbf{x}\|_2^2$$

Theorem (The Fast JL Lemma)

Let $\mathbf{\Pi} = \mathbf{SHD} \in \mathbb{R}^{m \times n}$ 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 \mathbf{x} ,

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

with probability $(1 - \delta)$.

Very little loss in embedding dimension compared to full random matrix, and $\mathbf{\Pi}$ can be multiplied by \mathbf{x} in $O(n \log n)$ (nearly linear) time.

SHRT mixing lemma proof: Need to prove $(z_i)^2 \leq \frac{c \log(n/\delta)}{n} \|\mathbf{z}\|_2^2$ for all i .

Let \mathbf{h}_i^T be the i^{th} row of \mathbf{H} . $z_i = \mathbf{h}_i^T \mathbf{D} \mathbf{x}$ where:

$$\mathbf{h}_i^T \mathbf{D} = \frac{1}{\sqrt{n}} \begin{bmatrix} 1 & 1 & \dots & -1 & -1 \end{bmatrix} \begin{bmatrix} D_1 & & & & \\ & D_2 & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & D_n \end{bmatrix}$$

where D_1, \dots, D_n are random ± 1 's.

This is equivalent to

$$\mathbf{h}_i^T \mathbf{D} = \frac{1}{\sqrt{n}} \begin{bmatrix} R_1 & R_2 & \dots & R_n \end{bmatrix},$$

where R_1, \dots, R_n are random ± 1 's.

So we have, for all i , $\mathbf{z}_i = \mathbf{h}_i^T \mathbf{D}\mathbf{x} = \frac{1}{\sqrt{n}} \sum_{j=1}^n R_{ij}x_j$.

- \mathbf{z}_i is a random variable with mean 0 and variance $\frac{1}{n} \|\mathbf{x}\|_2^2$, and is a sum of independent random variables.
- By Central Limit Theorem, we expect that:

$$\Pr[|\mathbf{z}_i| \geq t \cdot \frac{\|\mathbf{x}\|_2}{\sqrt{n}}] \leq e^{-O(t^2)}.$$

- Setting $t = \sqrt{\log(n/\delta)}$, we have for constant c ,

$$\Pr \left[|\mathbf{z}_i| \geq c \sqrt{\frac{\log(n/\delta)}{n}} \|\mathbf{y}\|_2 \right] \leq \frac{\delta}{n}$$

- Applying a union bound to all n entries of \mathbf{z} gives the SHRT mixing lemma.

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 ± 1 's). Then for any vector $\mathbf{a} \in \mathbb{R}^n$,

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

This is call the Khintchine Inequality. It is specialized to sums of scaled ± 1 's, and is a bit tighter and easier to apply than using a generic Bernstein bound.

With probability $1 - \delta$, we have that all $\mathbf{z}_i \leq \sqrt{\frac{c \log(n/\delta)}{n}} \|\mathbf{c}\|_2$.

As shown earlier, we can thus guarantee that:

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

as long as $\mathbf{S} \in \mathbb{R}^{m \times n}$ is a random sampling matrix with

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

$\|\mathbf{S}\mathbf{z}\|_2^2 = \|\mathbf{S}\mathbf{H}\mathbf{D}\mathbf{x}\|_2^2 = \|\mathbf{\Pi}\mathbf{x}\|_2^2$ and $\|\mathbf{z}\|_2^2 = \|\mathbf{x}\|_2^2$, so we are done.

Theorem (The Fast JL Lemma)

Let $\mathbf{\Pi} = \mathbf{SHD} \in \mathbb{R}^{m \times n}$ 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 \mathbf{x} ,

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

with probability $(1 - \delta)$.

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

$O(nd \log n)$ 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:

$$O(\text{nnz}(\mathbf{A})) \text{ time.}$$

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

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

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

WHAT WERE AILON AND CHAZELLE THINKING?

Simple, inspired algorithm that has been used for accelerating:

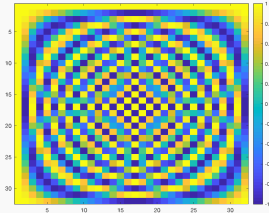
- Vector dimensionality reduction
- Linear algebra
- Locality sensitive hashing (SimHash)
- Randomized kernel learning methods (we will discuss after Thanksgiving)

```
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.

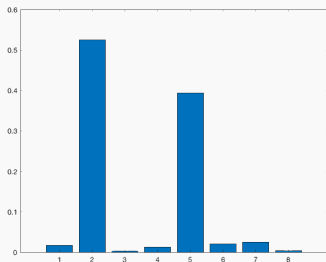
$$F_{j,k} = e^{-2\pi i \frac{j \cdot k}{n}}, \quad F^*F = I.$$



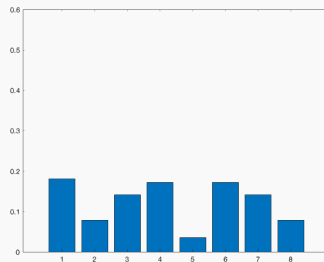
Real part of $F_{j,k}$.

Fy computes the Discrete 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 (informal): A function and its Fourier transform cannot both be concentrated.



Vector y .



Fourier transform Fy .

What do we know?

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

Taking a Fourier transform exactly eliminates this hard case, without changing \mathbf{y} 's norm.

One of the central tools in the field of **sparse recovery** aka **compressed sensing**.