

# CS-GY 9223 D: Lecture 13

## Compressed Sensing + Sparse Recovery

---

NYU Tandon School of Engineering, Prof. Christopher Musco

# SPARSE RECOVERY/COMPRESSED SENSING

Euclidean into  $\ell_0$

driven mostly.

What do we know?

## BASIC PROBLEM SETUP

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

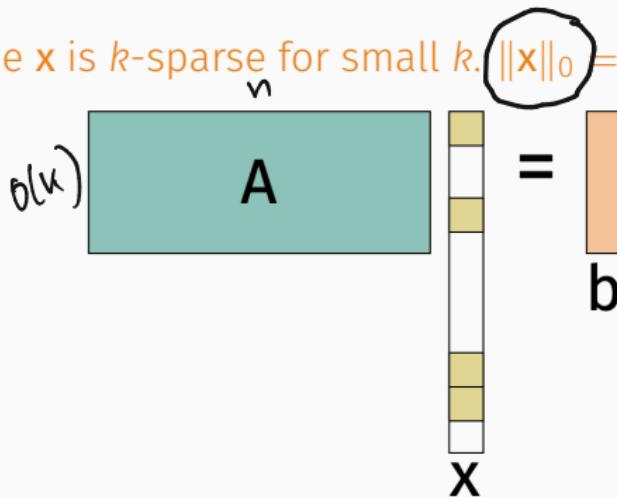
The diagram illustrates the system of equations  $Ax = b$ . On the left, there is a green box labeled 'A' representing the matrix. To its right is a vertical bracket containing two vertical vectors: 'x' (the parameter vector) and 'b' (the observation vector). An equals sign is positioned between the matrix 'A' and the vector 'b'. Below the vector 'x' is a red circle containing the letter 'x', indicating that there are multiple possible solutions for the parameter vector.

- Infinite possible solutions  $x$ . In general, impossible to recover parameter vector.

## SPARSITY RECOVERY/COMPRESSED SENSING

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  $m \ll n$  rows. In fact, often  $m = \sim O(k)$  suffice.
- Need additional (strong) assumptions about  $A$ !

## QUICK ASIDE

- In machine learning, we typically 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 many settings, we will get to choose A's rows  $\mathbf{a}_1, \dots, \mathbf{a}_m$ .  
I.e. each  $b_i = \underline{\mathbf{a}_i^T \mathbf{x}}$  for some vector  $\mathbf{a}_i$  that we select.
- In this setting, we often call  $b_i$  a linear measurement of  $\mathbf{x}$  and we call  $\mathbf{A}$  a measurement matrix.

## ASSUMPTIONS ON MEASUREMENT MATRIX

$$\Theta(n \log n) \quad n=2k$$

When should this problem be difficult?

$$\Theta(n \log n)$$

$$\begin{matrix} \Theta(n) \\ \left[ \begin{array}{c|c} +1 & +1 \\ -1 & -1 \\ -1 & -1 \\ +1 & +1 \end{array} \right] \end{matrix} \xrightarrow{n} A = \begin{pmatrix} \text{yellow} \\ \text{white} \\ \text{yellow} \\ \text{white} \end{pmatrix} = b$$

$$A \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = A \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

## ASSUMPTIONS ON MEASUREMENT MATRIX

$$A\mathbf{x} = \mathbf{f}_j \rightarrow A(\mathbf{x} - \mathbf{d}) = \mathbf{0}$$

*2k sparse*

Many ways to formalize our intuition

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

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

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

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

- 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).

## FIRST SPARSE RECOVERY RESULT

### 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  $\hat{x}$  is the unique minimizer of:

$$\begin{array}{lll} \min \|z\|_0 & \text{subject to} & Az = b. \end{array}$$

- Establishes that information theoretically we can recover  $x$ . Solving the  $\ell_0$ -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  $Ay = b$ ,  $\|y\|_0 \leq \|x\|_0$ .

$$Ax = b$$

$$\underbrace{A(y-x)}_{{\leq} 2k \text{ sparse}} = 0$$

$$\|A(y-x)\|_2 \in (1 \pm \epsilon) \|y-x\|_2$$

$\downarrow$   
must be  
 $\|y-x\|_2$

$$\underbrace{\|A(y-x)\|_2}_{{\geq} 2k \text{ sparse}} = 0 \quad \|y-x\| \neq 0$$

## ROBUSTNESS

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

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

or even

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

## ROBUSTNESS

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  $(k, \epsilon)$ -RIP.

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

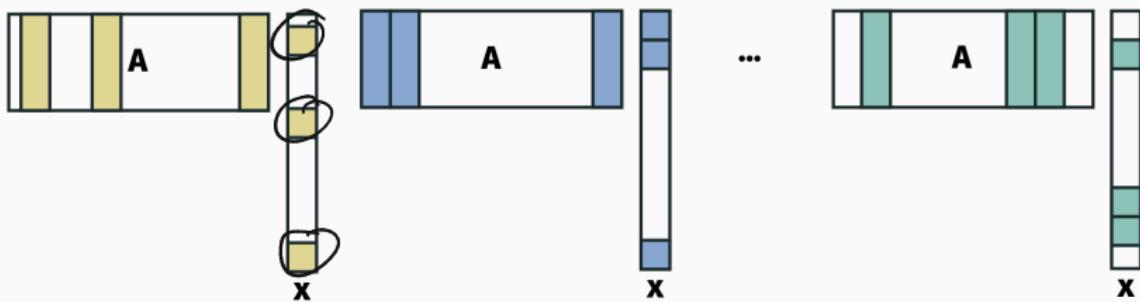
## RESTRICTED ISOMETRY PROPERTY

Definition  $(q, \epsilon)$ -Restricted Isometry Property – Candes, Tao '05)

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

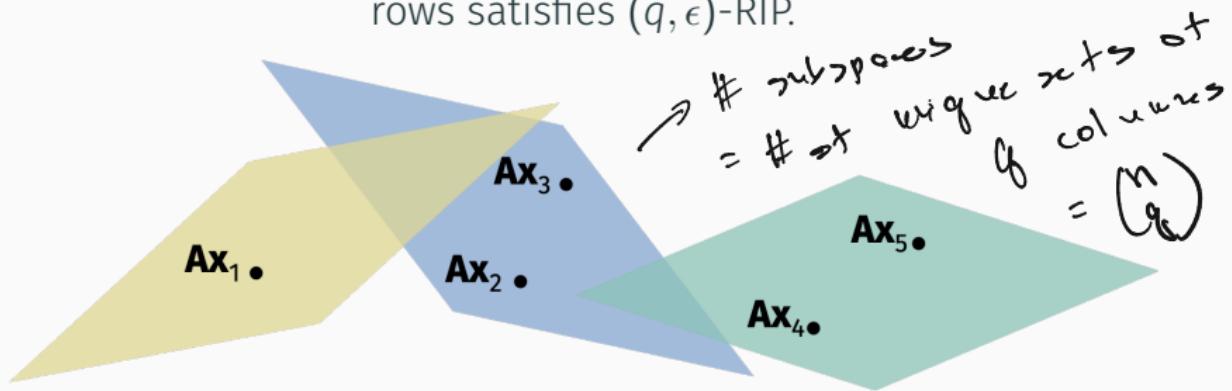
$$(1 - \epsilon)\|\mathbf{x}\|_2^2 \leq \|\mathbf{Ax}\|_2^2 \leq (1 + \epsilon)\|\mathbf{x}\|_2^2.$$

The vectors that can be written as  $\mathbf{Ax}$  for  $\mathbf{x}$  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  $(q, \epsilon)$ -RIP.



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  $A \in \mathbb{R}^{m \times n}$  is chosen from any distribution  $\mathcal{D}$  satisfying the Distributional JL Lemma, then with probability  $1 - \delta$ ,

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

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

Quick argument:

$$\log\left(\binom{n}{q}\right) \approx \log(n^q) = q \log(n)$$

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

$$O\left(\frac{q \log(n/q)}{\epsilon^2}\right)$$

## APPLICATION: HEAVY HITTERS IN DATA STREAMS

Suppose you view a stream of numbers in  $1, \dots, n$ :

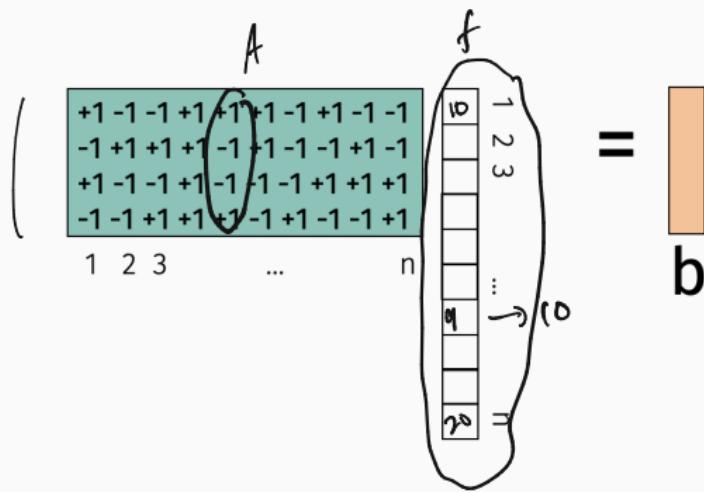
4, 18, 4, 1, 2, 24, 6, 4, 3, 18, 18, ...

After some time, you want to report which  $k$  items appeared most frequently in the stream.

E.g. Amazon is monitoring web-logs to see which product pages people view. They want to figure out which products are viewed most frequently.  $n \approx 500$  million.

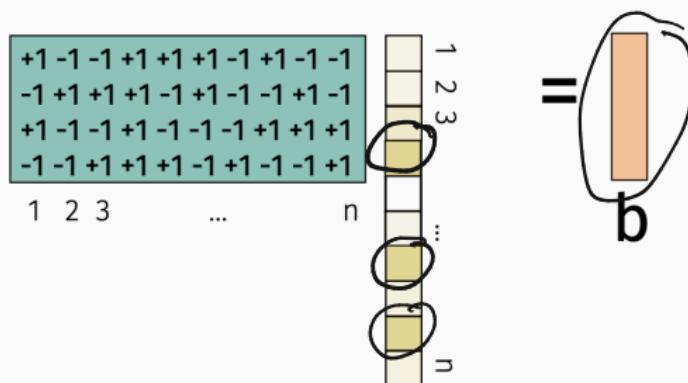
How can you do this quickly in small space?

## APPLICATION: HEAVY HITTERS IN DATA STREAMS



- Every time we receive a number  $i$  in the stream, add column  $A_i$  to  $b$ .

## APPLICATION: HEAVY HITTERS IN DATA STREAMS



- At the end  $\mathbf{b} = \mathbf{A}\mathbf{x}$  for an approximately sparse  $\mathbf{x}$  if there were only a few “heavy hitters”. Recover  $\mathbf{x}$  from  $\mathbf{b}$  using a sparse recovery method (like  $\ell_0$  minimization).

## APPLICATION: HEAVY HITTERS IN DATA STREAMS

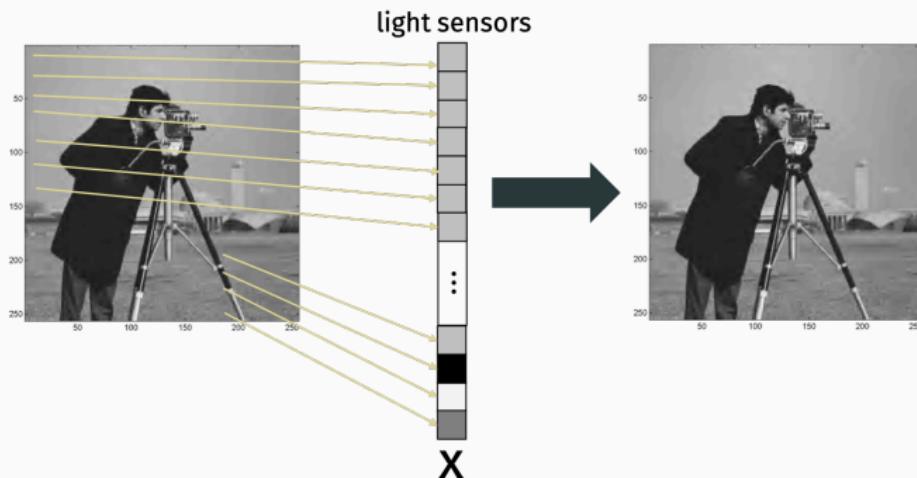
Naturally handles both insertions or deletions.

*insert(4), insert(18), remove(4), insert(1), insert(2), remove(2) ...*

E.g. Amazon is monitoring what products people add to their “wishlist” and wants a list of most tagged products. Wishlists can be changed over time, including by removing items.

## APPLICATION: SINGLE PIXEL CAMERA

Typical acquisition of image by camera:

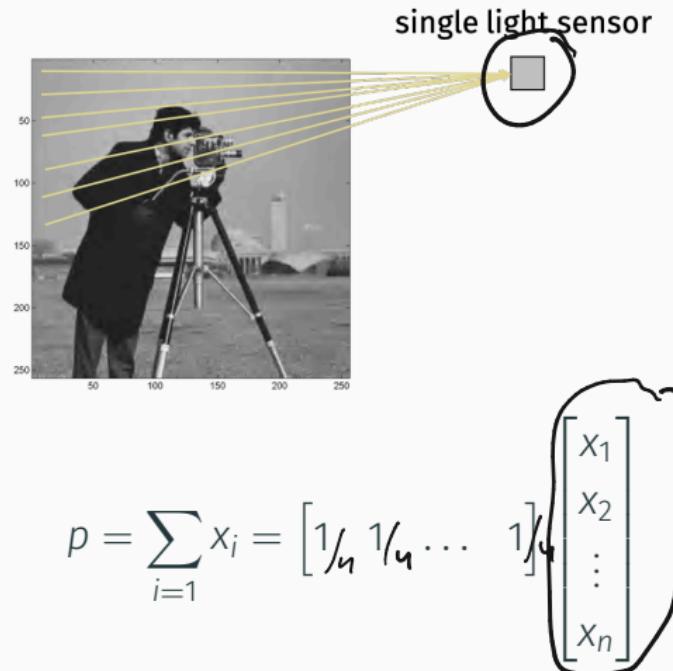


Requires one image sensor per pixel captured.

## APPLICATION: SINGLE PIXEL CAMERA

Compressed acquisition of image:

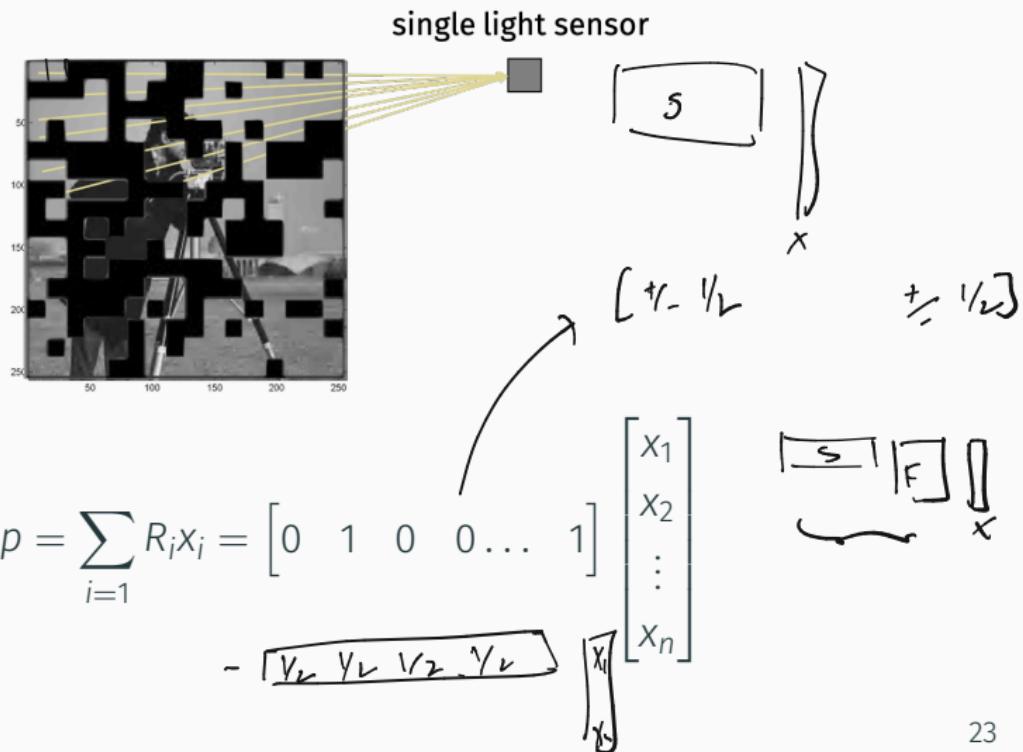
(0, 1)



Does not provide very much information about the image.

## APPLICATION: SINGLE PIXEL CAMERA

But several random linear measurements do!



## APPLICATION: SINGLE PIXEL CAMERA

Applications in:

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

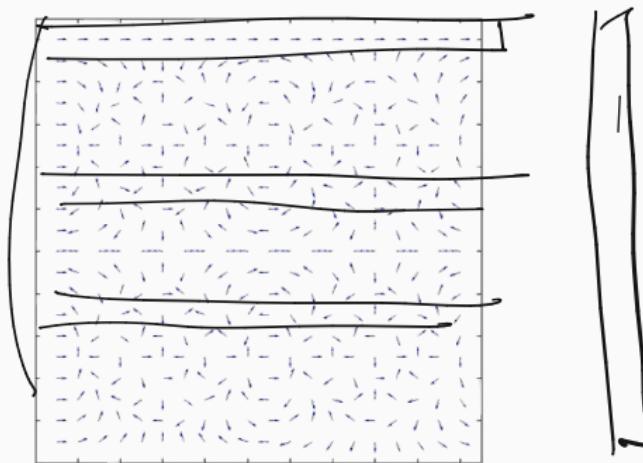
Compressed sensing theory does not exactly describe these problems, but has been very valuable in modeling them.

## THE DISCRETE FOURIER MATRIX

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

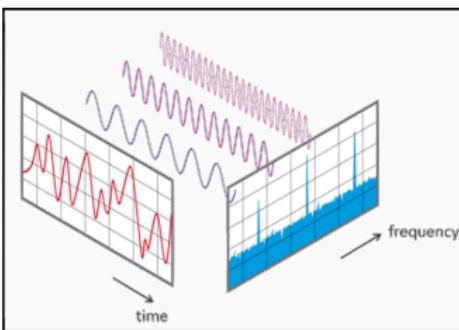
$$F_{j,k} = e^{\frac{-2\pi i}{n} j k},$$

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



## THE DISCRETE FOURIER MATRIX

$\mathbf{F}\mathbf{x}$  is the Discrete Fourier Transform of the vector  $\mathbf{x}$  (what an FFT computes).



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

Because  $\mathbf{F}^*\mathbf{F} = \mathbf{I}$ ,  $\mathbf{F}^*\mathbf{F}\mathbf{x} = \mathbf{x}$ , so we can recover  $\mathbf{x}$  if we have access to its DFT.  $\mathbf{F}\mathbf{x}$ .

## RESTRICTED ISOMETRY PROPERTY

$$(\ell_0, \|\cdot\|_1)$$

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  $(k, \epsilon)$ -RIP. [Haviv, Regev, 2016].

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

You have seen some of the tools used to prove this when we proved that a subsampled Hadamard matrix, which is a type of Fourier matrix, can be used to give a  $JL$  guarantee.

## THE DISCRETE FOURIER MATRIX

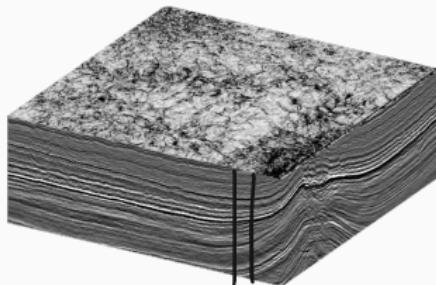
If  $A$  is a subset of  $q$  rows from  $F$ , then  $Ax$  is a subset of random frequency components from  $x$ 's discrete Fourier transform.

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

## APPLICATION: GEOPHYSICS

Warning: very cartoonish explanation of very complex problem.

Understanding what material is beneath the crust:



Think of vector  $\mathbf{x}$  as scalar values of the density/reflectivity in a single vertical core of the earth.

How do we measure entries of Fourier transform  $\mathbf{F}\mathbf{x}$ ?

## APPLICATION: GEOPHYSICS

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



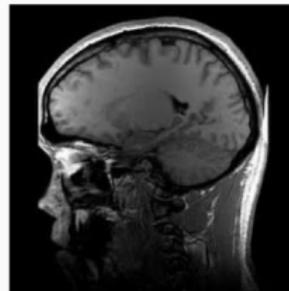
Vibroseis Truck

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.

**Killer app: Oil Exploration.**

Warning: very cartoonish explanation of very complex problem.

### Medical Imaging (MRI)



Vector  $\mathbf{x}$  here is a 2D image. Everything works with 2D Fourier transforms.

How do we measure entries of Fourier transform  $\mathbf{F}\mathbf{x}$ ?

## APPLICATION: GEOPHYSICS

Blast the body with sounds waves of varying frequency.



The fewer measurements we need from  $\mathbf{F}_x$ , the faster we can acquire an image.

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

### Definition $((q, \epsilon)\text{-Restricted Isometry Property})$

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

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

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:

$$\min \|z\|_0 \quad \text{subject to} \quad Az = b.$$

**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.)

$$\min_z \|z\|_1 \quad \text{subject to} \quad (Az = b)$$

- Objective is convex.
- Optimizing over convex set.

What is one method we know for solving this problem?

## BASIS PURSUIT LINEAR PROGRAM

Equivalent formulation:

Problem (Basis Pursuit Linear Program.)

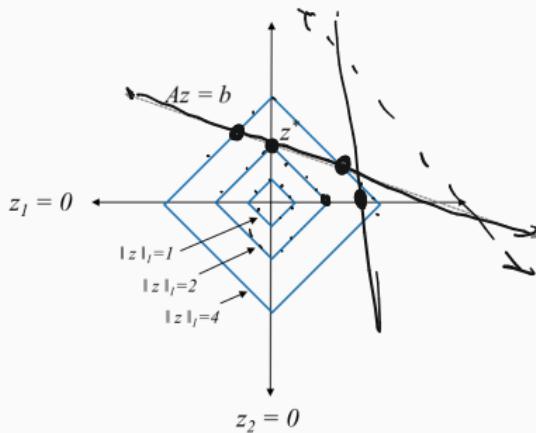
$$\min_{w,z} \cancel{1^T w} \quad \text{subject to} \quad \left( Az = b, w \geq 0, -w \leq z \leq w. \right)$$

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

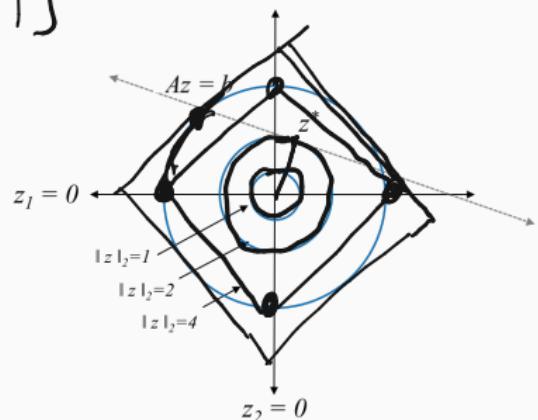
## BASIS PURSUIT INTUITION

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

$$A = \begin{bmatrix} 1 & 1 \end{bmatrix}, \quad b = \begin{bmatrix} b \end{bmatrix}$$



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**

$$b - Ax$$

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

Similar proof to  $\ell_0$  minimization:

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

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

We will argue that  $\Delta$  is approximately sparse.

## TOOLS NEEDED

First tool:

$$\sum_{i=1}^q |w_i|^2$$

For any  $q$ -sparse vector  $w$ ,

$$\left( \sum_{i=1}^q |w_i| \right)^2 =$$

$$\sum_{i,j} |w_i| |w_j|$$

$$[1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0]$$

$$\sqrt{q}$$

$$q$$

$$\|w\|_2^2 \leq \|w\|_1^2 \leq \sqrt{q} \|w\|_2$$

$$\sum_{i,j} |w_i|^2 + |w_j|^2 = q \sum_{i=1}^q |w_i|^2$$

$\leq$

Second tool:

$$a^2 + b^2 \geq 2ab$$

$$(a-b)^2 \geq 0$$

$$a^2 + b^2 - 2ab \geq 0$$

For any norm and vectors  $a, b$ ,

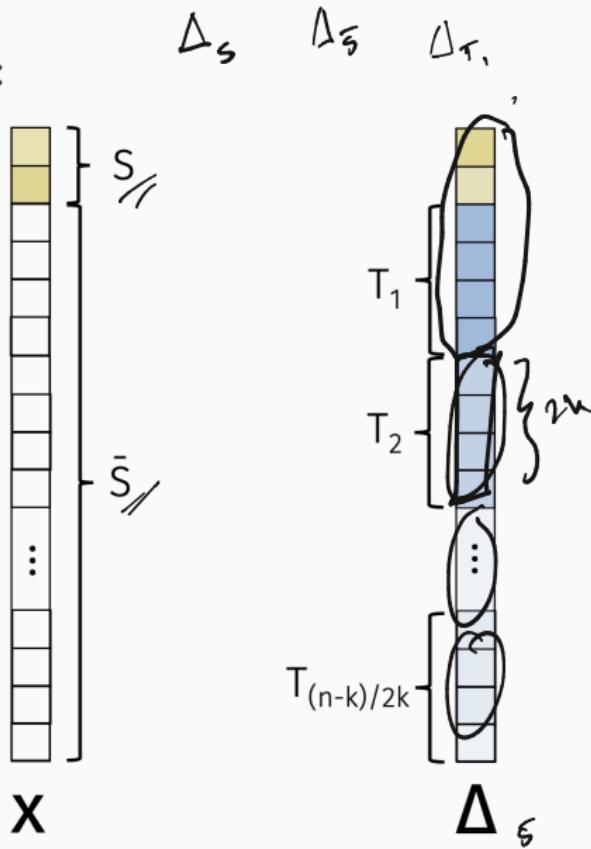
$$\|a + b\| \geq \underbrace{\|a\| - \|b\|}_{\text{ }}$$

$$\|a + b - b\| \leq \|a + b\| + \|b\|$$

$$\|a\| \leq \|a + b\| + \|b\|$$

## BASIS PURSUIT ANALYSIS

Some definitions:



## BASIS PURSUIT ANALYSIS

Claim 1:  $\|\Delta_S\|_1 \geq \|\Delta_{\bar{S}}\|_1$        $\|x + \Delta\|_1 \leq \|x\|_1$

$\|x + \Delta_S\|_1 + \|\Delta_{\bar{S}}\|_1 \leq \|x\|_1$

$(\cancel{\|x\|_1} - \|\Delta_S\|_1) + \|\Delta_{\bar{S}}\|_1 \leq \cancel{\|x\|_1}$

$\|\Delta_{\bar{S}}\|_1 - \|\Delta_S\|_1 \leq 0$

## BASIS PURSUIT ANALYSIS

Claim 2:  $\|\Delta_S\|_2 \geq \sqrt{2} \sum_{j \geq 2} \|T_j\|_2$ :

$$\|\Delta_S\|_2 \geq \frac{1}{\sqrt{k}} \|\Delta_S\|_1 \geq \frac{1}{\sqrt{n}} \|\Delta_{\bar{S}}\|_1 = \frac{1}{\sqrt{n}} \left( \sum_{j \geq 1} \|\Delta_{T_j}\|_1 \right)$$

$$\|\Delta_{T_j}\|_1 \geq \sqrt{2k} \|\Delta_{T_{j+1}}\|_2 \quad \text{where } \max(|\Delta_{T_{j+1}}|)$$

$$\begin{matrix} \downarrow \\ 2k \cdot 2 \end{matrix} \quad \sqrt{2k \cdot 2^2} = \sqrt{2n} \cdot 2$$

$$\|\Delta_S\|_2 \geq \frac{1}{\sqrt{n}} \sum_{j \geq 1} \sqrt{2n} \|\Delta_{T_{j+1}}\|_2 \geq \sqrt{2} \sum_{j \geq 2} \|T_j\|_2$$

## BASIS PURSUIT ANALYSIS

Finish up proof by contradiction:

$$\begin{aligned}
 0 &= \|A\Delta\|_2 \geq \underbrace{\|A\Delta_{S \cup T_i}\|_2}_{\geq (1-\varepsilon)\|\Delta_{S \cup T_i}\|_2} - \sum_{j \in T_i} (1+\varepsilon) \|A\Delta_{T_j}\|_2 \\
 &\geq (1-\varepsilon)\|\Delta_{S \cup T_i}\|_2 - \underbrace{\sum_{j \in T_i} (1+\varepsilon) \|A\Delta_{T_j}\|_2}_{\geq (1-\varepsilon)\|\Delta_S\|_2 - (1+\varepsilon) \frac{1}{\sqrt{2}} \|\Delta_S\|_2} \\
 &= \underbrace{\left( (1-\varepsilon) - (1+\varepsilon) \frac{1}{\sqrt{2}} \right)}_{\downarrow} \circled{ \|\Delta_S\|_2 } \\
 \varepsilon &= .17
 \end{aligned}$$

## FASTER METHODS

A lot of interest in developing even faster algorithms that avoid using the “heavy hammer” of linear programming and run in even faster than  $\underline{O(n^{3.5})}$  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  $\sim \underline{O(nk \log n)}$  for Gaussian measurement matrices and  $O(n \log n)$  for subsampled Fourier matrices.
- Other “first order” type methods: Orthogonal Matching Pursuit, CoSaMP, Subspace Pursuit, etc.

## FASTER METHODS

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

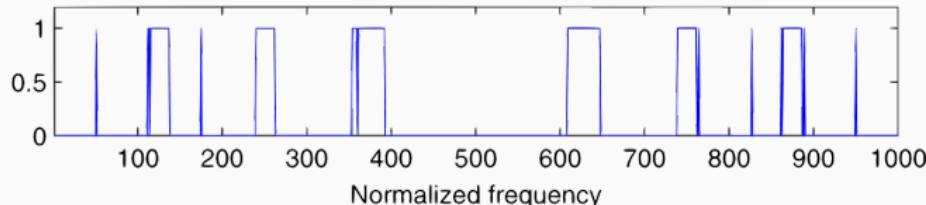
Hold up...

## (SPARSE FOURIER TRANSFORM)

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

- Randomly subsample  $\mathbf{F}x$ .
- Feed that input into our sparse recovery algorithm to extract  $x$ .

Fourier and inverse Fourier transforms in sublinear time when the output is sparse.



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