

CS-GY 9223 I: Lecture 13

Compressed Sensing + Sparse Recovery

NYU Tandon School of Engineering, Prof. Christopher Musco

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 underdetermined linear regression problem. It shows the matrix equation $Ax = b$, the minimization of the L2 norm of the residual $\|Ax - b\|_2$, and the dimensions of the matrices and vectors.

On the left, a handwritten sketch shows a matrix A (represented by a vertical rectangle) multiplied by a vector x (represented by a vertical rectangle) equals a vector b (represented by a vertical rectangle).

In the center, the handwritten text reads: $\min \|Ax - b\|_2 = 0$.

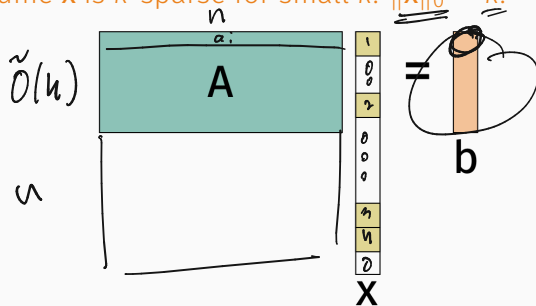
On the right, a diagram shows a matrix A (represented by a horizontal rectangle) multiplied by a vector x (represented by a vertical rectangle) equals a vector b (represented by a vertical rectangle).

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

QUICK ASIDE

- In the past, we have thought about \mathbf{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 \mathbf{A} 's rows. I.e. each $b_i = \mathbf{x}^T \mathbf{a}_i$ 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

When should this problem be difficult?

Handwritten diagrams illustrating the Strassen matrix multiplication algorithm. The top part shows the recursive splitting of matrices A and B into 2×2 blocks of size $n/2$. Matrix A is split into A_1, A_2, A_3 and A_4, A_5, A_6 . Matrix B is split into B_1, B_2, B_3 and B_4, B_5, B_6 . The equation $A \cdot B = C$ is shown. The bottom part shows the calculation of the seven products P_1 through P_7 using the recursive formula. For example, $P_1 = A_1 \cdot B_1$, $P_2 = (A_1 + A_2) \cdot (B_3 + B_6)$, etc. The final result C is then calculated from these products.

$$b = \underline{c} \lambda_i \text{ for some column } \lambda_i \in A$$

Many ways to formalize our intuition

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

RESTRICTED ISOMETRY PROPERTY

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

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

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

$q = \text{integer} \leq n$

If q is 2k sparse

- Johnson-Lindenstrauss type condition.

$$\|A\|_F^2 = 100$$
$$\|x\|_2 = 0$$

- 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 (ℓ_0 -minimization)

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 $(\underline{2k}, \underline{\epsilon})$ -RIP for any $\epsilon < 1$ then \mathbf{x} is the unique minimizer of:

$$\min \|\mathbf{z}\|_0 \quad \text{subject to} \quad \underline{\mathbf{A}\mathbf{z} = \mathbf{b}.}$$

- 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 later in the lecture.

$$\mathbf{z}^* = \mathbf{x}$$

FIRST SPARSE RECOVERY RESULT

Proof: By contradiction:

Assume by way of contradiction that $z^* \neq x$.

$$\|z^*\|_0 \leq \|x\|_0 = k.$$

Consider $(z^* - x) \rightarrow 2k$ -sparse vector.

$$Az^* = b \quad Ax = b$$

$$\|A(z^* - x)\|_2 = 0$$

$$\|A(z^* - x)\|_2 \geq \underbrace{(1 - \epsilon)}_{2k\text{-sparse}} \underbrace{\|z^* - x\|_2}_{2k\text{-RIP}} \neq 0$$

contradiction
if $\epsilon < 1$.

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 $\underline{\mathbf{b}} = \mathbf{A}(\underline{\mathbf{x}} + \mathbf{e})$ 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

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

Suppose \mathbf{A} is
 $(O(n), 1/2)$ -RIP

We will not discuss robustness in detail, but 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 $(O(k), \epsilon)$ -RIP.

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

APPLICATION: HEAVY HITTERS IN DATA STREAMS

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

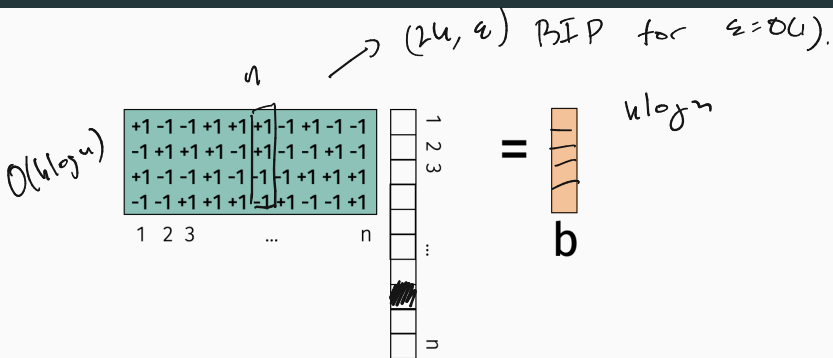
4, 18, 4, 1, 2, 24, 6, 4, 3, 18, 18, ...

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 .

$$\text{update}(i) : b = b + A_i$$

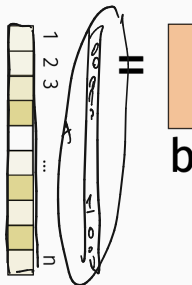
$$\text{remove}(i) : b = b - A_i$$

APPLICATION: HEAVY HITTERS IN DATA STREAMS

$n \cdot \log n$ random numbers

$$A = \begin{bmatrix} +1 & -1 & -1 & +1 & +1 & +1 & -1 & +1 & -1 & -1 \\ -1 & +1 & +1 & +1 & -1 & +1 & -1 & -1 & +1 & -1 \\ +1 & -1 & -1 & +1 & -1 & -1 & -1 & +1 & +1 & +1 \\ -1 & -1 & +1 & +1 & +1 & -1 & +1 & -1 & -1 & +1 \end{bmatrix}$$

1 2 3 ... n



Random Hash

from $\{1, \dots, n \log n\} \rightarrow \{-1, 1\}$

- 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 ℓ_0 minimization).

\mathbf{b}

→ recover vector \mathbf{x} .

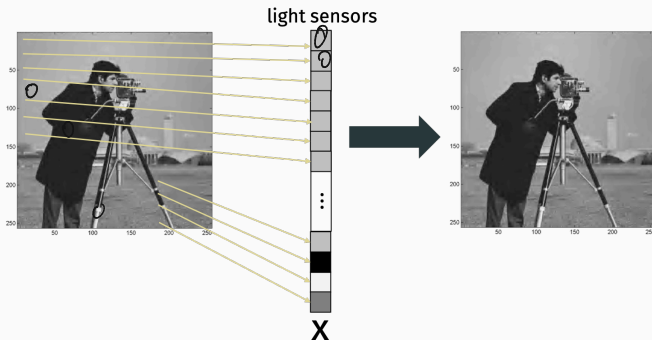
How about when there are insertions or deletions?

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

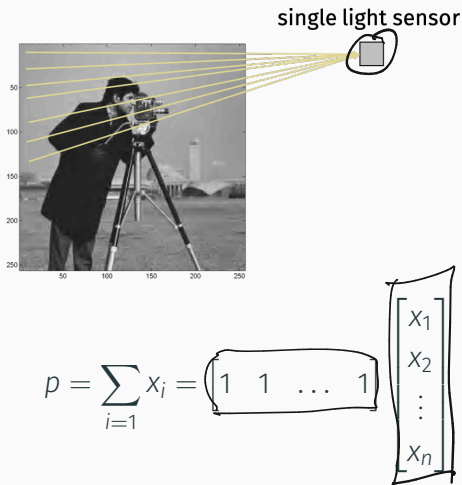
Typical acquisition of image by camera:



Requires one image sensor per pixel captured.

APPLICATION: SINGLE PIXEL CAMERA

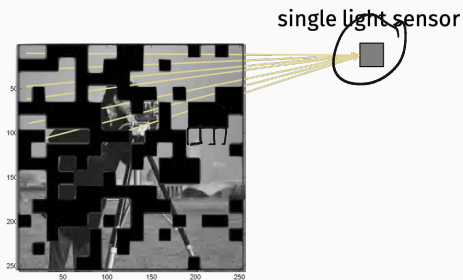
Compressed acquisition of image:



Does not provide very much information about the image.

APPLICATION: SINGLE PIXEL CAMERA

But several random linear measurements do!



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

Applications in:

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

Compressed sensing theory does not exactly describe the problem, but has been very valuable in modeling it.

RESTRICTED ISOMETRY PROPERTY

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

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

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

Uniformly subsampled Fourier matrices with
 $m \sim O\left(\frac{k \log^2 k \log n}{\epsilon^2}\right)$ rows $(O(k), \epsilon)$ -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 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.

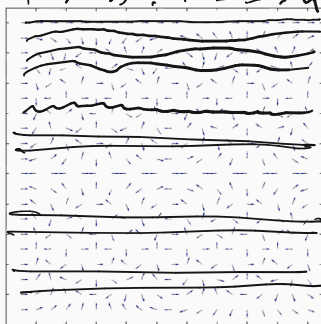
THE DISCRETE FOURIER MATRIX

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

$$F_{j,k} = e^{\frac{-2\pi i}{n} j \cdot k} \quad i = \sqrt{-1}$$

Recall that $e^{\frac{-2\pi i}{n} j \cdot k} = \underbrace{\cos(2\pi jk/n)}_{1 \ 2 \ 3 \ \dots \ n} - i \underbrace{\sin(2\pi jk/n)}_{1 \ 2 \ 3 \ \dots \ n}$.

F_X → Discrete Fourier Transform
 what FFT computes



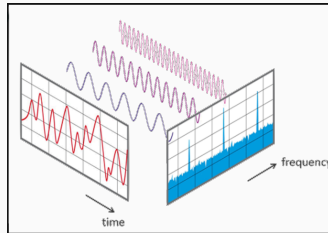
$\cos(2\pi j/n)$
 $\cos(2\pi 2j/n)$
 $\cos(2\pi 3j/n)$
 \vdots

A

Set A to contain a random $\tilde{O}(k \log n)$ rows of this matrix.

THE DISCRETE FOURIER MATRIX

Fx is the Discrete Fourier Transform of the vector x (what an FFT computes).

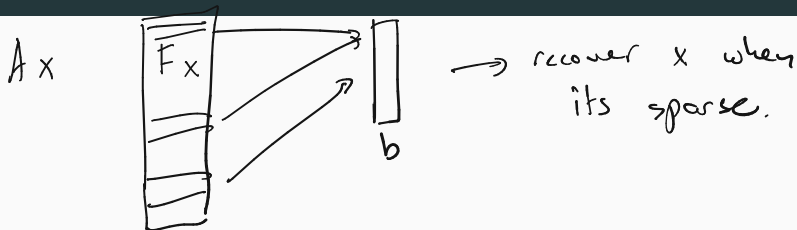


Decomposes x into different frequencies: $[F_x]_j$ is the component with frequency j/n .

Because $F^*F = I$, $F^*[F_x]$ = x , so we can recover x if we have access to its DFT. Fx .

"Inverse Discrete Fourier Transform"

THE DISCRETE FOURIER MATRIX



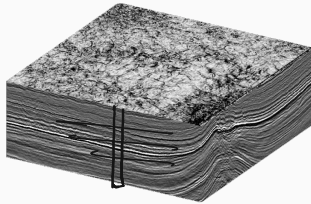
If A is a subset of q rows from F , then Ax is a subset of random frequency components from x 's discrete Fourier transform.

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 $\mathbf{F}\mathbf{x}$?

$\{F\mathbf{x}\}_j \rightarrow$ spec spec freq j .

APPLICATION: GEOPHYSICS

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



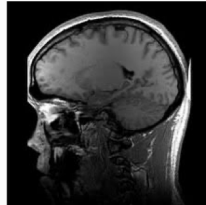
Vibroseis Truck

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

Killer app: Oil Exploration.

Warning: very cartoonish explanation of very complex problem.

Medical Imaging (MRI, ~~Ultrasound~~, etc.)



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 ~~or ultrasonic waves~~ of varying frequencies.



The fewer measurements we need from F_x , the faster we can acquire and 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

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

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

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

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 Bourgain et al.).

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

Theorem (ℓ_0 -minimization)

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} . If \mathbf{A} is $(2k, \epsilon)$ -RIP for any $\epsilon < 1$ then \mathbf{x} is the unique minimizer of:

$O(k \log n)$ rows

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

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

Convex relaxation of the ℓ_0 minimization problem: $b = Ax$

Problem (Basis Pursuit, i.e. ℓ_1 minimization.)

$$\min_z \|z\|_1$$

subject to

$$\underline{Az} = b.$$

• Objective is convex: $f(z) = \|z\|_1$

$$f(\lambda x + (1-\lambda)y) = \|\lambda x + (1-\lambda)y\|_1$$

• Optimizing over convex set:

$$\leq \lambda \|x\|_1 + (1-\lambda)\|y\|_1$$

$$\text{If } Ax_1 = b \text{ and } Ax_2 = b$$

$$= \lambda f(x) + (1-\lambda)f(y)$$

What is one method we know for solving this problem?

$$A(\lambda z_1 + (1-\lambda)z_2) = b$$

BASIS PURSUIT LINEAR PROGRAM

Equivalent formulation:

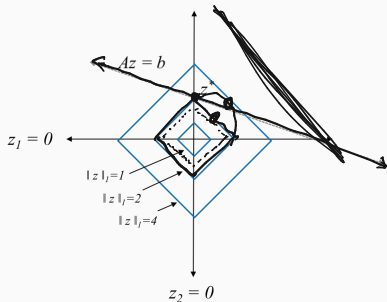
Problem (Basis Pursuit Linear Program.)

$$\min_{w,z} \mathbf{1}^T \mathbf{w} \quad \text{subject to} \quad \mathbf{Az} = \mathbf{b}, \quad -\mathbf{w} \leq \mathbf{z} \leq \mathbf{w}.$$

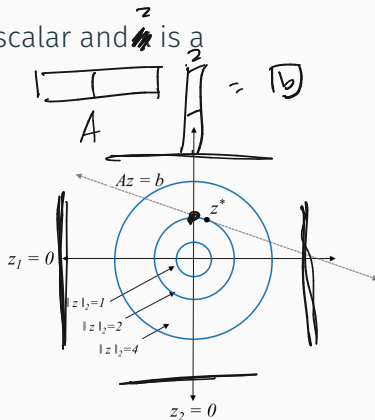
Can be solved using any algorithm for linear programming. An Interior Point Method will run in at worst $\sim O(n^{3.5})$ time.

BASIS PURSUIT INTUITION

Suppose \mathbf{A} is 1×2 , so \mathbf{b} is just a scalar and \mathbf{z} is a 2-dimensional vector.



Vertices of level sets of ℓ_1 norm correspond to sparse solutions.



This is not the case e.g. for the ℓ_2 norm.

BASIS PURSUIT ANALYSIS

(2k, ε) for ℓ_0 minimization

Theorem

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

Similar proof to ℓ_0 minimization:

- By way of contradiction, assume x is not the optimal solution. Then there exists some non-zero Δ 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 Δ is sparse.

Only one tool needed:

For any q -sparse vector w ,

$$w = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad \|w\|_1 = \|w\|_2 = 1 \quad w = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} \quad \begin{matrix} \text{is in the } k \text{ locations} \\ \|w\|_1 = 2 \quad \|w\|_2 = \sqrt{2} \end{matrix}$$

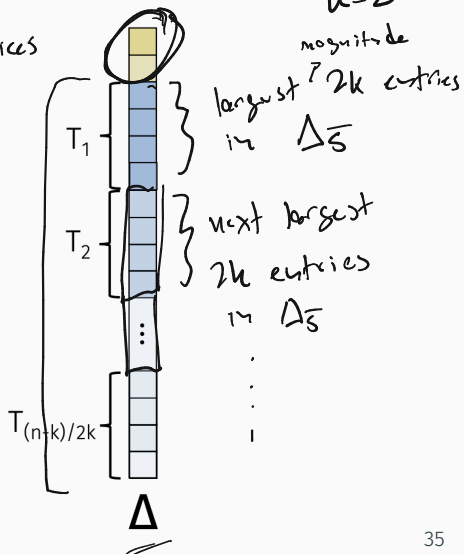
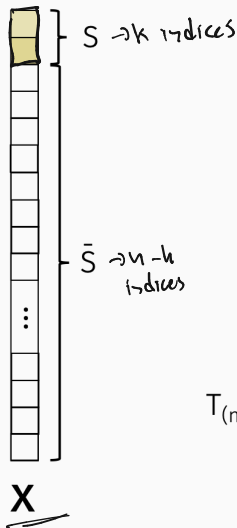
BASIS PURSUIT ANALYSIS

Some definitions:

$$Ax = b$$

2-sparse recovery

$$u=2$$

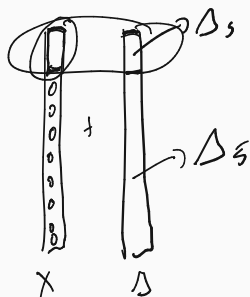


BASIS PURSUIT ANALYSIS

Claim 1: $\|\Delta_S\|_1 \geq \|\Delta_{\bar{S}}\|_1$ $A\Delta = 0$

$x + \Delta$ is a solution better than x , conjectured to exist by way of contradiction

$\|x + \Delta\|_1 \leq \|x\|_1$



$$\begin{aligned} \|x + \Delta\|_1 &= \|x + \Delta_S\|_1 + \|\Delta_{\bar{S}}\|_1 \\ &\geq \|x\|_1 - \|\Delta_S\|_1 + \|\Delta_{\bar{S}}\|_1 \\ &= \|x\|_1 + (\underbrace{\|\Delta_{\bar{S}}\|_1 - \|\Delta_S\|_1}_{\text{is positive}}) \end{aligned}$$

BASIS PURSUIT ANALYSIS

Claim 2: $\|\Delta_S\|_2 \geq \sqrt{2} \sum_{j \geq 2} \|\Delta_{T_j}\|_2$ $\sqrt{n} \|\Delta_S\|_2 \geq \|\Delta_S\|_1$

$$\|\Delta_S\|_1 \geq \sqrt{2k} \sum_{j \geq 2} \|\Delta_{T_j}\|_2$$

1. $\|\Delta_S\|_1 = \sum_{j \geq 1} \|\Delta_{T_j}\|_1$

2. For all j , $\|\Delta_{T_j}\|_1 \geq \sqrt{2k} \|\Delta_{T_{j+1}}\|_2$

$$\begin{aligned} \|\Delta_{T_j}\|_1 &\geq \underbrace{2k \cdot \min(|\Delta_{T_j}|)}_{\geq \frac{\|\Delta_{T_{j+1}}\|_2}{\sqrt{2k}}} \\ \|\Delta_{T_{j+1}}\|_2 &\leq \sqrt{2k \cdot \min(|\Delta_{T_j}|)} \\ &= \sqrt{2k} \cdot \min(|\Delta_{T_j}|) \end{aligned}$$

$$\|\Delta_{T_j}\|_2 \geq \sqrt{2k} \|\Delta_{T_{j+1}}\|_2 \quad 37$$

(3k, \epsilon) - BIP

Finish up proof by contradiction:

$$\underline{\|\Delta_S\|_2} \geq \sqrt{2} \sum_{j \geq 2} \|\Delta_{T_j}\|_2$$

$$\begin{aligned} \underline{0} = \|A\Delta\|_2 &\geq \underline{\|A\Delta_{S \cup T_1}\|_2} - \sum_{j \geq 2} \underline{\|A\Delta_{T_j}\|_2} \\ &\geq (1-\epsilon)\|\Delta_{S \cup T_1}\|_2 - (1+\epsilon) \sum_{j \geq 2} \|\Delta_{T_j}\|_2 \\ &\geq (1-\epsilon)\|\Delta_S\|_2 - (1+\epsilon) \frac{1}{\sqrt{2}} \|\Delta_S\|_2 \\ &= \left(1-\epsilon - \frac{(1+\epsilon)}{\sqrt{2}}\right) \|\Delta_S\|_2 \\ &\quad \underbrace{\hspace{10em}}_{\text{contradiction when positive}} \end{aligned}$$

A lot 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 $\min_z \|\mathbf{A}z - \mathbf{b}\|$ 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 Fourier matrices.
- Other “first order” type methods: Orthogonal Matching Pursuit, CoSaMP, Subspace Pursuit, etc.

When \mathbf{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+].

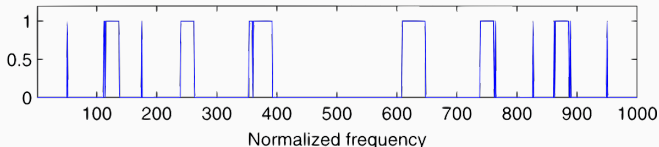
Hold up...

SPARSE FOURIER TRANSFORM

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

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