# CS-GY 6763: Lecture 12
# Fast Johnson-Lindenstrauss Transform, Sparse Recovery and Compressed Sensing

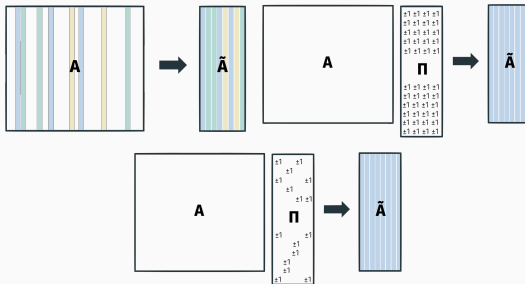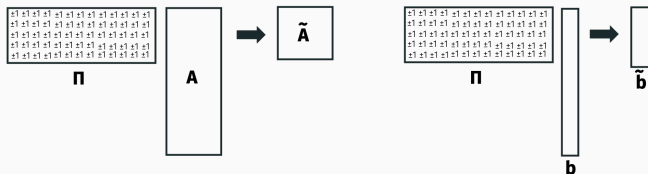NYU Tandon School of Engineering, Prof. Christopher Musco

**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 $A$.

Randomized approximate regression using a
Johnson-Lindenstrauss Matrix:



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

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

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

### Theorem (Randomized Linear Regression)

*Let $\boldsymbol{\Pi}$ be a properly scaled JL matrix (random Gaussian, sign, sparse random, etc.) with $m = \tilde{O}\left(\frac{d}{\epsilon^2}\right)$ rows. Then with probability $(1 - \delta)$, for any $\mathsf{A} \in \mathbb{R}^{n \times d}$ and $\mathsf{b} \in \mathbb{R}^n$,*

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

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

- Showed that for <u>all</u> x, $\|A\tilde{x} - b\|_2^2 = (1 \pm \epsilon) \min_x \|Ax - b\|_2^2$
- Easy to prove for a single x using JL lemma.
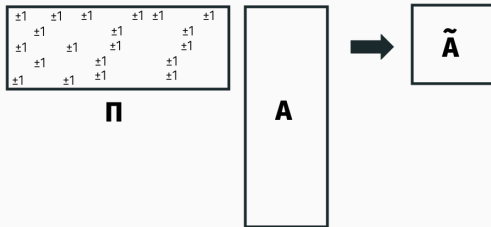- To extend to all x (an finite set) used an $\epsilon$-net argument.

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

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

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

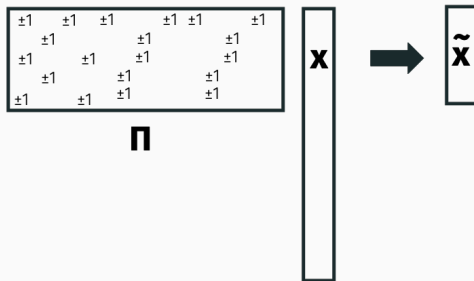Goal: Develop faster Johnson-Lindenstrauss projections.



Typically using <u>sparse</u> or <u>structured</u> matrices instead of fully random JL matrices.

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

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

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



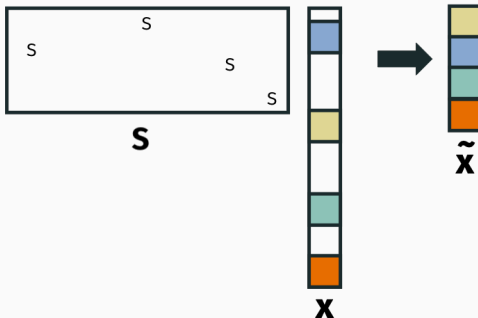Recall that once the bound above is proven, linearity lets use preserve things like $\|y - z\|_2^2$ or $\|Aw - b\|_2^2$.

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.



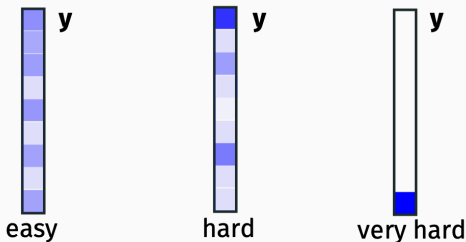If we take $m$ samples, $\tilde{x}$ can be computed in $O(m)$ time. Woohoo!

What is the problem with this approach?

Uniform sampling only works well if $y = Ax$ is "flat".



easy      hard      very hard

### 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.

10

Subsampled Randomized Hadamard Transform[1] (SHRT)
(Ailon-Chazelle, 2006)

## Theorem (The Fast JL Lemma)

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

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

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

Very little loss in embedding dimension compared to standard JL.
Leverages the simple sampling result from above.

---

[1]One of my top 3 favorite randomized algorithms.

Key idea: First multiply x by a "mixing matrix" M which ensures it cannot be too concentrated in one place.

M will have the properties that

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

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

$$\Pi x = SMx$$

Good mixing matrices should look random:

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

$$\textbf{M} \qquad\qquad \textbf{x}$$

I claim to mix any **x** with high probability, **M** underline{needs} to be chosen randomly. Why?

Recall that $\|\mathbf{M}\mathbf{x}\|_2 = \|\mathbf{x}\|_2$, so **M** is orthogonal.

Good mixing matrices should look random:



$$\begin{array}{cccccccc}
+1 & -1 & +1 & +1 & +1 & -1 & +1 & -1 \\
-1 & -1 & -1 & +1 & +1 & +1 & -1 & -1 \\
+1 & -1 & +1 & +1 & +1 & -1 & -1 & -1 \\
+1 & +1 & +1 & +1 & -1 & +1 & -1 & +1 \\
-1 & -1 & +1 & +1 & -1 & +1 & +1 & -1 \\
-1 & +1 & -1 & -1 & -1 & +1 & -1 & -1 \\
-1 & +1 & -1 & +1 & -1 & -1 & -1 & +1
\end{array}$$

**M**           **x**

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

Subsampled Randomized Hadamard Transform

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

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

The Hadarmard 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.

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

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

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

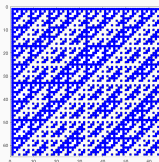The $n \times n$ Hadamard matrix has all entries as $\pm\frac{1}{\sqrt{n}}$.

**Property 1**: For any $k = 0, 1, \ldots$, we have $\|H_k v\|_2^2 = \|v\|_2^2$ for all $v$. I.e., $H_k$ is orthogonal.

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

**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.

19

**Lemma (SHRT mixing lemma)**

*Let $H$ be an $(n \times n)$ Hadamard matrix and $D$ a random $\pm 1$ 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:

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

Our main results then follows directly from our sampling result from earlier:

### 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 $(1 - \delta)$.*

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

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

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

where $D_1, \ldots, D_n$ are random $\pm 1$'s.

This is equivalent to

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

where $R_1, \ldots, R_n$ are random $\pm 1$'s.

22

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

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

$z_i$ 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:

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

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

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

.

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

Could use Hoeffding or Bernstein inequality, or a shift, need to use Bernstein type concentration inequality to prove the bound:

### Lemma (Rademacher Concentration)

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

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

This is call the <u>Khintchine Inequality</u>. It is specialized to sums of scaled $\pm 1$'s, and is a bit tighter and easier to apply than using a generic Bernstein bound.

Recall that $z_i = h_i^T Dx = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} R_i x_i$.

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

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

As shown earlier, we can thus guarantee that:

$$(1 - \epsilon)\|z\|_2^2 \leq \|Sz\|_2^2 \leq (1 + \epsilon)\|z\|_2^2$$

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

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

$\|Sz\|_2^2 = \|SHDx\|_2^2 = \|\Pi x\|_2^2$ and $\|z\|_2^2 = \|x\|_2^2$, so we are done.

**Upshot for regression:** Compute $\mathbf{\Pi A}$ 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 A when A is dense.

Clarkson-Woodruff 2013, STOC Best Paper: Let $O(\text{nnz}(A))$ be the number of non-zeros in A. It is possible to compute Ã with poly($d$) rows in:

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

Π is chosen to be an ultra-sparse random matrix. Uses totally different techniques (you can't do JL + $\epsilon$-net). Related to Danrong's reading group presentation.

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

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

Simple, inspired algorithm that has been used for accelerating:

- Vector dimensionality reduction
- Linear algebra
- Locality sensitive hashing (SimHash)
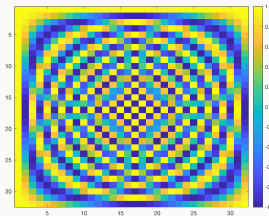- Randomized kernel learning methods.

```
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);
```

The Hadamard Transform is closely related to the Discrete Fourier Transform.
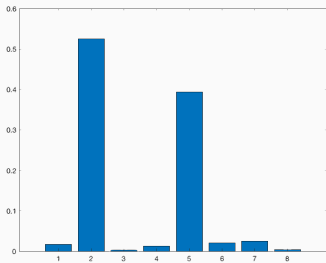
$$\mathsf{F}_{j,k} = e^{-2\pi i \frac{j \cdot k}{n}}, \qquad\qquad \mathsf{F}^*\mathsf{F} = \mathsf{I}.$$



Real part of $\mathsf{F}_{j,k}$.

$\mathsf{F}\mathsf{y}$ computes the Discrete Fourier Transform of the vector $\mathsf{y}$. 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 it's Fourier transform cannot both be concentrated.



Vector **y**.



Fourier transform **Fy**.

What do we know?
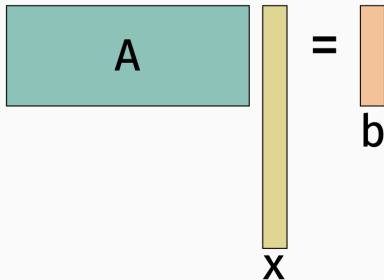
Sampling does not preserve norms, i.e. $\|\mathsf{Sy}\|_2 \not\approx \|\mathsf{y}\|_2$ when $\mathsf{y}$ has a few large entries.

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

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

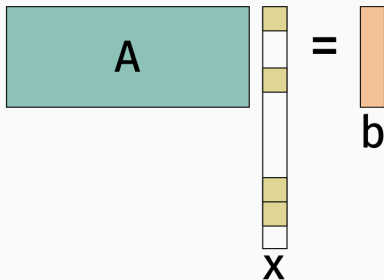**Underdetermined linear regression:** Given $A \in \mathbb{R}^{m \times n}$ with $m < n$, $b \in \mathbb{R}^m$. Assume $b = Ax$ for some $x \in \mathbb{R}^n$.



- Infinite possible solutions $y$ to $Ay = b$, so in general, it is impossible to recover parameter vector $x$ from the data $A$, $b$.

**Underdetermined linear regression:** Given $A \in \mathbb{R}^{m \times n}$ with $m < n$, $b \in \mathbb{R}^m$. Solve $Ax = b$ for $x$.

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



- In many cases can recover $x$ with $\ll n$ rows. In fact, often $\sim O(k)$ suffice.
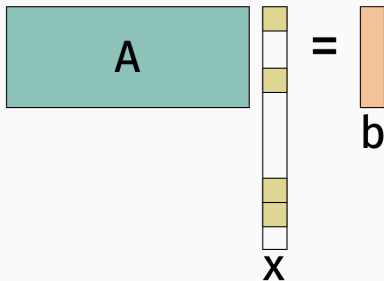- Need additional assumptions about $A$!

- In statistics and machine learning, we often think about $A$'s rows as data drawn from some universe/distribution:

| | bedrooms | bathrooms | sq.ft. | floors | list price | sale price |
|---|---|---|---|---|---|---|
| home 1 | 2 | 2 | 1800 | 2 | 200,000 | 195,000 |
| home 2 | 4 | 2.5 | 2700 | 1 | 300,000 | 310,000 |
| . | . | . | . | . | . | . |
| . | . | . | . | . | . | . |
| . | . | . | . | . | . | . |
| home n | 5 | 3.5 | 3600 | 3 | 450,000 | 450,000 |

- In other settings, we will get to <u>choose</u> $A$'s rows. I.e. each $b_i = \mathbf{x}^T \mathbf{a}_i$ for some vector $\mathbf{a}_i$ that we select.

- In the later case, we often call $b_i$ a <u>linear measurement</u> of $\mathbf{x}$ and we call $A$ a measurement matrix.

When should this problem be difficult?

Many ways to formalize our intuition

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

Definition (($q, \epsilon$)-Restricted Isometry Property)

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

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

- Johnson-Lindenstrauss type condition.
- **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 **A** with $\sim O(q \log(n/q))$ rows satisfies RIP.

Theorem ($\ell_0$-minimization)

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

$$\min\|z\|_0 \qquad subject\ to \qquad Az = b.$$

- Establishes that <u>information theoretically</u> we can recover $x$. Solving the $\ell_0$-minimization problem is computationally difficult, requiring $O(n^k)$ time. We will address faster recovery shortly.

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

Proof: By contradiction, assume there is some $y \neq x$ such that $Ay = b$, $\|y\|_0 \leq \|x\|_0$.

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

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

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

  or even

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

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**.
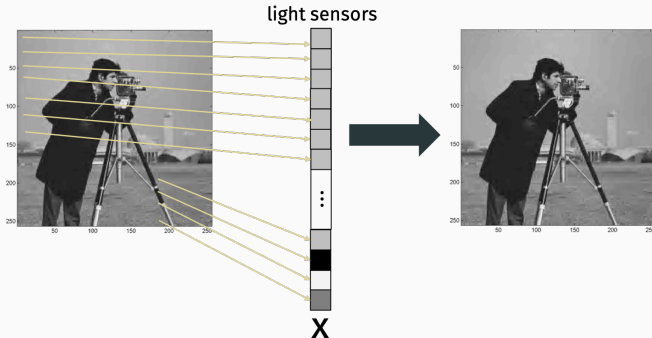
What matrices satisfy this property?

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

Some real world data may look random, but this is also a useful observation algorithmically when we want to design A.
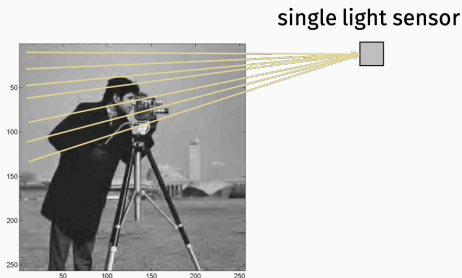
Typical acquisition of image by camera:



Requires one image sensor per pixel captured.

Compressed acquisition of image:
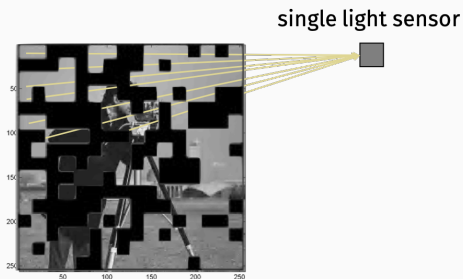
single light sensor



$$p = \sum_{i=1} x_i = \begin{bmatrix} \frac{1}{n} & \frac{1}{n} & \cdots & \frac{1}{n} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

Does not provide very much information about the image.

46

But several random linear measurements do!

single light sensor



$$p = \sum_{i=1} R_i x_i = \begin{bmatrix} 0 & 1 & 0 & 0 \dots & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

Applications in:

- Imaging outside of the visible spectrum (more expensive sensors).
- Microscopy.
- Other scientific imaging.

Compressed sensing theory does not exactly describe these problems, but has been very valuable in modeling them.
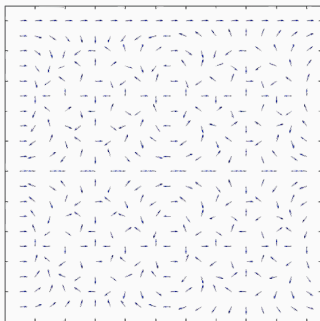
## THE DISCRETE FOURIER MATRIX
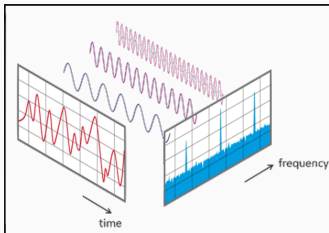
The $n \times n$ discrete Fourier matrix **F** is defined:

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

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

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



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

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

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

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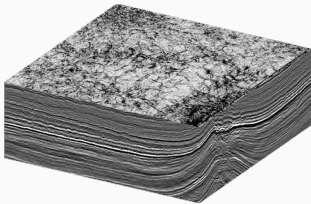
Medical Imaging (MRI)



How do we measure entries of Fourier transform **Fx**? Blast the body with sounds 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 explorations, vibrations from drilling, etc. The fewer measurements we need from Fx, the cheaper and faster our data acquisition process becomes.

Setting **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 **F** yields a matrix that with high probability satisfies $(k, \epsilon)$-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!
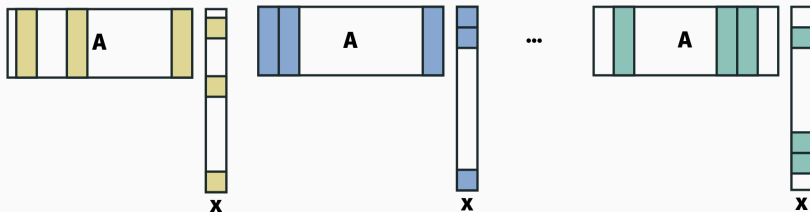
Definition ($(q, \epsilon)$-Restricted Isometry Property – Candes, Tao '05)

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

$$(1 - \epsilon)\|\mathbf{x}\|_2^2 \leq \|\mathbf{A}\mathbf{x}\|_2^2 \leq (1 + \epsilon)\|\mathbf{x}\|_2^2.$$
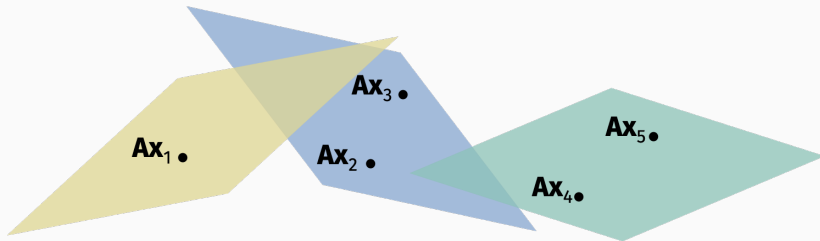
The vectors that can be written as **Ax** for $q$ sparse **x** lie in a union of $q$ dimensional linear subspaces:

**Candes, Tao 2005**: A random JL matrix with $O(q \log(n/q)/\epsilon^2)$ rows satisfies $(q, \epsilon)$-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?

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

Quick argument:

Definition (($q, \epsilon$)-Restricted Isometry Property)

A matrix $A$ satisfies ($q, \epsilon$)-RIP if, for all $x$ with $\|x\|_0 \leq q$,

$$(1 - \epsilon)\|x\|_2^2 \leq \|Ax\|_2^2 \leq (1 + \epsilon)\|x\|_2^2.$$

Lots of other random matrices satisfy RIP as well.

One major theoretical question is if we can deterministically construct good RIP matrices. Interestingly, if we want ($O(k), O(1)$) RIP, we can only do so with $O(k^2)$ rows (now very slightly better – thanks to Bourgain et al.).

Whether or not a linear dependence on $k$ is possible with a deterministic construction is unknown.

### Theorem ($\ell_0$-minimization)

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

$$\min \|z\|_0 \qquad subject\ to \qquad Az = b.$$

Algorithm question: Can we recover x using a faster method? Ideally in polynomial time.

Convex relaxation of the $\ell_0$ minimization problem:

**Problem (Basis Pursuit, i.e. $\ell_1$ minimization.)**

$$\min_z \|z\|_1 \qquad \textit{subject to} \qquad Az = b.$$

- Objective is convex.

- Optimizing over convex set.

What is one method we know for solving this problem?

Equivalent formulation:

Problem (Basis Pursuit Linear Program.)

$$\min_{w,z} \mathbf{1}^T w \qquad \textit{subject to} \qquad Az = b, w \geq 0, -w \leq z \leq w.$$

Can be solved using any algorithm for linear programming. An Interior Point Method will run in $\sim O(n^{3.5})$ time.

### Theorem

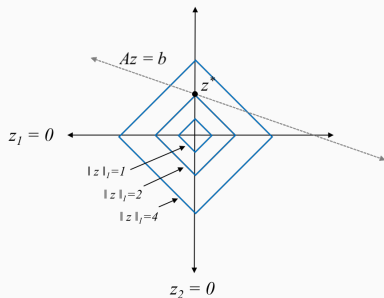*If $A$ is $(3k, \epsilon)$-RIP for $\epsilon < .17$ and $\|x\|_0 = k$, then $x$ is the unique optimal solution of the Basis Pursuit LP).*

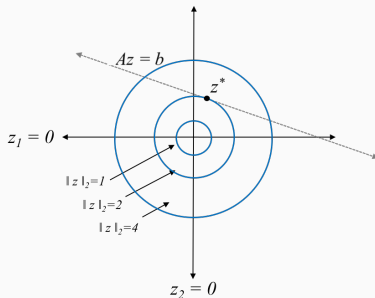### Two surprising things about this result:

- Exponentially improve computational complexity with only a <u>constant factor</u> overhead in measurement complexity.
- Typical "relax-and-round" algorithm, but rounding is not even necessary! Just return the solution of the relaxed problem.

Suppose $A$ is $2 \times 1$, so $b$ is just a scalar and $x$ is a 2-dimensional vector.



Vertices of level sets of $\ell_1$ norm correspond to sparse solutions.

This is not the case e.g. for the $\ell_2$ norm.

> **Theorem**
>
> *If $A$ is $(3k, \epsilon)$-RIP for $\epsilon < .17$ and $\|x\|_0 = k$, then $x$ is the unique optimal solution of the Basis Pursuit LP).*

Similar proof to $\ell_0$ minimization:

- By way of contradiction, assume $x$ is <u>not the optimal solution</u>. Then there exists some non-zero $\Delta$ such that:
  - $\|x + \Delta\|_1 \leq \|x\|_1$
  - $A(x + \Delta) = Ax$. I.e. $A\Delta = 0$.

Difference is that we can no longer assume that $\Delta$ is sparse.

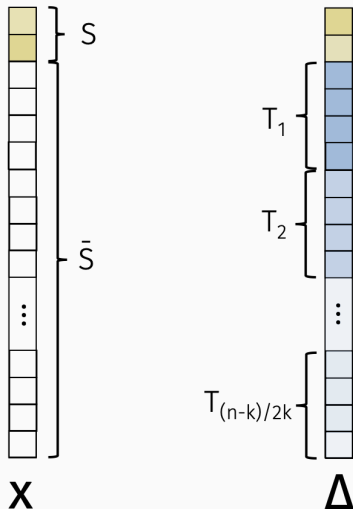<span style="color:orange">We will argue that $\Delta$ is approximately sparse.</span>

First tool:

$$\text{For any } q\text{-sparse vector } \mathbf{w}, \qquad \|\mathbf{w}\|_2 \leq \|\mathbf{w}\|_1 \leq \sqrt{q}\|\mathbf{w}\|_2$$

Second tool:

$$\text{For any norm and vectors } \mathbf{a}, \mathbf{b}, \qquad \|\mathbf{a} + \mathbf{b}\| \geq \|\mathbf{a}\| - \|\mathbf{b}\|$$

Some definitions:

Claim 1: $\|\Delta_S\|_1 \geq \|\Delta_{\bar{S}}\|_1$

**Claim 2:** $\|\Delta_S\|_2 \geq \sqrt{2} \sum_{j \geq 2} \|\mathbf{\Delta}_{T_j}\|_2$:

$$\|\mathbf{\Delta}_S\|_2 \geq \frac{1}{\sqrt{k}}\|\mathbf{\Delta}_S\|_1 \geq \frac{1}{\sqrt{k}}\|\mathbf{\Delta}_{\bar{S}}\|_1 = \frac{1}{\sqrt{k}} \sum_{j \geq 1} \|\mathbf{\Delta}_{T_j}\|_1.$$

Claim: $\|\mathbf{\Delta}_{T_j}\|_1 \geq \sqrt{2k}\|\mathbf{\Delta}_{T_{j+1}}\|_2$

Finish up proof by contradiction: Recall that $A$ is assumed to have the $(3k, \epsilon)$ RIP property.

$$0 = \|A\Delta\|_2 \geq \|A\Delta_{S \cup T_1}\|_2 - \sum_{j \geq 2} \|A\Delta_{T_j}\|_2$$

A lot of interest in developing even faster algorithms that avoid using the "heavy hammer" of linear programming and run in even faster than $O(n^{3.5})$ time.

- **Iterative Hard Thresholding**: Looks a lot like projected gradient descent. Solve $\min_z \|Az - b\|$ with gradient descent while continually projecting $z$ back to the set of $k$-sparse vectors. Runs in time $\sim O(nk \log n)$ for Gaussian measurement matrices and $O(n \log n)$ for subsampled Fourer matrices.

- Other "first order" type methods: Orthogonal Matching Pursuit, CoSaMP, Subspace Pursuit, etc.

When **A** is a subsampled Fourier matrix, there are now methods that run in $O(k \log^c n)$ time [Hassanieh, Indyk, Kapralov, Katabi, Price, Shi, etc. 2012+].
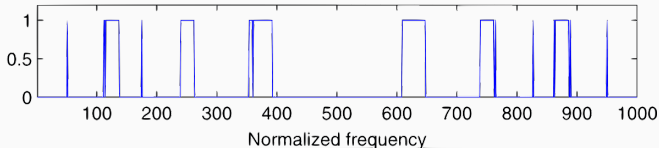
Hold up...

**Corollary:** When **x** is $k$-sparse, we can compute the inverse Fourier transform **F**\***Fx** of **Fx** in $O(k \log^c n)$ time!

- Randomly subsample **Fx**.
- Feed that input into our sparse recovery algorithm to extract **x**.

Fourier and inverse Fourier transforms in <u>sublinear time</u> when the output is sparse.



**Applications in:** Wireless communications, GPS, protein imaging, radio astronomy, etc. etc.

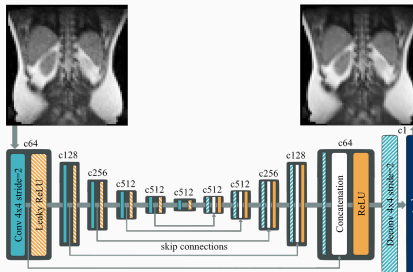## Compressed Sensing using Generative Models

Ashish Bora[*]      Ajil Jalal[†]      Eric Price[‡]      Alexandros G. Dimakis[§]

**Abstract**

The goal of compressed sensing is to estimate a vector from an underdetermined system of noisy linear measurements, by making use of prior knowledge on the structure of vectors in the relevant domain. For almost all results in this literature, the structure is represented by sparsity in a well-chosen basis. We show how to achieve guarantees similar to standard compressed sensing but without employing sparsity at all. Instead, we suppose that vectors lie near the range of a generative model $G : \mathbb{R}^k \to \mathbb{R}^n$. Our main theorem is that, if $G$ is $L$-Lipschitz, then roughly $O(k \log L)$ random Gaussian measurements suffice for an $\ell_2/\ell_2$ recovery guarantee. We demonstrate our results using generative models from published variational autoencoder and generative adversarial networks. Our method can use 5-10x fewer measurements than Lasso for the same accuracy.

Original