CS-GY 9223 D: Lecture 13 Compressed Sensing + Sparse Recovery

NYU Tandon School of Engineering, Prof. Christopher Musco

SPARSE RECOVERY/COMPRESSED SENSING

What do we know?

BASIC PROBLEM SETUP

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



• Infinite possible solutions **x**. In general, impossible to recover parameter vector.

SPARSITY RECOVERY/COMPRESSED SENSING

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

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



- 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 <u>choose</u> A's rows $\mathbf{a}_1, \dots, \mathbf{a}_m$. I.e. each $b_i = \mathbf{a}_i^T \mathbf{x}$ for some vector \mathbf{a}_i that we select.
- In this setting, we often call b_i a <u>linear measurement</u> of 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 μ -incoherent. $|\mathbf{A}_i^T \mathbf{A}_j| \le \mu \|\mathbf{A}_i\|_2 \|\mathbf{A}_j\|_2$ for all columns $\mathbf{A}_i, \mathbf{A}_j, i \ne j$.
 - Recover vectors **x** with sparsity $k = 1/\mu$.
- Focus today: A obeys the <u>Restricted Isometry Property</u>.

Definition ((q, ϵ)-Restricted Isometry Property) A matrix **A** satisfies (q, ϵ)-RIP if, for all **x** with $||\mathbf{x}||_0 \le q$, $(1 - \epsilon)||\mathbf{x}||_2^2 \le ||\mathbf{A}\mathbf{x}||_2^2 \le (1 + \epsilon)||\mathbf{x}||_2^2$.

- 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).

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

```
\min \|\mathbf{z}\|_0 \qquad subject \ to \qquad \mathbf{A}\mathbf{z} = \mathbf{b}.
```

Establishes that <u>information theoretically</u> we can recover
x. Solving the l₀-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 <u>unique</u> minimizer of min_{Az=b} $||\mathbf{z}||_0$.

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

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

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

or even

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

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:

<u>Gaspard Riche de Prony</u>, 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, ϵ) -RIP.

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

RESTRICTED ISOMETRY PROPERTY

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

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

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

The vectors that can be written as **Ax** for *k* sparse **x** lie in a union of *q* dimensional linear subspaces:





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

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

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

Quick argument:

Suppose you view a stream of numbers in 1,..., n:

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

After some time, you want to report which *k* items appeared <u>most frequently</u> 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 b = Ax for an approximately sparse x if there were only a few "heavy hitters". Recover x from b using a sparse recovery method (like l₀ minimization).

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:



single light sensor

Does not provide very much information about the image.

APPLICATION: SINGLE PIXEL CAMERA

But several random linear measurements do!

single light sensor



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

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

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



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



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

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

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, ϵ) -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. 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.

Warning: very cartoonish explanation of very complex problem.

Understanding what material is beneath the crust:



Think of vector **x** as scalar values of the density/reflectivity in a single vertical core of the earth.

How do we measure entries of Fourier transform Fx?

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 **x** here is a 2D image. Everything works with 2D Fourier transforms.

How do we measure entries of Fourier transform Fx?

APPLICATION: GEOPHYSICS

Blast the body with sounds waves of varying frequency.



The fewer measurements we need from **Fx**, 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, ϵ)**-Restricted Isometry Property)** A matrix **A** satisfies (q, ϵ)-RIP if, for all **x** with $||\mathbf{x}||_0 \le q$,

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

Lots of other random matrices satisfy RIP as well.

One major theoretical question is if we can <u>deterministically</u> <u>construct</u> 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 (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:

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

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

Convex relaxation of the ℓ_0 minimization problem:

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

$$\min_{\mathbf{z}} \|\mathbf{z}\|_1 \qquad subject \ to \qquad \mathbf{A}\mathbf{z} = \mathbf{b}.$$

- Objective is convex.
- Optimizing over convex set.

What is one method we know for solving this problem?

Equivalent formulation:

Problem (Basis Pursuit Linear Program.)

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

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

Suppose A is 2×1 , so b is just a scalar and x 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.

Theorem

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

Similar proof to ℓ_0 minimization:

• By way of contradiction, assume **x** is <u>not the optimal</u> solution. Then there exists some non-zero Δ such that:

$$\cdot \|\mathbf{x} + \Delta\|_1 \le \|\mathbf{x}\|_1$$

•
$$A(x + \Delta) = Ax$$
. I.e. $A\Delta = 0$.

Difference is that we can no longer assume that Δ is sparse.

We will argue that Δ is approximately sparse.

First tool:

For any *q*-sparse vector \mathbf{w} , $\|\mathbf{w}\|_2 \le \|\mathbf{w}\|_1 \le \sqrt{q} \|\mathbf{w}\|_2$

Second tool:

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

Some definitions:



BASIS PURSUIT ANALYSIS

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

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

Finish up proof by contradiction:

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 ~ 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 <u>O(k log^c n)</u> 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.