# ECE 8201: Low-dimensional Signal Models for High-dimensional Data Analysis

Lecture 3: Sparse signal recovery: A RIPless analysis of $\ell_1$ minimization

Yuejie Chi

The Ohio State University

THE OHIO STATE UNIVERSITY

# Outline

- A RIPless theory for CS using $\ell_1$ minimization recovery

  **Reference:** E. J. Candes and Y. Plan. A probabilistic and RIPless theory of compressed sensing. 2010.

# Subgradient of $\ell_1$ function

Consider a convex function $f(\boldsymbol{x})$.

**Definition 1. [Subgradient]** $\boldsymbol{u} \in \partial f(\boldsymbol{x}_0)$ *is a subgradient of a convex $f$ at $\boldsymbol{x}_0$ if for all $\boldsymbol{x}$:*

$$f(\boldsymbol{x}) \geq f(\boldsymbol{x}_0) + \boldsymbol{u}^T(\boldsymbol{x} - \boldsymbol{x}_0)$$

**Remark:** if $f$ is differentiable at $\boldsymbol{x}_0$, the only subgradient is the gradient $\nabla f(\boldsymbol{x}_0)$.

**Example:** For the scalar absolute function $f(t) = |t|$, $t \in \mathbb{R}$, $u \in \partial f(t)$ iff

$$\begin{cases} u = \mathsf{sgn}(t), & t \neq 0 \\ u \in [-1, 1], & t = 0 \end{cases}$$

For $f(\boldsymbol{x}) = \|\boldsymbol{x}\|_1$, $\boldsymbol{x} \in \mathbb{R}^n$, $\boldsymbol{u} \in \partial f(\boldsymbol{x})$ iff

$$\begin{cases} u_i = \mathsf{sgn}(x_i), & x_i \neq 0 \\ u_i \in [-1, 1], & x_i = 0 \end{cases}$$

# Characterization of $\ell_1$ solution: an optimization viewpoint

**Proposition 1. [Necessary and Sufficient condition for $\ell_1$ recovery]** *Denote the support of $x$ as $T$. $x$ is the solution to BP if for all $h \in \mathrm{Null}(A)$,*

$$\sum_{i \in T} sign(x_i)h_i \leq \sum_{i \in T^c} |h_i|.$$

*Furthermore, $x$ is the unique solution if the equality holds iff $h = 0$.*

**Remark:** Recovery property only depends on the sign pattern of $x$, not the magnitudes!

Proof of Proposition 1: We first show it is a sufficient condition. Denote the solution of BP as $\hat{x} = x + h$. We have

$$Ah = A(\hat{x} - x) = 0,$$

i.e. $h \in \mathrm{Null}(A)$.

Since $\boldsymbol{x}$ is supported on $T$, we have

$$\|\boldsymbol{x}\|_1 \geq \|\hat{\boldsymbol{x}}\|_1 = \|\boldsymbol{x} + \boldsymbol{h}\|_1 = \sum_{i \in T} |x_i + h_i| + \sum_{i \in T^c} |h_i|$$

$$\geq \sum_{i \in T} |x_i| + \text{sign}(x_i)h_i + \sum_{i \in T^c} |h_i| \geq \sum_{i \in T} |x_i| = \|\boldsymbol{x}\|_1.$$

Therefore $\boldsymbol{h} = 0$ and $\hat{\boldsymbol{x}} = \boldsymbol{x}$. Next we show it is also a necessary condition. If there exists $\boldsymbol{h} \in \text{Null}(\boldsymbol{A})$ such that

$$\sum_{i \in T} \text{sign}(x_i)h_i > \sum_{i \in T^c} |h_i|$$

then we can verify

$$\|\boldsymbol{x} - \boldsymbol{h}\|_1 = \sum_{i \in T} |x_i - h_i| + \sum_{i \in T^c} |h_i| < \sum_{i \in T}(|x_i| - \text{sign}(x_i)h_i) + \sum_{i \in T^c} |h_i|$$

$$< \sum_{i \in T} |x_i| = \|\boldsymbol{x}\|_1.$$

# Dual certificate

Denote the support of $\boldsymbol{x}$ as $T$.

**Proposition 2.** *$\boldsymbol{x}$ is an optimal solution of BP iff there exists $\boldsymbol{u} = \boldsymbol{A}^{\mathsf{T}}\boldsymbol{\lambda}$ such that*

$$
\begin{cases}
u_i = \mathsf{sgn}(x_i), & i \in T \\
u_i \in [-1, 1], & i \in T^c
\end{cases}
$$

*In addition, if $|u_i| < 1$ for $i \in T^c$ and $\boldsymbol{A}_T$ has full columns rank, $\boldsymbol{x}$ is the <span style="color:red">unique</span> solution.*
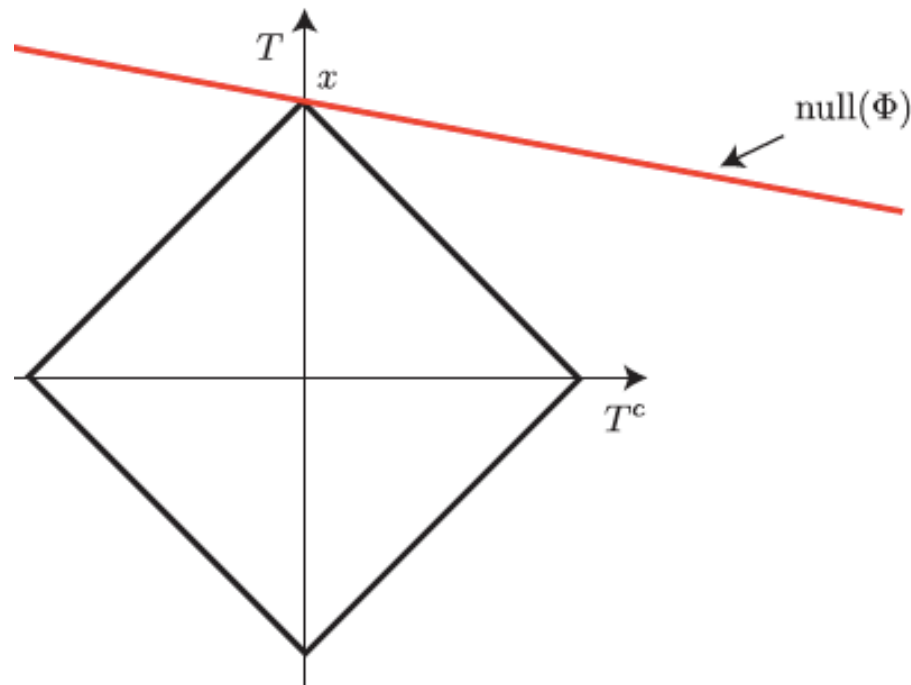
**Remarks:**

- We call $\boldsymbol{u}$ or $\boldsymbol{\lambda}$ the *(exact) dual certificate*. If we can find such a dual certificate, we can verify the optimality of BP.

- Note that $\boldsymbol{u} \perp \mathrm{Null}(\boldsymbol{A}))$, which is also a subgradient of $\|\boldsymbol{x}\|_1$ at $\boldsymbol{x}$.

# Dual certificate: geometric interpretation

Geometric interpretation of the dual certificate: there exists a subgradient $\boldsymbol{u}$ of the objective function $\|\boldsymbol{x}\|_1$ at the ground truth $\boldsymbol{x}$ such that

$$\boldsymbol{u} \perp \mathrm{Null}(\boldsymbol{A})$$

# Unicity

If supp$(\boldsymbol{x}) \subset T$, $\boldsymbol{A}\boldsymbol{x} = \boldsymbol{A}_T \boldsymbol{x}_T$.

Note for any $\boldsymbol{h} \in \mathrm{Null}(\boldsymbol{A})$,

$$\sum_{i \in T} \mathsf{sgn}(x_i) h_i = \sum_{i \in T} u_i h_i = \langle \boldsymbol{u}, \boldsymbol{h} \rangle - \sum_{i \in T^c} u_i h_i$$

$$= -\sum_{i \in T^c} u_i h_i \qquad (\text{since } \boldsymbol{u} \perp \mathrm{Null}(\boldsymbol{A}))$$

$$< \sum_{i \in T^c} |h_i| \qquad (\text{since } |u_i| < 1 \text{for } i \in T^c)$$

unless $\boldsymbol{h}_{T^c} \neq 0$. If $\boldsymbol{h}_{T^c} = 0$, since $\boldsymbol{A}_T$ has full column rank,

$$\boldsymbol{A}\boldsymbol{h} = \boldsymbol{A}_T \boldsymbol{h}_T = 0$$

which indicates $\boldsymbol{h}_T = 0$ as well. In summary $\boldsymbol{h} = \boldsymbol{h}_T + \boldsymbol{h}_{T^c} = 0$, and $\boldsymbol{x}$ is the unique solution.

# A Probabilistic Approach with Gaussian matrices

Our goal is to develop a theory of compressed sensing that 1) does not require RIP; and 2) admits near-optimal performance guarantees.

Let $A$ be composed of i.i.d. $\mathcal{N}(0,1)$ entries.

**Question:** How well does BP perform for an arbitrary but fixed sparse signal?

$$\text{(BP:)} \quad \hat{x} = \operatorname*{argmin}_{x} \ \|x\|_1 \quad \text{subject to} \quad y = Ax.$$

**Theorem 1.** *Let $x \in \mathbb{R}^n$ be an arbitrary* fixed *vector that is $k$-sparse. Assume $A$ is composed of i.i.d. $\mathcal{N}(0,1)$ entries. As long as $m \geq C_1 k \log n$ for some large enough constant $C_1$, $x$ is the unique solution to BP with probability at least $1 - n^{-C_2}$ for some constant $C_2$.*

Remark: Compare this result with the earlier RIP-based result.

# Proof by certifying the dual certificate

Denote the support of $x$ as $T$.

We first verify that $A_T$ is full column rank with high probability.

Since $A_T$ is a fixed $m \times k$ matrix with i.i.d. $\mathcal{N}(0,1)$ entries, random matrix theory tells us (we'll just take for granted)

$$\mathbb{P}\left(\frac{1}{\sqrt{m}}\sigma_{\max}(A) > 1 + \sqrt{\frac{k}{m}} + t\right) \le e^{-mt^2/2}$$

$$\mathbb{P}\left(\frac{1}{\sqrt{m}}\sigma_{\min}(A) < 1 - \sqrt{\frac{k}{m}} - t\right) \le e^{-mt^2/2}.$$

Then as long as $m \ge c_1 k$ for some large constant $c_1$, we have

$$\left\|\frac{1}{m}A_T^\mathsf{T}A_T - I\right\| \le \frac{1}{2}$$

with probability at least $1 - e^{-c_2 m}$ for some $c_2$. Call this event $\mathcal{A}$.

# Construction of the dual certificate

We need to find a dual certificate $\boldsymbol{u} = \boldsymbol{A}^\mathsf{T}\boldsymbol{\lambda}$ such that

$$\begin{cases} u_i = \mathsf{sgn}(x_i), & i \in T \\ |u_i| < 1, & i \in T^c \end{cases}$$

Consider the solution to the following $\ell_2$ minimization problem:

$$\min \|\boldsymbol{u}\|_2 \quad \text{s.t.} \quad \boldsymbol{u} = \boldsymbol{A}^\mathsf{T}\boldsymbol{\lambda}, \quad u_i = \mathsf{sgn}(x_i), \quad i \in T$$

which can be written explicitly as

$$\boldsymbol{u} = \boldsymbol{A}^\mathsf{T}\boldsymbol{A}_T(\boldsymbol{A}_T^\mathsf{T}\boldsymbol{A}_T)^{-1}\mathsf{sgn}(\boldsymbol{x}_T).$$

Note that under event $\mathcal{A}$, $\boldsymbol{A}_T^\mathsf{T}\boldsymbol{A}_T$ is invertible, and

$$\left\|(\boldsymbol{A}_T^\mathsf{T}\boldsymbol{A}_T)^{-1}\right\| \leq \frac{2}{m}.$$

We will show the above choice is a valid dual certificate.

# Validation of the dual certificate

The only condition that needs extra work is to establish

$$|u_i| < 1, \quad \forall i \in T^c.$$

This amounts to bound

$$\max_{i \in T^c} |u_i| = \max_{i \in T^c} \left| \left\langle \boldsymbol{a}_i, \underbrace{\boldsymbol{A}_T(\boldsymbol{A}_T^\mathsf{T}\boldsymbol{A}_T)^{-1}\mathsf{sgn}(\boldsymbol{x}_T)}_{\boldsymbol{w}} \right\rangle \right|$$

where $\boldsymbol{a}_i$ is the $i$th column of $\boldsymbol{A}$.

Note that $\boldsymbol{a}_i$ and $\boldsymbol{w}$ are independent for $i \in T^c$. For a fixed index $i \in T^c$,

- Conditioned on $\boldsymbol{w}$, $u_i \sim \mathcal{N}(0, \|\boldsymbol{w}\|_2^2)$, we have the Chernoff bound for the tail of a Gaussian rv:

$$\mathbb{P}(|u_i| \geq 1|\boldsymbol{w}) \leq 2\exp\left(-\frac{1}{2\|\boldsymbol{w}\|_2^2}\right)$$

- Under the event $\mathcal{A}$, we could also bound $\|\boldsymbol{w}\|_2$ as

$$
\begin{aligned}
\|\boldsymbol{w}\|_2 &\leq \|\boldsymbol{A}_T(\boldsymbol{A}_T^\mathsf{T}\boldsymbol{A}_T)^{-1}\| \cdot \|\mathsf{sgn}(\boldsymbol{x}_T)\|_2 \\
&\leq \|(\boldsymbol{A}_T^\mathsf{T}\boldsymbol{A}_T)^{-1}\|^{1/2} \cdot \|\mathsf{sgn}(\boldsymbol{x}_T)\|_2 \\
&\leq \sqrt{\frac{2k}{m}}
\end{aligned}
$$

since $(\boldsymbol{A}_T(\boldsymbol{A}_T^\mathsf{T}\boldsymbol{A}_T)^{-1})^\mathsf{T}\boldsymbol{A}_T(\boldsymbol{A}_T^\mathsf{T}\boldsymbol{A}_T)^{-1} = (\boldsymbol{A}_T^\mathsf{T}\boldsymbol{A}_T)^{-1}$.

We have

$$
\begin{aligned}
\mathbb{P}(\max_{i \in T^c} |u_i| \geq 1) &\leq |T^c| \cdot \mathbb{P}(|u_i| > 1) \quad \text{union bound} \\
&\leq n \int_{\boldsymbol{w}} \mathbb{P}(|u_i| \geq 1 | \boldsymbol{w}) d\mu(\boldsymbol{w}).
\end{aligned}
$$

Note that

$$\int_{\boldsymbol{w}} \mathbb{P}(|u_i| \geq 1|\boldsymbol{w})d\mu(\boldsymbol{w})$$

$$= \int_{\|\boldsymbol{w}\|_2 \leq \sqrt{\frac{2k}{m}}} \mathbb{P}(|u_i| \geq 1|\boldsymbol{w})d\mu(\boldsymbol{w}) + \int_{\|\boldsymbol{w}\|_2 > \sqrt{\frac{2k}{m}}} \mathbb{P}(|u_i| \geq 1|\boldsymbol{w})d\mu(\boldsymbol{w})$$

$$\leq \int_{\|\boldsymbol{w}\|_2 \leq \sqrt{\frac{2k}{m}}} \mathbb{P}(|u_i| \geq 1|\boldsymbol{w})d\mu(\boldsymbol{w}) + \mathbb{P}\left(\|\boldsymbol{w}\|_2 > \sqrt{\frac{2k}{m}}\right)$$

$$\leq \int_{\|\boldsymbol{w}\|_2 \leq \sqrt{\frac{2k}{m}}} 2e^{-\frac{1}{2\|\boldsymbol{w}\|_2^2}}d\mu(\boldsymbol{w}) + \mathbb{P}\left(\mathcal{A}^c\right)$$

$$\leq 2e^{-\frac{m}{4k}} + e^{-c_2 m} \leq 3e^{-\frac{m}{4k}},$$

which gives

$$\mathbb{P}(\max_{i \in T^c} |u_i| \geq 1) \leq 3ne^{-\frac{m}{4k}}$$

Set $m = 4(\gamma + 1)k \log n$ for some $\gamma > 0$, we have $\|\boldsymbol{u}_{T^c}\|_\infty < 1$ with probability at least $1 - n^{-\gamma}$.

# A General RIPless Theory for CS

Approach: Consider a sampling model that allows dependent entries across the row entries.

Let the signal be denoted as $\boldsymbol{x} \in \mathbb{R}^n$. The $i$th measurement is given as

$$y_i = \langle \boldsymbol{a}_i, \boldsymbol{x} \rangle, \quad i = 1, \ldots, m,$$

where each sampling/measurement vector is drawn from a distribution $F$:

$$\boldsymbol{a}_i \sim F, \quad \text{i.i.d.}$$

We will make a few assumptions on F such that it provides the incoherent sampling we want.

# Incoherence sampling

We define $F$ to satisfy two key properties:

- *Isometry property:* for $\boldsymbol{a} \sim F$,

$$\mathbb{E}\boldsymbol{a}\boldsymbol{a}^\mathsf{T} = \boldsymbol{I}$$

- *Incoherence property:* we let $\mu$ to be the smallest number such at $\boldsymbol{a} = [a_1, \ldots, a_n]^\mathsf{T} \sim F$,

$$\max_{1 \leq i \leq n} |a_i|^2 \leq \mu.$$

Remark:

- Both conditions may be relaxed a little, see [Candes and Plan, 2010]. In particular, we could allow the incoherence property holds with high probability, to accommodate the case $\boldsymbol{a} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$.

- $\mu \geq 1$ since $\mathbb{E}|a_i|^2 = 1$. On the other hand, $\mu$ could be as large as $n$. To get good performance, we would like to have $\mu$ small.

# Examples of incoherence sampling



concentrated vector          incoherent measurements

- Denote the (scaled) DFT matrix $\boldsymbol{\Phi}$ with entries $\phi_{l,i} = e^{j2\pi li/n}$. Let $\boldsymbol{a} \sim F$ be obtained by sampling a row of $\boldsymbol{\Phi}$ uniformly at random, we have

$$\mathbb{E}[\boldsymbol{a}\boldsymbol{a}^{\mathsf{H}}] = \sum_{l=1}^{n} \frac{1}{n}\phi_l^{\mathsf{H}}\phi_l = \boldsymbol{I}$$

and $\max_i |a_i|^2 = 1 := \mu$ .

# Other Examples of incoherent sampling

- Binary sensing: $\mathbb{P}(a_i = \pm 1) = \frac{1}{2}$,

$$\mathbb{E}[\boldsymbol{a}\boldsymbol{a}^\mathsf{T}] = \boldsymbol{I}, \qquad \max_i |a_i|^2 = 1.$$

- Gaussian sensing: $a_i \sim \mathcal{N}(0, 1)$, we have

$$\mathbb{E}[\boldsymbol{a}\boldsymbol{a}^\mathsf{T}] = \boldsymbol{I}, \qquad \max_i |a_i|^2 \approx 2 \log n.$$

- Partial Fourier transform (useful in MRI): pick a frequency $\omega \sim \mathsf{Unif}[0, 1]$, and set $a_i = e^{j2\pi\omega i}$. We have

$$\mathbb{E}[\boldsymbol{a}\boldsymbol{a}^\mathsf{T}] = \boldsymbol{I}, \qquad \max_i |a_i|^2 = 1.$$

# Performance Guarantees for BP

**Theorem 2. [Noise-free, Basis Pursuit]** *Let $x \in \mathbb{R}^n$ be an arbitrary* fixed *vector that is $k$-sparse. Then $x$ is the unique solution to BP with high probability, as long as*

$$m \geq C\mu k \log n$$

*for some constant $C$.*

- The proof is based on slightly different methods, by constructing an *inexact* dual certificate using the *golfing scheme*. This technique is developed first by D. Gross for analyzing matrix completion. We will discuss this method later in the course in more details.

- The result is near-optimal for the general class of incoherence sampling models. It is clear that the "oversampling ratio" depends on the coherence parameter $\mu$.

- When specializing to the Gaussian case, this result is sub-optimal by $\log n$.

# Performance Guarantees for the General Case using LASSO

Consider noisy observation with Gaussian noise:

$$y = Ax + w$$

where $w \sim \mathcal{N}(0, \sigma^2 I)$. Consider the LASSO algorithm:

$$\hat{x} = \underset{x}{\operatorname{argmin}} \frac{1}{2}\|y - Ax\|_2^2 + \lambda\|x\|_1$$

**Theorem 3. [Candes and Plan, 2010]** *Set $\lambda = 10\sigma\sqrt{\log n}$. Then with high probability, we have*

$$\|\hat{x} - x\|_2 \lesssim \frac{\|x - x_k\|_1}{\sqrt{k}} + \sigma\sqrt{\frac{k \log n}{m}}$$

*provided $m \gtrsim \mu k \log n$.*