CS-GY 6763: Lecture 14 Finish Sparse Recovery and Compressed Sensing, Introduction to Spectral Sparsification

NYU Tandon School of Engineering, Prof. Christopher Musco

- Final project due next Wednesday, same day as final exam.
- Exam study guide will be released tonight.
- Solutions for last problem sets will be reviewed in office hours.

SPARSE RECOVERY/COMPRESSED SENSING PROBLEM SETUP

- Design a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $m < n, \mathbf{b} \in \mathbb{R}^{m}$.
- "Measure" $\mathbf{b} = \mathbf{A}\mathbf{x}$ for some <u>k-sparse</u> $\mathbf{x} \in \mathbb{R}^{n}$.



• Recover **x** from **b**.

Dual goals: Minimize <u>sample complexity</u> (number of rows in **A** and <u>computational complexity</u> to recover **x** from **A**, **b**.

APPLICATIONS

This simple to state problem models a lot of important real-world applications!



• Sample complexity usually corresponds to some application-dependent cost (e.g. length of time to acquire MRI, number of experiments needed to image below the earths surface).

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.

Typically design **A** with as few rows as possible that fulfills some desired property.

- 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 $k = 1/\mu$.

A obeys the (q, ϵ) -Restricted Isometry Property.

• Recover vectors **x** with sparsity k = O(q).

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

Can argue this property holds for random matrices (JL matrices) and subsampled Fourier matrices with roughly $m = O\left(\frac{q \log n}{\epsilon^2}\right)$ rows.

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 <u>unique</u> minimizer of:

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

 Establishes that information theoretically we can recover x in O(n^k) time from O(k log n) measurements.

Proof:

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 *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, 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?

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:

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:



Problem: This optimization problem naively takes $O(n^k)$ time to solve.

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.

Can be solved in poly(n) time using a linear program or using e.g. projected gradient descent. Other very relaxations also work. E.g. Lasso regularization $\min_{\mathbf{z}} \|\mathbf{A}\mathbf{z} - \mathbf{b}\|_2 + \lambda \|\mathbf{z}\|_1$.

Theorem

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

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.

Why ℓ_1 norm instead of ℓ_2 norm?

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 **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:
 - $\cdot \ \|x+\Delta\|_1 \leq \|x\|_1$
 - $A(x + \Delta) = Ax$. I.e. $A\Delta = 0$.

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

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: *S* is the set of *k* non-zero indices in **x**. \overline{T}_1 is the set of 2*k* indices <u>not in *S*</u> with largest magnitude in Δ . \overline{T}_2 is the set of 2*k* indices <u>not in *S*</u> with next largest magnitudes, etc.



Recall: By way of contradiction, if **x** is not the minimizer of the ℓ_1 problem, then there is some Δ such that $A(x + \Delta) = b$ and $||x + \Delta||_1 \le ||x||_1$.

Claim 1 (approximate sparsity of Δ): $\|\Delta_S\|_1 \ge \|\Delta_{\overline{S}}\|_1$

Claim 2 (ℓ_2 approximate sparsity): $\|\Delta_S\|_2 \ge \sqrt{2} \sum_{j \ge 2} \|\Delta_{T_j}\|_2$: We have:

$$\|\Delta_{S}\|_{2} \geq \frac{1}{\sqrt{k}} \|\Delta_{S}\|_{1} \geq \frac{1}{\sqrt{k}} \|\Delta_{\overline{S}}\|_{1} = \frac{1}{\sqrt{k}} \sum_{j \geq 1} \|\Delta_{T_{j}}\|_{1}.$$

So it suffices to show that: $\|\Delta_{T_j}\|_1 \ge \sqrt{2k} \|\Delta_{T_{j+1}}\|_2$

Finish up proof by contradiction: Recall that A is assumed to have the $(3k, \epsilon)$ RIP property. And by way of contradiction $A(x + \Delta) = b$.

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

We have that $(1-\epsilon) - \frac{1+\epsilon}{\sqrt{2}} \ge 0$ whenever $\epsilon < .17$.



Theorem

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

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+].

Wait a minute...

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 for image data is based on the idea that "natural images" are sparse if <u>some basis</u>. E.g. the DCT or Wavelet basis.



I.e. there is some representation of the image that requires many fewer numbers than explicitly writing down the pixels.

COMPRESSED SENSING RELATED TO MODERN DEEP LEARNING METHOD METHODS

Compressed Sensing using Generative Models

Ashish Bora*

Eric Price[‡]

Ajil Jalal[†]

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 lite near the range of a generative model $G: \mathbb{R}^k \to \mathbb{R}^n$. Our main theorem is that, if G is L-lapshitz, then roughly $O(k \log L)$ random Gaussian measurements suffice for an ℓ_2/ℓ_2 recovery guarantee. We demonstrate our results using generative models from published variational autoencoder and generative adversarial networks. Our method can use 5-10K rever measurements than Lass for the same accurecy.



Reconstruction using the same number of samples. Last row is method based on a GAN generative model.



Process: measure image x by computing b = Ax for a random matrix A. Use gradient descent to find $z \in \mathbb{R}^k$ to minimize:

 $\|A\mathcal{G}(z)-b\|.$

Return $\mathcal{G}(\mathbf{z})$.

A LITTLE ABOUT MY RESEARCH

Theorem (Subspace Embedding)

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

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

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

Implies regression result, and more.

Example: Any singular value $\tilde{\sigma}_i$ of **IIA** is a $(1 \pm \epsilon)$ approximation to the true singular value σ_i of **B**.

Recurring research interest: Replace random projection methods with <u>random sampling methods</u>. Prove that for essentially all problems of interest, can obtain same asymptotic runtimes.



Sampling has the added benefit of <u>preserving matrix sparsity</u> or structure, and can be applied in a <u>wider variety of settings</u> where random projections are too expensive. **Goal:** Can we use sampling to obtain subspace embeddings? I.e. for a given **A** find **Ã** whose rows are a (weighted) subset of rows in **A** and:



EXAMPLE WHERE STRUCTURE MATTERS

Let **B** be the edge-vertex incidence matrix of a graph *G* with vertex set *V*, |V| = d. Recall that $\mathbf{B}^T \mathbf{B} = \mathbf{L}$.



Recall that if $\mathbf{x} \in \{-1, 1\}^n$ is the <u>cut indicator vector</u> for a cut *S* in the graph, then $\frac{1}{4} \|\mathbf{B}\mathbf{x}\|_2^2 = \operatorname{cut}(S, V \setminus S)$.

LINEAR ALGEBRAIC VIEW OF CUTS



 $\mathbf{x} \in \{-1, 1\}^d$ is the <u>cut indicator vector</u> for a cut *S* in the graph, then $\frac{1}{4} \|\mathbf{Bx}\|_2^2 = \mathsf{cut}(S, V \setminus S)$ Extends to weighted graphs, as long as square root of weights is included in **B**. Still have the $\mathbf{B}^T \mathbf{B} = \mathbf{L}$.



And still have that if $\mathbf{x} \in \{-1, 1\}^d$ is the <u>cut indicator vector</u> for a cut S in the graph, then $\frac{1}{4} ||\mathbf{Bx}||_2^2 = \operatorname{cut}(S, V \setminus S)$.

Goal: Approximate **B** by a weighted subsample. I.e. by \tilde{B} with $m \ll |E|$ rows, each of which is a scaled copy of a row from **B**.



Natural goal: \tilde{B} is a subspace embedding for **B**. In other words, \tilde{B} has $\approx O(d)$ rows and for all **x**,

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

B is itself an edge-vertex incidence matrix for some <u>sparser</u> graph *G*! *G* is called a <u>spectral sparsifier</u> for *G*.



For example, we have that for any set S,

 $(1 - \epsilon) \operatorname{cut}_{G}(S, V \setminus S) \leq \operatorname{cut}_{\widetilde{G}}(S, V \setminus S) \leq (1 + \epsilon) \operatorname{cut}_{G}(S, V \setminus S).$

So \tilde{G} can be used in place of G in solving e.g. max/min cut problems, balanced cut problems, etc.

In contrast **ΠB** would look nothing like an edge-vertex incidence matrix if **Π** is a JL matrix.

Spectral sparsifiers were introduced in 2004 by Spielman and Teng in an influential paper on faster algorithms for solving Laplacian linear systems.

- Generalize the cut sparsifiers of Benczur, Karger '96.
- Further developed in work by Spielman, Srivastava + Batson, '08.
- Have had huge influence in algorithms, and other areas of mathematics – this line of work lead to the 2013 resolution of the Kadison-Singer problem in functional analysis by Marcus, Spielman, Srivastava.

Rest of class: Learn about an important random sampling algorithm for constructing spectral sparsifiers, and subspace embeddings for matrices more generally.

Goal: Find \tilde{A} such that $\|\tilde{A}x\|_2^2 = (1 \pm \epsilon) \|Ax\|_2^2$ for all x.

Possible Approach: Construct à by <u>uniformly sampling</u> rows from A. 6



Can check that this approach fails even for the special case of a graph vertex-edge incidence matrix.

Key idea: <u>Importance sampling</u>. Select some rows with higher probability.

Suppose A has *n* rows $\mathbf{a}_1 \dots, \mathbf{a}_n$. Let $p_1, \dots, p_n \in [0, 1]$ be sampling probabilities. Construct $\tilde{\mathbf{A}}$ as follows:

- For i = 1, ..., n
 - Select \mathbf{a}_i with probability p_i .
 - If \mathbf{a}_i is selected, add the scaled row $\frac{1}{\sqrt{p_i}}\mathbf{a}_i$ to $\tilde{\mathbf{A}}$.

Remember, ultimately want that $\|\tilde{\mathbf{A}}\mathbf{x}\|_2^2 = (1 \pm \epsilon)\|\mathbf{A}\mathbf{x}\|_2^2$ for all \mathbf{x} . Claim 1: $\mathbb{E}[\|\tilde{\mathbf{A}}\mathbf{x}\|_2^2] = \|\mathbf{A}\mathbf{x}\|_2^2$.

Claim 2: Expected number of rows in \tilde{A} is $\sum_{i=1}^{n} p_i$.

How should we choose the probabilities p_1, \ldots, p_n ?

MAIN RESULT

For i = 1, ..., n,

$$\tau_i = \mathbf{a}_i^T (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{a}_i.$$

Theorem (Subspace Embedding from Subsampling)

For each *i*, and fixed constant *c*, let $p_i = \min\left(1, \frac{c\log d}{\epsilon^2} \cdot \tau_i\right)$. Let \tilde{A} have rows sampled from A with probabilities p_1, \ldots, p_n . With probability 9/10,

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

and \tilde{A} has $O(d \log d/\epsilon^2)$ rows in expectation.

How should we choose the probabilities p_1, \ldots, p_n ?

As usual, consider a single vector \mathbf{x} and understand how to sample to preserve norm of $\mathbf{y} = \mathbf{A}\mathbf{x}$:

$$\|\mathbf{\tilde{A}}\mathbf{x}\|_{2}^{2} = \|\mathbf{S}\mathbf{A}\mathbf{x}\|_{2}^{2} = \|\mathbf{S}\mathbf{y}\|_{2}^{2} \approx \|\mathbf{y}\|_{2}^{2} = \|\mathbf{A}\mathbf{x}\|_{2}^{2}.$$

Then we can union bound over an ϵ -net to extend to all **x**.

As discussed a few lectures ago, uniform sampling only works well if y = Ax is "flat".



Instead consider sampling with probabilities at least proportional to the magnitude of **y**'s entries:

$$p_i > c \cdot \frac{y_i^2}{\|y\|_2^2}$$
 for constant *c* to be determined.

Using a Bernstein bound (or Chebyshev's inequality if you don't care about the δ dependence) we have that if $c = \frac{\log(1/\delta)}{c^2}$ then:

$$\Pr[\left|\|\tilde{\mathbf{y}}\|_2^2 - \|\mathbf{y}\|_2^2\right| \ge \epsilon \|\mathbf{y}\|_2^2] \le \delta.$$

The number of samples we take in expectation is:

$$\sum_{i=1}^{n} p_i = \sum_{i=1}^{n} c \cdot \frac{y_i^2}{\|y_i\|_2^2} = \frac{\log(1/\delta)}{\epsilon^2}.$$

We don't know $y_1, \ldots, y_n!$ And in fact, these values aren't fixed. We wanted to prove a bound for $\mathbf{y} = \mathbf{A}\mathbf{x}$ for any \mathbf{x} .

Idea behind leverage scores: Sample row *i* from **A** using the worst case (largest necessary) sampling probability:

$$au_i = \max_{\mathbf{x}} \frac{y_i^2}{\|\mathbf{y}\|_2^2}$$
 where $\mathbf{y} = \mathbf{A}\mathbf{x}$

If we sample with probability $p_i = \frac{1}{\epsilon^2} \cdot \tau_i$, then we will be sampling by at least $\frac{1}{\epsilon^2} \cdot \frac{y_i^2}{\|\mathbf{y}\|_2^2}$, <u>no matter what **y** is</u>.

$$au_i = \max_{\mathbf{x}} \frac{y_i^2}{\|\mathbf{y}\|_2^2}$$
 where $\mathbf{y} = \mathbf{A}\mathbf{x}$.

A little messy algebra shows that $\mathbf{x}^* = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{a}_i$.

Two concerns:

1) How to efficiently compute τ_1, \ldots, τ_n ?

2) The number of samples we take will be roughly $\sum_{i=1}^{n} \tau_i$. How do we bound this?

Topic for another day!

In many applications, computational costs are second order to data collection costs. We have a huge range of possible data points $\mathbf{a}_1, \ldots, \mathbf{a}_n$ that we can collect labels/values b_1, \ldots, b_n for. Goal is to learn \mathbf{x} such that:

 $\mathbf{a}_i^T \mathbf{x} \approx b_i$.

Want to do so after observing as few b_1, \ldots, b_n as possible. Applications include healthcare, environmental science, etc.



Can be solved via random sampling for linear models.



Claim: Let \tilde{A} is an O(1)-factor subspace embedding for A(obtained via leverage score sampling). Then $\tilde{x} = \arg \min \|\tilde{A}x - \tilde{b}\|_2^2$ satisfies:

 $\|\mathbf{A}\tilde{\mathbf{x}} - \mathbf{b}\|_2^2 \le O(1)\|\mathbf{A}\mathbf{x}^* - \mathbf{b}\|_2^2$

Computing $\tilde{\mathbf{x}}$ only requires collecting $O(d \log d)$ labels.

Lots of applications:

- Robust bandlimited and multiband interpolation [STOC 2019].
- Active learning for Gaussian process regression [NeurIPS 2020].
- \cdot Active learning beyond the ℓ_2 norm [FOCS 2022]
- Active learning for polynomial regression [SODA 2023]
- Active learning for 1 layer neural nets [NeurIPS 2023]
- DOE Grant on "learning based" algorithms for solving parametric partial differential equations.