

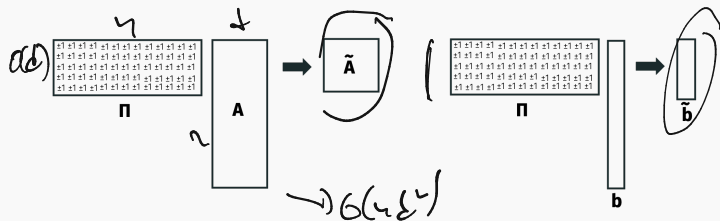
CS-GY 6763: Lecture 13

Fast Johnson-Lindenstrauss Transform, Sparse
Recovery and Compressed Sensing

NYU Tandon School of Engineering, Prof. Christopher Musco

RANDOMIZED NUMERICAL LINEAR ALGEBRA

Main idea: Speed up classical linear algebra problems using randomization.



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

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

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

Theorem (Example: Randomized Linear Regression)

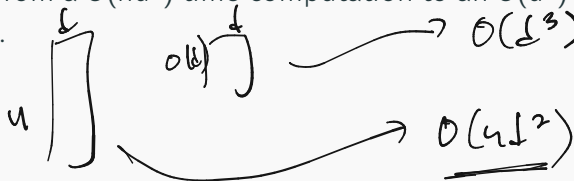
Let Π be a properly scaled JL matrix (random Gaussian, sign, sparse random, etc.) with $\underline{m} = O\left(\frac{d}{\epsilon^2}\right)$ rows. Then with probability $9/10$, for any $\mathbf{A} \in \mathbb{R}^{n \times d}$ and $\mathbf{b} \in \mathbb{R}^n$,

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

$$\epsilon = 1/2$$

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

Reduce from a $O(nd^2)$ time computation to an $O(d^3)$ time problem.



Theorem (Second Example: Randomized Low-Rank Approximation¹)

Let Π be a properly scaled JL matrix (random Gaussian, sign, sparse random, etc.) with $m = O\left(\frac{k}{\epsilon}\right)$ rows. Then with probability $9/10$, for any $\mathbf{A} \in \mathbb{R}^{n \times d}$,

$$\hat{\mathbf{A}} = \Pi \mathbf{A}$$

$$\|\mathbf{A} - \mathbf{A}\tilde{\mathbf{V}}_k\tilde{\mathbf{V}}_k^T\|_2^2 \leq (1 + \epsilon)\|\mathbf{A} - \mathbf{A}_k\|_2^2$$

where $\tilde{\mathbf{V}}_k$ contains the top k right singular vectors of $\hat{\mathbf{A}}$.

Reduce from a $O(ndk)$ time computation to an $O(dk^2)$ time problem.

¹See e.g. Sarlos, 2006 or Halko, Martinson, Tropp, 2011.

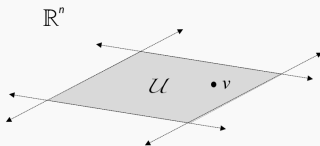
Key Ingredient:

Theorem (Subspace Embedding 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 $1 - \delta$,

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

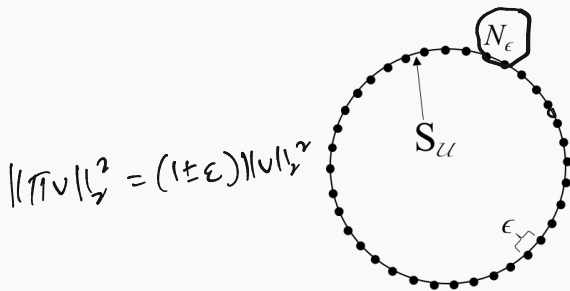
for all $\mathbf{v} \in \mathcal{U}$ as long as $m = O\left(\frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon^2}\right)$.



SUBSPACE EMBEDDING PROOF

Proof idea: Construct ϵ -net, N_ϵ , for the unit sphere, S .

1. Prove that $\|\Pi \mathbf{w}\|_2^2 = (1 \pm \epsilon)\|\mathbf{w}\|_2^2$ for all $\mathbf{w} \in N_\epsilon$ using union bound.
2. Use a direct argument to extend to the rest of sphere.



Lemma (ϵ -net for the sphere)

Let S be a d dimensional ^{unit} sphere. For any $\epsilon \leq 1$, there exists a set $N_\epsilon \subset S$ with $|N_\epsilon| \leq \left(\frac{3}{\epsilon}\right)^d$ such that $\forall \mathbf{v} \in S$,

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

We skipped the proof of this last time.

We will prove it using a common technique known as a “volume” argument.

$$\log\left(\left(\frac{3}{\epsilon}\right)^d\right) = d \log\left(\frac{3}{\epsilon}\right)$$

Lemma (ϵ -net for the sphere)

Let S be a d dimensional union sphere. For any $\epsilon \leq 1$, there exists a set $N_\epsilon \subset S$ with $|N_\epsilon| = \left(\frac{3}{\epsilon}\right)^d$ such that $\forall \mathbf{v} \in S$,

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

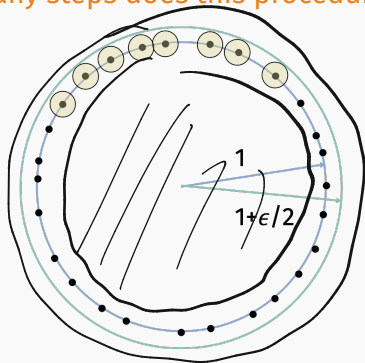


Imaginary algorithm for constructing N_ϵ :

- Set $N_\epsilon = \{\}$
- While such a point exists, choose an arbitrary point $\mathbf{v} \in S$ where there is no $\mathbf{w} \in N_\epsilon$ with $\|\mathbf{v} - \mathbf{w}\| \leq \epsilon$.
- Add \mathbf{v} to N_ϵ .

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 $\mathbf{v} \in S$ 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 $1 + \epsilon/2$.

Volume of d dimensional ball of radius r is

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

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

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

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

$$\frac{c \cdot \left(\frac{1 + \epsilon/2}{\epsilon/2}\right)^d}{c \cdot \left(\frac{\epsilon}{2}\right)^d} \leq \left(\frac{3}{\epsilon}\right)^d$$

$$\hookrightarrow \frac{2 + \epsilon}{\epsilon} \leq \frac{3}{\epsilon}$$

Theorem (Example: Randomized Linear Regression)

Let $\mathbf{\Pi}$ be a properly scaled JL matrix (random Gaussian, sign, sparse random, etc.) with $m = O\left(\frac{d}{\epsilon^2}\right)$ rows. 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 $\tilde{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{\Pi}\mathbf{A}\mathbf{x} - \mathbf{\Pi}\mathbf{b}\|_2^2$.

RUNTIME CONSIDERATION

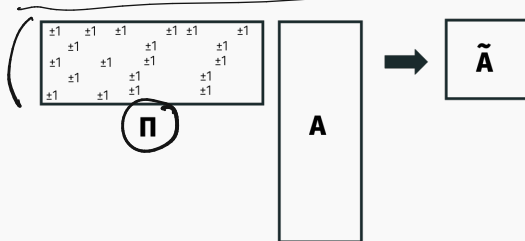
For $\epsilon, \delta = O(1)$, we need Π to have $m = O(d)$ rows.

- 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 $\|\Pi \mathbf{Ax} - \Pi \mathbf{b}\|_2^2$:
 - $O(d^3)$ 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: $O(\underline{mnd}) = O(nd^2)$ time. $O(ndk)$

Goal: Develop faster Johnson-Lindenstrauss projections.



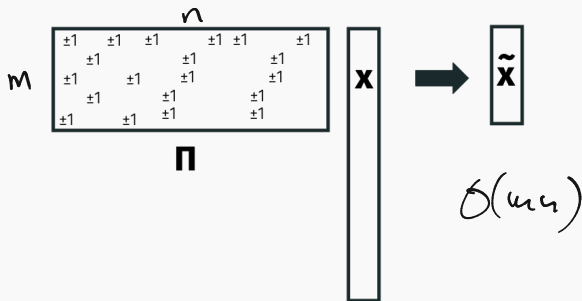
Typically using sparse or structured matrices instead of fully random JL matrices.

Useful in many other applications too. For example, faster methods are often used in LSH systems to implement SimHash.

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:

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



Recall that once the bound above is proven, linearity lets us preserve things like $\|\mathbf{y} - \mathbf{z}\|_2^2$ or $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$ for all \mathbf{x} .

Subsampled Randomized Hadamard Transform² (SHRT) (Ailon-Chazelle, 2006)

Theorem (The Fast JL Lemma)

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

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

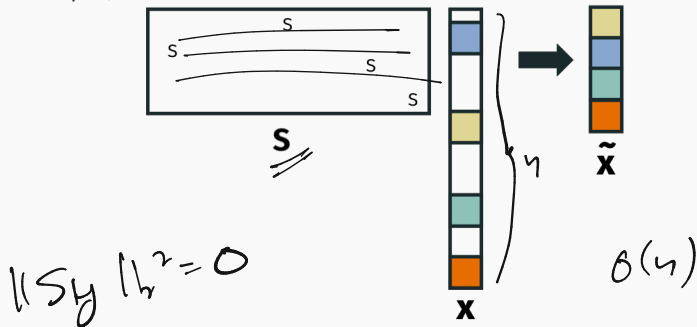
with probability $(1 - \delta)$ and $\mathbf{\Pi}\mathbf{x}$ can be computed in $O(n \log n)$ (nearly linear) time.

Very little loss in embedding dimension compared to standard JL.

²One of my favorite randomized algorithms.

SOLUTION FOR “FLAT” VECTORS

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

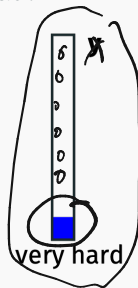
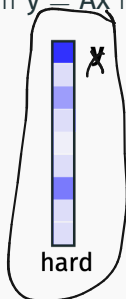
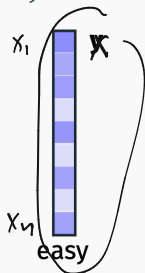


If we take m samples, \tilde{x} can be computed in $O(m)$ time.
Woohoo!

What is the problem with this approach?

VECTOR SAMPLING

Sampling only works well if $y = Ax$ is “flat”.



$$\left[\begin{array}{l} + \frac{1}{\sqrt{n}} \\ - \frac{1}{\sqrt{n}} \end{array} \right]$$
$$x_i^2 \leq \frac{1}{n} \|x\|_2^2$$

Claim

If $x_i^2 \leq \frac{c}{n} \|x\|_2^2$ for all i then $m = O(c \log(1/\delta)/\epsilon^2)$ samples suffices to ensure the $(1 - \epsilon)\|x\|_2^2 \leq \|Sx\|_2^2 \leq (1 + \epsilon)\|x\|_2^2$ with probability $1 - \delta$.

This just follows from standard Hoeffding inequality.

THE FAST JOHNSON-LINDENSTRAUSS TRANSFORM

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

\mathbf{M} will have the properties that

1. $\|\underline{\mathbf{M}\mathbf{x}}\|_2^2 = \|\mathbf{x}\|_2^2$ exactly.
2. Every entry in $\mathbf{M}\mathbf{x}$ is bounded. I.e. $\underline{[\mathbf{M}\mathbf{x}]_i^2} \leq \frac{c}{n} \|\underline{\mathbf{M}\mathbf{x}}\|_2^2$ for some factor c to be determined. $\leq \frac{c}{n} \|\mathbf{x}\|_2^2$
3. We will be able to multiply by \mathbf{M} in $O(\underline{n \log n})$ time.

Then we will multiply by a subsampling matrix \mathbf{S} to do the actual dimensionality reduction:

$$\underline{\Pi\mathbf{x} = \mathbf{S}\mathbf{M}\mathbf{x}} \quad \begin{matrix} \nearrow n \times n \\ O(n^2) \end{matrix}$$

THE FAST JOHNSON-LINDENSTRAUSS TRANSFORM

Good mixing matrices should look random:

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

M **x**

$$Mx = e_1$$

$$x = \underline{\underline{M^{-1}e_1}}$$

In fact, I claim to mix any x with high probability, M needs to be chosen randomly. Why?

Hint: Recall that $\|Mx\|_2 = \|x\|_2$, so M is orthogonal.

THE FAST JOHNSON-LINDENSTRAUSS TRANSFORM

Good mixing matrices should look random:

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

M **x**

But for this approach to work, we need to be able to compute \mathbf{Mx} very quickly. So we will use a pseudorandom matrix instead.

Subsampled Randomized Hadamard Transform

$$\Pi = \underline{S} \underline{M} \text{ where } \underline{M} = \underline{H} \underline{D}$$



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

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

1. $\|\underline{H}\underline{v}\|_2^2 = \|\underline{v}\|_2^2$ exactly. Thus $\|\underline{H}\underline{D}\underline{x}\|_2^2 = \|\underline{x}\|_2^2$
2. $\|\underline{H}\underline{v}\|_2$ can be computed in $O(n \log n)$ time.
3. All of the entries in \underline{H} have the same magnitude. I.e. the matrix is “flat”/

HADAMARD MATRICES RECURSIVE DEFINITION

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

$$2^0 = 1$$

$$2^1$$

$$\underline{H_0 = 1}$$

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

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

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

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

$$\left(\frac{1}{\sqrt{2}}\right)^k = \frac{1}{\sqrt{n}}$$

HADAMARD MATRICES ARE ORTHOGONAL

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

$$H_n^T H_n = I$$

Assume $H_{n-1}^T H_{n-1} = I$, prove inductively.

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

HADAMARD MATRICES

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

Assume:

using $c \cdot \frac{n}{2} \log_2(n/2)$ operations can compute

$H_{n-1}v$ where $k = \log_2(n)$.

$$H_k(x) \begin{bmatrix} H_{k-1} & H_{k-1} \\ H_{k-1} & -H_{k-1} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} a+b \\ a-b \end{bmatrix} \rightarrow \leq cn \log_2(n)$$

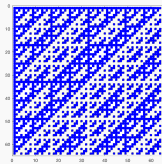
$$\frac{H_{k-1}x_1}{\downarrow} \\ a$$

$$H_{k-1}x_2 \\ \downarrow \\ b$$

$$\begin{aligned} \text{cost} &= c \cdot n \log_2(n/2) + n \\ &= c \cdot n \log_2(n) - cn + n \end{aligned}$$

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.

Pseudorandom objects like this appear all the time in computer science! Error correcting codes, efficient hash functions, etc.

RANDOMIZED HADAMARD ANALYSIS

Lemma (SHRT mixing lemma)

$$z = \Pi x \quad \Pi = HD$$

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

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

for some fixed constant c .

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

$$\|\Pi x\|_2^2 = \|SHDx\|_2^2 \approx \|x\|_2^2$$

The main result then follows directly from our sampling result from earlier:

Theorem (The Fast JL Lemma)

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

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

with probability $(1 - \delta)$.

RANDOMIZED HADAMARD ANALYSIS

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

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

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

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

This is equivalent to

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

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

RANDOMIZED HADAMARD ANALYSIS

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

- z_i is a random variable with mean 0 and variance $\frac{1}{n} \|\mathbf{x}\|_2^2$, which is a sum of independent random variables.

$$z_i = \mathbf{h}_i^T \mathbf{D}\mathbf{x} = \frac{1}{\sqrt{n}} \sum_{j=1}^n \underbrace{(R_{ij} x_j)}_{\text{random } \pm 1} \quad \frac{1}{\sqrt{n}} \|\mathbf{x}\|_2^2$$

$$\mathbb{E}[z_i] = \frac{1}{\sqrt{n}} \sum_{j=1}^n \underbrace{\mathbb{E}[R_{ij} x_j]}_0 = 0$$

$$\begin{aligned} \text{Var}[z_i] &= \frac{1}{n} \sum_{j=1}^n \text{Var}[R_{ij} x_j] = \frac{1}{n} \sum_{j=1}^n x_j^2 \cdot \mathbf{I} \\ &= \frac{1}{n} \|\mathbf{x}\|_2^2 \end{aligned}$$

RANDOMIZED HADAMARD ANALYSIS

z_i is a random variable with mean 0 and variance $\frac{1}{n}\|\mathbf{x}\|_2^2$, which is a sum of independent random variables.

- By Central Limit Theorem, we expect that:

$$\Pr\left[|z_i| \geq t \cdot \frac{\|\mathbf{x}\|_2}{\sqrt{n}}\right] \leq e^{-O(t^2)} \approx \frac{\delta}{3}$$

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

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

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

RADEMACHER CONCENTRATION

Can 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 called the Khinchine 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.

Recall that $\mathbf{z} = \mathbf{H}\mathbf{D}\mathbf{x}$.

With probability $1 - \delta$, we have that for all i ,

$$z_i \leq \sqrt{\frac{c \log(n/\delta)}{n}} \|\mathbf{x}\|_2 = \sqrt{\frac{c \log(n/\delta)}{n}} \|\mathbf{z}\|_2.$$

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

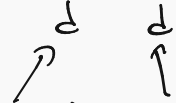
$$m = O\left(\frac{\log(n/\delta) \log(1/\delta)}{\epsilon^2}\right) \text{ rows.}$$

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

LINEAR REGRESSION WITH SHRTS


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

$$\log(N(\epsilon)) = \log\left(\left(\frac{1}{\epsilon}\right)^d\right)$$

$$\frac{\log(4/\delta) \log(1/\delta)}{\epsilon^2} = \frac{\log(4(\frac{1}{\epsilon})^d) \log((\frac{1}{\epsilon})^d)}{\epsilon^2}$$


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

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: Let $O(\text{nnz}(\mathbf{A}))$ be the number of non-zeros in \mathbf{A} . It is possible to compute $\tilde{\mathbf{A}}$ with $\text{poly}(d)$ rows in:

$O(\text{nnz}(\mathbf{A}))$ 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:

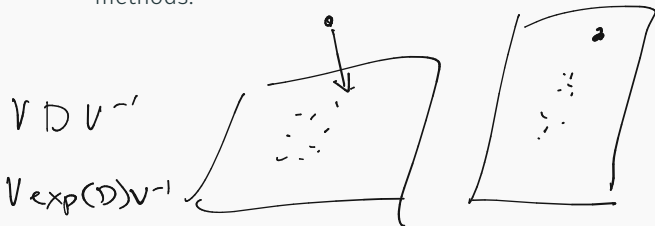
$O(\text{nnz}(\mathbf{A})) + \text{poly}(d, \epsilon)$.

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.

```
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);
```



$V D V^{-1}$
 $V \exp(D) V^{-1}$

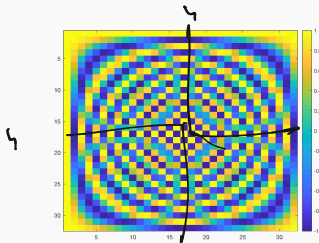
BREAK

WHAT WERE AILON AND CHAZELLE THINKING?

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

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

$$F^*F = \textcircled{I}$$



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

Fy computes the Discrete Fourier Transform of the vector y.

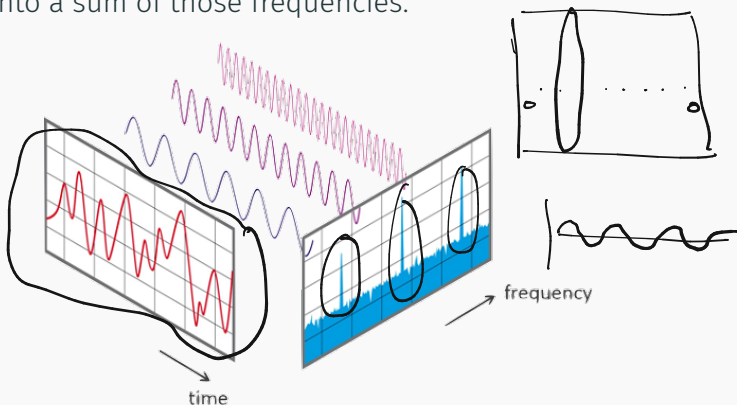
Can be computed in $O(n \log n)$ time using a divide and conquer algorithm (the Fast Fourier Transform).

FOURIER TRANSFORM

The real part of $e^{-2\pi i \frac{j \cdot k}{n}}$ equals $\cos(2\pi j \cdot k)$. So, the j^{th} row of F looks like a cosine wave with frequency $2\pi j$.

$$e^{i\theta} = \cos(\theta) + i\sin(\theta)$$

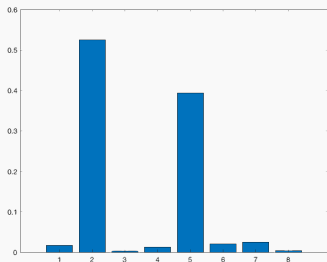
Computing Fx computes inner products of x with a bunch of different frequencies, which can be used to decompose the vector into a sum of those frequencies.



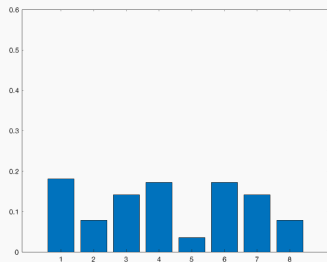
THE UNCERTAINTY PRINCIPAL



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



Vector y .



Fourier transform Fy .



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

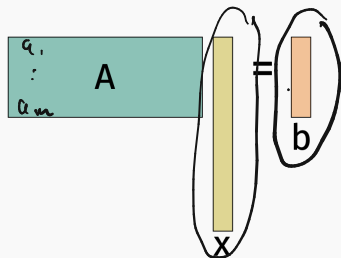
Taking a Fourier transform exactly eliminates this hard case, without changing \mathbf{y} 's norm.

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

SPARSE RECOVERY/COMPRESSED SENSING PROBLEM SETUP

Goal: Recover a vector x from linear measurements.

Choose $A \in \mathbb{R}^{m \times n}$ with $m < n$. Assume we can access $b = Ax$ via some black-box measurement process. Try to recover x from the information in b .



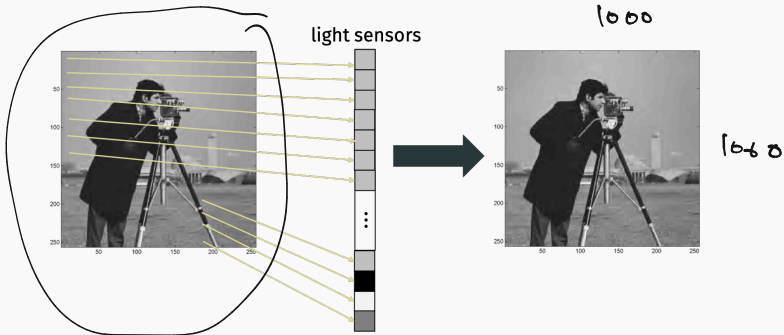
$$\begin{array}{c} \langle a_1, x \rangle \\ \vdots \\ \langle a_m, x \rangle \end{array}$$

$$Ax = Ay$$

- Infinite possible solutions y to $Ay = b$, so in general, it is impossible to recover x from b .
- Can often be possible if x has additional structure!

EXAMPLE APPLICATION: SINGLE PIXEL CAMERA

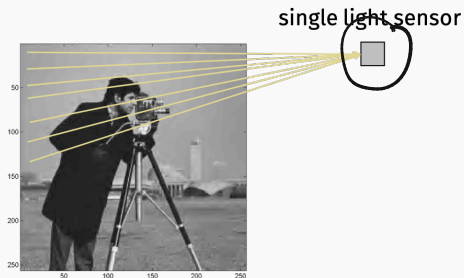
Typical acquisition of image by camera:



Requires one image sensor per pixel captured.

EXAMPLE APPLICATION: SINGLE PIXEL CAMERA

Compressed acquisition of image:

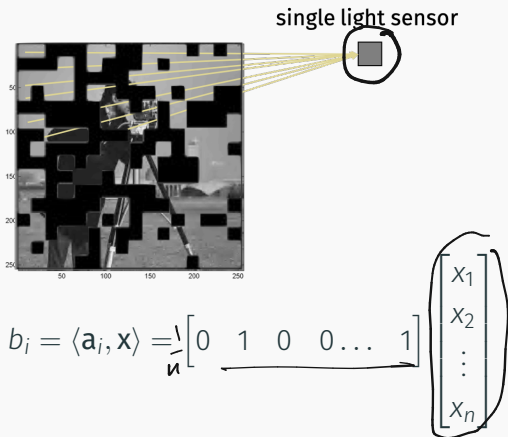


$$b = \sum_{i=1}^n x_i = \begin{bmatrix} \frac{1}{n} & \frac{1}{n} & \dots & \frac{1}{n} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \begin{matrix} \rightarrow 0, 1 \\ 0, 1 \\ \vdots \\ 0, 1 \end{matrix}$$

Does not provide very much information about the image.

EXAMPLE APPLICATION: SINGLE PIXEL CAMERA

But you can get more information from other linear measurements via masking!

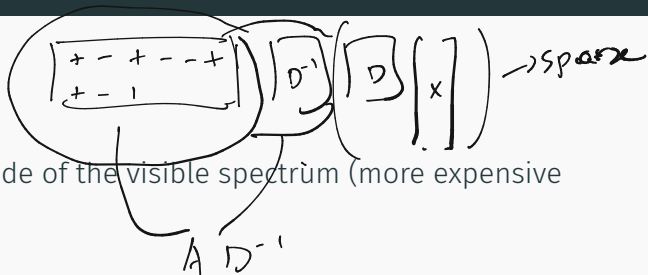


Piece together many of these masked measurements, and can recover the whole image!

EXAMPLE APPLICATION: SINGLE PIXEL CAMERA

Applications in:

- Imaging outside of the visible spectrum (more expensive sensors).
- Microscopy.
- Other scientific imaging.
- (We will discuss other applications shortly)

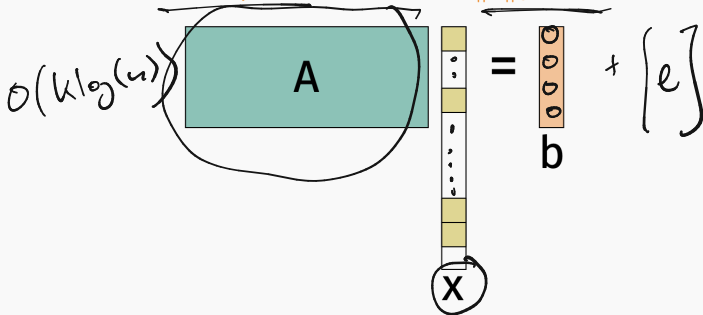


The theory we will discuss does not exactly describe these problems, but has been very valuable in modeling them.

SPARSITY RECOVERY/COMPRESSED SENSING

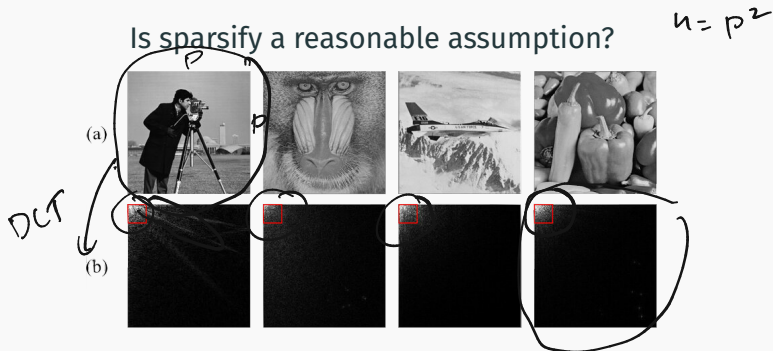
Need to make some assumption to solve the problem. Given $A \in \mathbb{R}^{m \times n}$ with $m < n$, $\mathbf{b} \in \mathbb{R}^m$, want to recover \mathbf{x} .

- Assume \mathbf{x} is k -sparse for small k . $\|\mathbf{x}\|_0 = k$.



- In many cases can recover \mathbf{x} with $\ll n$ rows. In fact, often $\sim \underline{O(k)}$ suffice.

SPARSITY ASSUMPTION



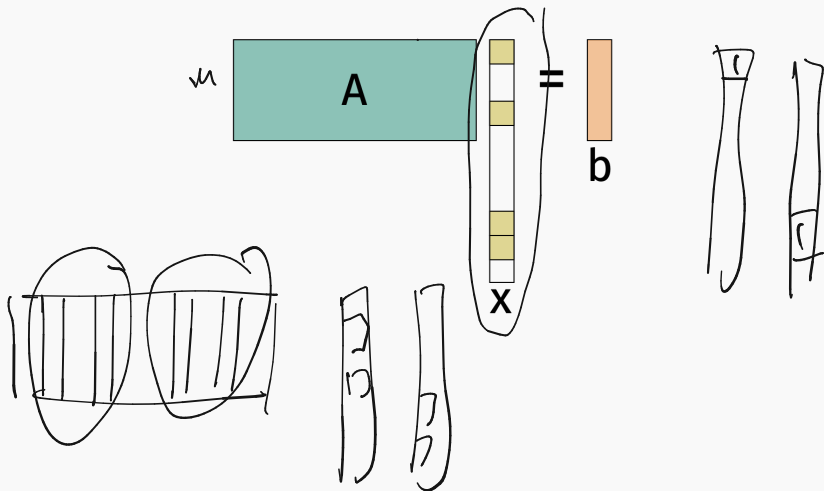
For some of the approaches we will discuss, it suffices to assume that \mathbf{x} is sparse in any fixed (and known) basis. I.e. that $\underline{\mathbf{V}}\mathbf{x}$ is sparse for some $n \times n$ orthogonal \mathbf{V} . E.g. images are sparse in the Discrete Cosine Transform basis.

Sparsity is a starting point for considering other more complex structure.

REQUIREMENTS FOR MEASUREMENT MATRIX

$$Ax = b$$

What matrices A would definitely not allow us to recover x ?



$r < d$

Many ways to formalize our intuition

i, j column

- A has Kruskal rank r . All sets of r columns in A are linearly independent.
 - Recover vectors x with sparsity $k = r/2$.
- A is μ -incoherent. $|A_i^T A_j| \leq \underline{\underline{\mu}} \|A_i\|_2 \|A_j\|_2$ for all columns $A_i, A_j, i \neq j$.
 - Recover vectors x with sparsity $k = 1/\mu$.
- (Focus today: A obeys the Restricted Isometry Property.)

RESTRICTED ISOMETRY PROPERTY

Definition ((q, ϵ) -Restricted Isometry Property)

A matrix \mathbf{A} satisfies (q, ϵ) -RIP if, for all \mathbf{x} with $\|\mathbf{x}\|_0 \leq q$,

$$\left((1 - \epsilon)\|\mathbf{x}\|_2^2 \leq \|\mathbf{Ax}\|_2^2 \leq (1 + \epsilon)\|\mathbf{x}\|_2^2. \right)$$

- Johnson-Lindenstrauss type condition.)
- \mathbf{A} preserves the norm of all q sparse vectors, instead of the norms of a fixed discrete set of vectors, or all vectors in a subspace (as in subspace embeddings).
- **Preview:** A random matrix \mathbf{A} with $\sim O(q \log(n/q))$ rows satisfies RIP.

$$O(k \log(n/k))$$

FIRST SPARSE RECOVERY RESULT

Theorem (ℓ_0 -minimization)

$$q = 2k$$

Suppose we are given $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} = \mathbf{A}\mathbf{x}$ for an unknown k -sparse $\mathbf{x} \in \mathbb{R}^n$. If \mathbf{A} is $(2k, \epsilon)$ -RIP for any $\epsilon < 1$ then \mathbf{x} is the unique minimizer of:

$$\min \|\mathbf{z}\|_0$$

subject to

$$\mathbf{A}\mathbf{z} = \mathbf{b}$$

- Establishes that information theoretically we can recover \mathbf{x} . Solving the ℓ_0 -minimization problem is computationally difficult, requiring $O(n^k)$ time. We will address faster recovery shortly.

FIRST SPARSE RECOVERY RESULT

Claim: If A is $(2k, \epsilon)$ -RIP for any $\epsilon < 1$ then x is the unique minimizer of $\min_{Az=b} \|z\|_0$.

Proof: By contradiction, assume there is some $y \neq x$ such that

$Ay = b$, $\|y\|_0 \leq \|x\|_0$.

x has $\|x\|_0 = k$ and $Ax = b$.

\downarrow
 $\|y\|_0 \leq k$

$\|A(y-x)\|_2 = 0$

$Ay - Ax = b - b = 0$

\downarrow
at most $2k$ sparse.

$A(y-x) = 0$

$\|y-x\|_2 > 0$

$(1-\epsilon)\|y-x\|_2 \leq \|A(y-x)\|_2 \leq (1+\epsilon)\|y-x\|_2 \rightarrow$ contradiction

Important note: There are robust versions of this theorem and the others we will discuss. These are much more important practically. Here's a flavor of a robust result:

- Suppose $\underline{\mathbf{b}} = \underline{\mathbf{A}}(\underline{\mathbf{x}} + \mathbf{e})$ where $\underline{\mathbf{x}}$ is k -sparse and \mathbf{e} is dense but has bounded norm.
- Recover some k -sparse $\tilde{\mathbf{x}}$ such that:

$$\|\tilde{\mathbf{x}} - \underline{\mathbf{x}}\|_2 \leq \|\mathbf{e}\|_1$$

or even

$$\|\tilde{\mathbf{x}} - \underline{\mathbf{x}}\|_2 \leq O\left(\frac{1}{\sqrt{k}}\right) \|\mathbf{e}\|_1.$$

We will not discuss robustness in detail, but along with computational considerations, it is a big part of what has made compressed sensing such an active research area in the last 20 years. Non-robust compressed sensing results have been known for a long time:

Gaspard Riche de Prony, *Essay experimental et analytique: sur les lois de la dilatabilite de fluides elastique et sur celles de la force expansive de la vapeur de l'alcool, a differentes temperatures*. Journal de l'Ecole Polytechnique, 24–76. 1795.

What matrices satisfy this property?

- Random Johnson-Lindenstrauss matrices (Gaussian, sign, etc.) with $m = O\left(\frac{k \log(n/k)}{\epsilon^2}\right)$ rows are (k, ϵ) -RIP.

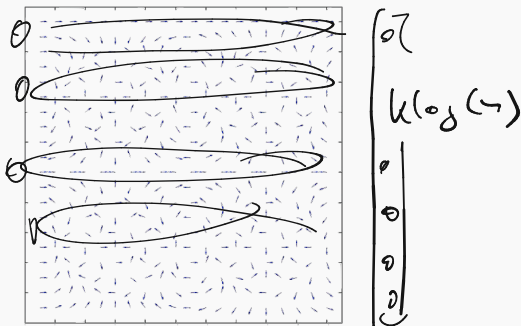
Some real world data may look random, but this is also a useful observation algorithmically when we want to design A.

THE DISCRETE FOURIER MATRIX

The $n \times n$ discrete Fourier matrix \mathbf{F} is defined:

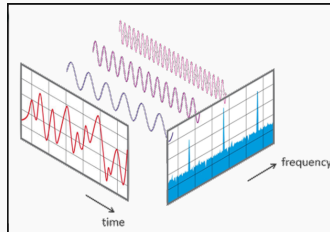
$$F_{j,k} = e^{-\frac{2\pi i}{n} j \cdot k},$$

where $i = \sqrt{-1}$. Recall $e^{-\frac{2\pi i}{n} j \cdot k} = \cos(2\pi jk/n) - i \sin(2\pi jk/n)$.



PSEUDORANDOM RIP MATRICES

In many applications can compute measurements of the form $\mathbf{Ax} = \mathbf{SFx}$, where \mathbf{F} is the Discrete Fourier Transform matrix (what an FFT computes) and \mathbf{S} is a subsampling matrix.

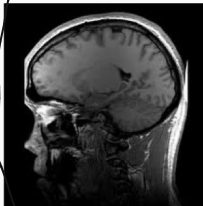


\mathbf{F} decomposes \mathbf{x} into different frequencies: $[\mathbf{Fx}]_j$ is the component with frequency j/n .

If $\mathbf{A} = \mathbf{SF}$ is a subset of rows from \mathbf{F} , then \mathbf{Ax} is a subset of random frequency components from \mathbf{x} 's discrete Fourier transform.

In many scientific applications, we can collect entries of \mathbf{Fx} one at a time for some unobserved data vector \mathbf{x} .

Warning: very cartoonish explanation of very complex problem.
Medical Imaging (MRI)

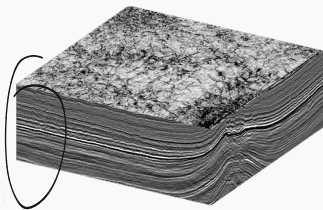


How do we measure entries of Fourier transform F_x ? Blast the body with sound waves of varying frequency.

- Especially important when trying to capture something moving (e.g. lungs, baby, child who can't sit still).
- Can also cut down on high power requirements.

Warning: very cartoonish explanation of very complex problem.

Understanding what material is beneath the crust:



Vibrate the earth at different frequencies! And measure the response.



Vibroseis Truck

Can also use airguns, controlled explosions, vibrations from drilling, etc. The fewer measurements we need from F_x , the cheaper and faster our data acquisition process becomes.

RESTRICTED ISOMETRY PROPERTY

$$k \log(n/k)$$

Setting \mathbf{A} to contain a random $m \sim O\left(\frac{k \log^2 k \log n}{\epsilon^2}\right)$ rows of the discrete Fourier matrix \mathbf{F} yields a matrix that with high probability satisfies (k, ϵ) -RIP. [Haviv, Regev, 2016].

Improves on a long line of work: Candès, Tao, Rudelson, Vershynin, Cheraghchi, Guruswami, Velingker, Bourgain.

Proving this requires similar tools to analyzing subsampled Hadamard transforms!

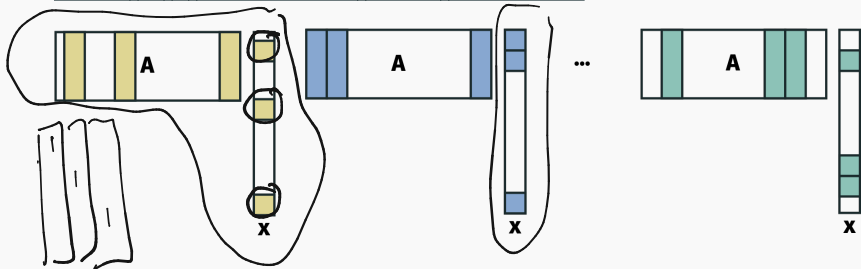
RESTRICTED ISOMETRY PROPERTY

Definition ((q, ϵ)-Restricted Isometry Property – Candes, Tao '05)

A matrix \mathbf{A} satisfies (q, ϵ)-RIP if, for all \mathbf{x} with $\|\mathbf{x}\|_0 \leq q$,

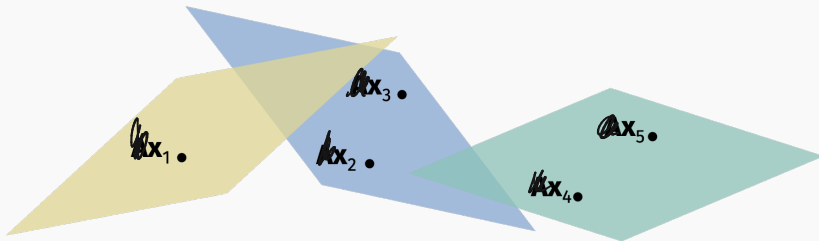
$$(1 - \epsilon)\|\mathbf{x}\|_2^2 \leq \|\mathbf{A}\mathbf{x}\|_2^2 \leq (1 + \epsilon)\|\mathbf{x}\|_2^2. \quad O(q) \text{ (4/4)}$$

The vectors that can be written as $\mathbf{A}\mathbf{x}$ for q sparse \mathbf{x} lie in a union of q dimensional linear subspaces:



RESTRICTED ISOMETRY PROPERTY

Candes ^{Tao} 2005: A random JL matrix with $O(q \log(n/q)/\epsilon^2)$ rows satisfies (q, ϵ) -RIP with high probability.



Any ideas for how you might prove this? I.e. prove that a random matrix preserves the norm of every x in this union of subspaces?

RESTRICTED ISOMETRY PROPERTY FROM JL

Theorem (Subspace Embedding from JL)

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

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

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



Quick argument:

$$\binom{n}{q} \approx \frac{n^q}{q!} \approx \left(\frac{n}{q}\right)^q$$

$$\frac{q + \log\left(\frac{n}{q}\right)^q}{\epsilon^2} = \frac{q + q \log\left(\frac{n}{q}\right)}{\epsilon^2}$$