Computational Principles for High-dim Data Analysis (Lecture Four)

Yi Ma

EECS Department, UC Berkeley

September 7, 2021



EECS208, Fall 2021

Convex Methods for Sparse Signal Recovery

1 Geometric Intuition

2 A First Correctness Result via Incoherence

Coherence of a Matrix Correctness of ℓ^1 Minimization Constructing an Incoherent Matrix Limitations of Incoherence

"Algebra is but written geometry; geometry is but drawn algebra." – Sophie Germain

Geometric Intuition: Coefficient Space

Given $\boldsymbol{y} = \boldsymbol{A} \boldsymbol{x}_o \in \mathbb{R}^m$ with $\boldsymbol{x}_o \in \mathbb{R}^n$ sparse:

 $\min \|\boldsymbol{x}\|_1 \quad \text{subject to} \quad \boldsymbol{A}\boldsymbol{x} = \boldsymbol{y}. \tag{1}$

The space of all feasible solutions is an affine subspace:

$$S = \{ \boldsymbol{x} \mid \boldsymbol{A} \boldsymbol{x} = \boldsymbol{y} \} = \{ \boldsymbol{x}_o \} + \operatorname{null}(\boldsymbol{A}) \quad \subset \mathbb{R}^n.$$
 (2)



• • • • • • • • • • • •

ℓ^1 Minimization in the Coefficient Space

Gradually expand a ℓ^1 ball of radius t from the origin **0**:

$$\mathbf{x} \cdot \mathbf{B}_1 = \{ \boldsymbol{x} \mid \| \boldsymbol{x} \|_1 \le t \} \quad \subset \mathbb{R}^n,$$
(3)

till its boundary first touches the feasible set S:



Comparison between ℓ^1 and ℓ^2 Minimization Given $y = Ax_o$ with x_o sparse:

 $\mathbf{A}: \min \|\boldsymbol{x}\|_1 \text{ subject to } \boldsymbol{A}\boldsymbol{x} = \boldsymbol{y}. \tag{4}$

versus

 $\mathbf{B}: \quad \min \|\boldsymbol{x}\|_2 \quad \text{subject to} \quad \boldsymbol{A}\boldsymbol{x} = \boldsymbol{y} \tag{5}$



Sparsity Promoting with Different ℓ^p Norms



Figure: Intersection between the ℓ^p -ball and the feasible set S, for p = 2, 1.5, 1 and 0.7, respectively. (Some argue p = 0.5 is somewhat special.)

Figure from Sparse and Redundant Representations, Michael Elad, Springer, 2010.

EECS208, Fall 2021

Geometric Intuition: High-dimensional Polytopes

Neighborly Polytopes

(vertices from a Gaussian matrix):

$$oldsymbol{A} = [oldsymbol{a}_1, oldsymbol{a}_2, \dots, oldsymbol{a}_n] \in \mathbb{R}^{m imes n}$$

The "correct" visualization of high-dimensional convex polytopes,¹ including the ℓ^1 ball:





A general convex set

The ℓ_1 ball

¹Lectures on Discrete Geometry, Jiri Matousek, Springer 2002

Geometric Intuition: Observation Space

The matrix $A \in \mathbb{R}^{m \times n}$ can be viewed as a linear projection from \mathbb{R}^n to \mathbb{R}^m :

 $\boldsymbol{A}: \ \mathsf{B}_1 \to \mathsf{P} = \boldsymbol{A}(\mathsf{B}_1), \quad \text{(6)}$

which maps a convex polytope to a convex polytope. Similarly, $\forall t \ge 0$:

 $t \cdot \mathsf{B}_1 \to t \cdot \mathbf{A}(\mathsf{B}_1).$



< □ > < 同 > < 三 > < 三 >

Geometric Intuition: Observation Space

All k-faces of B_1 cannot be mapped to the inside of the polytope $A(B_1)$:

A Million Dollar Question: When $\hat{x} = x_o$?

Coherence of a Matrix

Definition (Mutual Coherence)

For a matrix $\mathbf{A} = \begin{bmatrix} \mathbf{a}_1 \mid \mathbf{a}_2 \mid \cdots \mid \mathbf{a}_n \end{bmatrix} \in \mathbb{R}^{m \times n}$ with nonzero columns, the *mutual coherence* $\mu(\mathbf{A})$ is the largest normalized inner product between two distinct columns:

$$u(\boldsymbol{A}) = \max_{i \neq j} \left| \left\langle \frac{\boldsymbol{a}_i}{\|\boldsymbol{a}_i\|_2}, \frac{\boldsymbol{a}_j}{\|\boldsymbol{a}_j\|_2} \right\rangle \right|.$$
(7)

Example:

Uniqueness of Sparse Solution

Proposition (Coherence Controls Kruskal Rank)

For any $oldsymbol{A} \in \mathbb{R}^{m imes n}$,

$$\operatorname{krank}(\boldsymbol{A}) \geq \frac{1}{\mu(\boldsymbol{A})}.$$
(8)

In particular, if $oldsymbol{y} = oldsymbol{A} oldsymbol{x}_o$ and

$$\|\boldsymbol{x}_o\|_0 \leq \frac{1}{2\mu(\boldsymbol{A})},\tag{9}$$

then x_o is the unique optimal solution to the ℓ^0 minimization problem

$$\min \|\boldsymbol{x}\|_0 \quad s.t. \quad \boldsymbol{A}\boldsymbol{x} = \boldsymbol{y}. \tag{10}$$

Proof:

$$1 - k\mu(\mathbf{A}) < \sigma_{\min}(\mathbf{A}_{\mathsf{I}}^*\mathbf{A}_{\mathsf{I}}) \leq \sigma_{\max}(\mathbf{A}_{\mathsf{I}}^*\mathbf{A}_{\mathsf{I}}) < 1 + k\mu(\mathbf{A}).$$
(11)

Theorem (ℓ^1 Succeeds under Incoherence)

Let A be a matrix whose columns have unit ℓ^2 norm, and let $\mu(A)$ denote its mutual coherence. Suppose that $y = Ax_o$, with

$$\|\boldsymbol{x}_o\|_0 \le rac{1}{2\mu(\boldsymbol{A})}.$$
 (12)

Then x_o is the unique optimal solution to the problem

$$\min \|\boldsymbol{x}\|_1 \quad s.t. \quad \boldsymbol{y} = \boldsymbol{A}\boldsymbol{x}. \tag{13}$$

Tightness: there exist examples of A and x_o with $||x_o||_0 > \frac{1}{2} \left(1 + \frac{1}{\mu(A)}\right)$ for which ℓ^1 minimization does not recover x_o .

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ののの

Given $oldsymbol{y} = oldsymbol{A} oldsymbol{x}_o$, try to find $oldsymbol{x}_o$ via ℓ^1 minimization:

$$\min \|\boldsymbol{x}\|_1 \quad \text{s.t.} \quad \boldsymbol{y} = \boldsymbol{A}\boldsymbol{x}. \tag{14}$$

Lagrangian formulation:

$$\min \|\boldsymbol{x}\|_1 + \boldsymbol{\lambda}^*(\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}), \quad \exists \boldsymbol{\lambda} \in \mathbb{R}^m.$$
(15)

Optimality condition: x_o is a minimum of f(x) if and only if 0 is in the subgradient $\partial f(x)$ at x_o :

$$f(\boldsymbol{x}) \geq f(\boldsymbol{x}_o) + \boldsymbol{0}^*(\boldsymbol{x} - \boldsymbol{x}_o).$$

Optimality condition for ℓ^1 Minimization:

$$\mathbf{0} \in \partial \|\boldsymbol{x}_o\|_1 - \boldsymbol{A}^* \boldsymbol{\lambda} \quad \Leftrightarrow \quad \boldsymbol{A}^* \boldsymbol{\lambda} \in \partial \|\boldsymbol{x}_o\|_1. \tag{16}$$

Proof (a sketch of key ideas):

Due to convexity of $\|\cdot\|_1$, for any $m{v}\in\partial\,\|\cdot\|_1\,(m{x}_o)$ and $m{x}'\in\mathbb{R}^n$,

$$\|\boldsymbol{x}'\|_1 \ge \|\boldsymbol{x}_o\|_1 + \langle \boldsymbol{v}, \boldsymbol{x}' - \boldsymbol{x}_o \rangle \tag{17}$$

For $\boldsymbol{v} = \boldsymbol{A}^* \boldsymbol{\lambda}$, we have: $\langle \boldsymbol{A}^* \boldsymbol{\lambda}, \boldsymbol{x}' - \boldsymbol{x}_o \rangle = \langle \boldsymbol{\lambda}, \boldsymbol{A}(\boldsymbol{x}' - \boldsymbol{x}_o) \rangle = 0$. Therefore $\| \boldsymbol{x}' \|_1 \ge \| \boldsymbol{x}_o \|_1$.

To find such an optimality certificate $A^*\lambda \in \partial \left\|\cdot\right\|_1(x_o)$, we need:

$$A_{\mathsf{I}}^* \lambda = \sigma, \quad \|A_{\mathsf{I}^c}^* \lambda\|_{\infty} \le 1.$$
 (18)

A natural "candidate":

$$\hat{\boldsymbol{\lambda}}_{\ell^2} \doteq \boldsymbol{A}_{\mathsf{I}} (\boldsymbol{A}_{\mathsf{I}}^* \boldsymbol{A}_{\mathsf{I}})^{-1} \boldsymbol{\sigma}.$$
(19)

The rest is to check this satisfies (18) under the given conditions.

Proof (continued):

By construction, $A_1^* \hat{\lambda}_{\ell^2} = \sigma$. We are just left to verify (18), by calculating

$$\|\boldsymbol{A}_{\mathsf{I}^c}^* \hat{\boldsymbol{\lambda}}_{\ell^2}\|_{\infty} = \|\boldsymbol{A}_{\mathsf{I}^c}^* \boldsymbol{A}_{\mathsf{I}} (\boldsymbol{A}_{\mathsf{I}}^* \boldsymbol{A}_{\mathsf{I}})^{-1} \boldsymbol{\sigma}\|_{\infty}.$$
 (20)

Consider a single element of this vector $(j \in I^c)$, which has the form:

$$\begin{aligned} |a_{j}^{*}A_{\mathsf{I}}(A_{\mathsf{I}}^{*}A_{\mathsf{I}})^{-1}\sigma| &\leq \underbrace{\|A_{\mathsf{I}}^{*}a_{j}\|_{2}}_{\leq\sqrt{k}\mu}\underbrace{\|(A_{\mathsf{I}}^{*}A_{\mathsf{I}})^{-1}\|_{2,2}}_{<\frac{1}{1-k\mu(A)}}\underbrace{\|\sigma\|_{2}}_{=\sqrt{k}} \end{aligned} (21) \\ &\leq \frac{k\mu(A)}{1-k\mu(A)} \\ &\leq \frac{1}{1-k\mu(A)} \end{aligned} (22) \\ &\leq \frac{1}{\mathsf{Provided }k\mu(A) \leq 1/2}. \end{aligned}$$

Example I. Consider a discrete Fourier transform matrix F. Let $I \subset [n]$ be a random set of m indices,

$$\boldsymbol{A} = \boldsymbol{F}_{\mathsf{I}}^* \in \mathbb{C}^{m \times n}. \tag{24}$$

Example II. For two orthogonal matrices Φ and Ψ ,

$$\boldsymbol{A} = \boldsymbol{\Phi}_{\mathsf{I}}^* \boldsymbol{\Psi}. \tag{25}$$

(日) (四) (日) (日) (日)

Example III. For two orthogonal matrices, say Φ is Fourier F and Ψ is the identify I or the Wavelet W,

$$\boldsymbol{A} = [\boldsymbol{\Phi} \mid \boldsymbol{\Psi}] \in \mathbb{C}^{n \times 2n}.$$
 (26)

Incoherence and Uncertainty Principle

Incoherence between I and F: $|\langle e_i, f_j \rangle| = \frac{1}{\sqrt{n}}$.

Facts: A signal cannot be sparse in both time I and frequency F. Let $\hat{x} = Fx \in \mathbb{C}^n$ be the discrete Fourier transform of $x \in \mathbb{C}^n$. Then the **Heisenberg uncertainty principle** states that:

$$\operatorname{Var}(|\boldsymbol{x}|^2)\operatorname{Var}(|\hat{\boldsymbol{x}}|^2) \ge \frac{1}{16\pi^2}.$$
(27)

Or a deterministic uncertainty principle:

$$\|\boldsymbol{x}\|_{0} \cdot \|\hat{\boldsymbol{x}}\|_{0} \ge n \quad \text{or} \quad \|\boldsymbol{x}\|_{0} + \|\hat{\boldsymbol{x}}\|_{0} \ge 2\sqrt{n}.$$
 (28)

Incoherence and Uncertainty Principle

Theorem (Uncertainty Principle I²)

For $\mathbf{A} = [\mathbf{\Phi} \mid \mathbf{\Psi}] \in \mathbb{C}^{n \times 2n}$ with two orthogonal matrices $\mathbf{\Phi}$ and $\mathbf{\Psi}$. For any $\mathbf{0} = \mathbf{\Phi} \mathbf{e} + \mathbf{\Psi} \hat{\mathbf{e}}$ with $\mathbf{\Phi} \mathbf{e} = -\mathbf{\Psi} \hat{\mathbf{e}} \neq \mathbf{0}$, we have $\|\mathbf{e}\|_0 + \|\hat{\mathbf{e}}\|_0 \ge \frac{2}{\mu(\mathbf{A})}.$ (29)

Corollary (Uncertainty Principle II)

For $A = [\Phi \mid \Psi] \in \mathbb{C}^{n \times 2n}$ with two orthogonal matrices Φ and Ψ . For any nonzero $y = Ax = A\hat{x}$ with x different from \hat{x} , we have $\|x\|_0 + \|\hat{x}\|_0 \ge \frac{2}{\mu(A)}.$ (30)

Question: What can you say about y = Ax with $||x||_0 < \frac{1}{\mu(A)}$?

²Sparse and Redundant Representations, Michael Elad, Springer, 2010. () Solution ()

18 / 26

Recall phenomena associated with random matrices:

• Measure Concentration ($\epsilon \sim O(n^{-1/2})$)

$$\mathsf{Area}\{\boldsymbol{x} \in \mathbb{S}^{n-1} : -\epsilon \le x_n \le \epsilon\} = 0.99 \cdot \mathsf{Area}(\mathbb{S}^{n-1}), \qquad \textbf{(31)}$$

• Neighborly Polytopes (vertices from a Gaussian matrix):

$$\boldsymbol{A} = [\boldsymbol{a}_1, \boldsymbol{a}_2, \dots, \boldsymbol{a}_n] \in \mathbb{R}^{m \times n}.$$

Theorem (Spherical Measure Concentration³)

Let $u \sim uni(\mathbb{S}^{m-1})$ be distributed according to the uniform distribution on the sphere. Let $f : \mathbb{S}^{m-1} \to \mathbb{R}$ be an 1-Lipschitz function:

$$\forall \boldsymbol{u}, \, \boldsymbol{u}', \quad \left| f(\boldsymbol{u}) - f(\boldsymbol{u}') \right| \leq 1 \cdot \left\| \boldsymbol{u} - \boldsymbol{u}' \right\|_2, \tag{32}$$

and let med(f) denote any median of the random variable Z = f(u). Then

$$\mathbb{P}[f(\boldsymbol{u}) > \operatorname{med}(f) + t] \leq 2 \exp\left(-\frac{mt^2}{2}\right), \quad (33)$$
$$\mathbb{P}[f(\boldsymbol{u}) < \operatorname{med}(f) - t] \leq 2 \exp\left(-\frac{mt^2}{2}\right). \quad (34)$$

³Lectures on Discrete Geometry, Jiri Matousek, Springer 2002 + (=) (=)

Theorem

Let $A = [a_1 | \cdots | a_n]$ with columns $a_i \sim \operatorname{uni}(\mathbb{S}^{m-1})$ chosen independently according to the uniform distribution on the sphere. Then with probability at least 3/4,

$$\mu(\mathbf{A}) \leq C \sqrt{\frac{\log n}{m}},\tag{35}$$

where C > 0 is a numerical constant.

Proof (a sketch): For any $v \in \mathbb{S}^{m-1}$, $\mathbb{E}[|v^*a|]^2 \leq [(v^*a)^2] \leq \frac{1}{m}$ implies

$$\operatorname{med}(|\boldsymbol{v}^*\boldsymbol{a}|) \leq 2\mathbb{E}[|\boldsymbol{v}^*\boldsymbol{a}|] \leq \frac{2}{\sqrt{m}} \quad (\mathsf{Markov inequality}).$$
$$\mathbb{P}\left[|\boldsymbol{v}^*\boldsymbol{a}| > \frac{2+t}{\sqrt{m}}\right] \leq 2\exp\left(-\frac{t^2}{2}\right). \tag{36}$$

Proof (continued):

As all the *n* columns $\{a_i\}$ are independent:

$$\mathbb{P}\left[|\boldsymbol{a}_{i}^{*}\boldsymbol{a}_{j}| > \frac{2+t}{\sqrt{m}}\right] \leq 2\exp\left(-\frac{t^{2}}{2}\right).$$
(37)

Summing the failure probability over all n(n-1)/2 pairs of $(\boldsymbol{a}_i, \boldsymbol{a}_j)$:

$$\mathbb{P}\left[\exists (i,j) : |\boldsymbol{a}_i^*\boldsymbol{a}_j| > \frac{2+t}{\sqrt{m}}\right] \leq n(n-1)\exp\left(-\frac{t^2}{2}\right).$$
(38)

Setting $t = 2\sqrt{\log 2n}$, the RHS probability is less than 1/4.

(日) (四) (日) (日) (日)

Limitations of Incoherence

Theorem (Welch Bound)

For any matrix $A = [a_1 | \cdots | a_n] \in \mathbb{R}^{m \times n}$, $m \le n$, and suppose that the columns a_i have unit ℓ^2 norm. Then

$$\mu(\boldsymbol{A}) = \max_{i \neq j} |\langle \boldsymbol{a}_i, \boldsymbol{a}_j \rangle| \geq \sqrt{\frac{n-m}{m(n-1)}} = \Omega\left(\frac{1}{\sqrt{m}}\right).$$
(39)

Guarantee based on $\mu(A)$

Limitations of Incoherence

Proof of the Welch bound.

Let $G = A^*A \in \mathbb{R}^{n \times n}$ and its eigenvalues satisfy: $\sum_{i=1}^m \lambda_i(G)$ = trace $(G) = \sum_{i=1}^n \|a_i\|_2^2 = n$. Using this fact, we have:

$$\frac{n^{2}}{m} \leq \frac{n^{2}}{m} + \sum_{i=1}^{m} \left(\lambda_{i}(G) - \frac{n}{m}\right)^{2}$$

$$= \frac{n^{2}}{m} + \sum_{i=1}^{m} \left\{\lambda_{i}^{2}(G) + \frac{n^{2}}{m^{2}} - 2\frac{n}{m}\lambda_{i}(G)\right\}$$

$$= \sum_{i=1}^{m} \lambda_{i}^{2}(G) = ||G||_{F}^{2} = \sum_{i,j} |a_{i}^{*}a_{j}|^{2} = n + \sum_{i \neq j} |a_{i}^{*}a_{j}|^{2}$$

$$\leq m + m(m-1)(mon |a^{*}a_{j}|)^{2}$$
(40)
(41)
(41)

$$\leq n+n(n-1)\big(\max_{i\neq j}|\boldsymbol{a}_i^*\boldsymbol{a}_j|\big)^2. \tag{43}$$

Ma (EECS Department, UC Berkeley)

イロト イヨト イヨト

Limitations of Incoherence

Incoherence ensures to recover k-sparse solution from

 $m \geq \tilde{O}(k^2)$

measurements.

Experimental results suggest m = O(k):

In a proportional growth setting $m \propto n$, $k \propto m$, ℓ^1 minimization succeeds with very high probability whenever the constants of proportionality n/m and k/m are small enough.

Next: how to sharpen the bound?

< 日 > < 同 > < 三 > < 三 >

Assignments

- Reading: Section 3.1 & 3.2 of Chapter 3.
- Programming Homework # 1.

< 1 k

э