

# ECE 18-898G: Special Topics in Signal Processing: Sparsity, Structure, and Inference

Sparse Recovery using L1 minimization

Yuejie Chi

Department of Electrical and Computer Engineering

**Carnegie Mellon University**

Spring 2018

# Outline

---

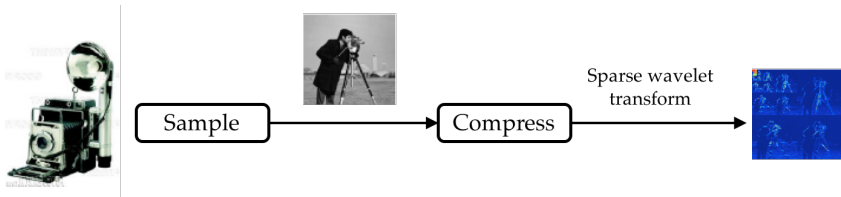
- $\ell_1$  minimization for sparse recovery
- Restricted isometry property (RIP)
- A RIPless theory

# Motivation of Compressed Sensing

---

Conventional paradigms for data acquisition:

- Measure full data
- Compress (by discarding a large fraction of coefficients)



**Problem:** data is often highly compressible

- Most of acquired data can be thrown away without any perceptual loss

# Blind sensing

---

Ideally, if we know *a priori* which coefficients are worth estimating, then we can simply measure these coefficients

- Unfortunately, we often have no idea which coefficients are most relevant

## **Compressed sensing: compression on the fly**

- mimic the behavior of the above ideal situation without pre-computing all coefficients
- often achieved by *random* sensing mechanism

*Why go to so much effort to acquire all the data when most of what we get will be thrown away?*

*Can't we just directly measure the part that won't end up being thrown away?*

— David Donoho

## Setup: sparse recovery

---

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

Recover  $\mathbf{x} \in \mathbb{R}^p$  given  $\mathbf{y} = \mathbf{A}\mathbf{x}$

where  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n]^\top \in \mathbb{R}^{n \times p}$  ( $n \ll p$ ): sampling matrix;  
 $\mathbf{a}_i$ : sampling vector;  $\mathbf{x}$ : sparse signal

# Restricted isometry properties

# Optimality for $\ell_0$ minimization

---

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^p} \|\mathbf{x}\|_0 \quad \text{s.t. } \mathbf{A}\mathbf{x} = \mathbf{y}$$

If instead  $\exists$  a sparser feasible  $\tilde{\mathbf{x}} \neq \mathbf{x}$  s.t.  $\|\tilde{\mathbf{x}}\|_0 \leq \|\mathbf{x}\|_0 = k$ , then

$$\mathbf{A}(\mathbf{x} - \tilde{\mathbf{x}}) = \mathbf{0}. \quad (6.1)$$

We don't want (6.1) to happen, so we hope

$$\mathbf{A}(\underbrace{\mathbf{x} - \tilde{\mathbf{x}}}_{2k\text{-sparse}}) \neq \mathbf{0}, \quad \forall \tilde{\mathbf{x}} \text{ with } \|\tilde{\mathbf{x}}\|_0 \leq k$$

To simultaneously account for all  $k$ -sparse  $\mathbf{x}$ , we hope  $\mathbf{A}_T$  ( $|T| \leq 2k$ ) to have full column rank, where  $\mathbf{A}_T$  consists of all columns of  $\mathbf{A}$  at indices from  $T$



# Restricted isometry property (RIP)

## Definition 6.1 (Restricted isometry constant)

Restricted isometry constant  $\delta_k$  of  $\mathbf{A}$  is smallest quantity s.t.

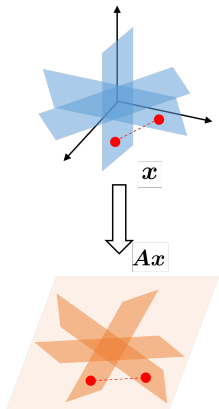
$$(1 - \delta_k)\|\mathbf{x}\|^2 \leq \|\mathbf{A}\mathbf{x}\|^2 \leq (1 + \delta_k)\|\mathbf{x}\|^2 \quad (6.2)$$

holds for all  $k$ -sparse vector  $\mathbf{x} \in \mathbb{R}^p$

- Equivalently, (6.2) says

$$\max_{S:|S|=k} \underbrace{\|\mathbf{A}_S^\top \mathbf{A}_S - \mathbf{I}_k\|}_{\text{near orthonormality}} = \delta_k$$

where  $\mathbf{A}_S$  consists of all columns of  $\mathbf{A}$  at indices from  $S$



# Restricted isometry property (RIP)

## Definition 6.1 (Restricted isometry constant)

Restricted isometry constant  $\delta_k$  of  $\mathbf{A}$  is smallest quantity s.t.

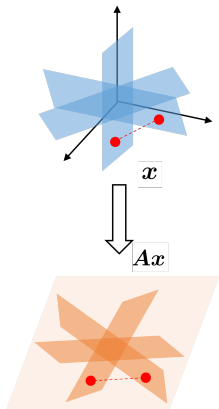
$$(1 - \delta_k)\|\mathbf{x}\|^2 \leq \|\mathbf{Ax}\|^2 \leq (1 + \delta_k)\|\mathbf{x}\|^2 \quad (6.2)$$

holds for all  $k$ -sparse vector  $\mathbf{x} \in \mathbb{R}^p$

- (Homework) For any  $\mathbf{x}_1, \mathbf{x}_2$  that are supported on disjoint subsets  $S_1, S_2$  with  $|S_1| \leq s_1$  and  $|S_2| \leq s_2$ :

$$|\langle \mathbf{Ax}_1, \mathbf{Ax}_2 \rangle| \leq \delta_{s_1+s_2} \|\mathbf{x}_1\|_2 \|\mathbf{x}_2\|_2 \quad (6.3)$$

angle-preserving! (consequence of parallelogram identity)



# RIP and $\ell_0$ minimization

---

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^p} \|\mathbf{x}\|_0 \quad \text{s.t. } \mathbf{Ax} = \mathbf{y}$$

## Fact 6.2

*Suppose  $\mathbf{x}$  is  $k$ -sparse. If  $\delta_{2k} < 1$ , then  $\ell_0$ -minimization is exact and unique.*

# RIP and $\ell_1$ minimization

---

## Theorem 6.3 (Candès 2008)

Suppose  $x$  is  $k$ -sparse. If  $\delta_{2k} < \sqrt{2} - 1$ , then  $\ell_1$ -minimization is exact and unique.

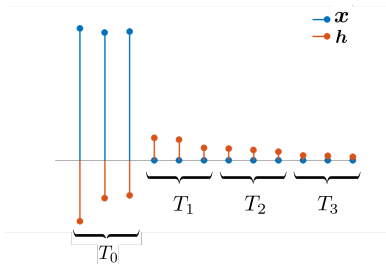
- RIP implies success of  $\ell_1$  minimization (also many other methods, as we'll see from later lectures)
- A universal result: works simultaneously for **all**  $k$ -sparse signals
- As we will see later, many random designs satisfy this condition with *near-optimal sample complexity*

$$m \sim O(k \log(n/k))$$

## Proof of Theorem 6.3

---

Suppose  $\mathbf{x} + \mathbf{h}$  is feasible and obeys  $\|\mathbf{x} + \mathbf{h}\|_1 \leq \|\mathbf{x}\|_1$ . The goal is to show that  $\mathbf{h} = \mathbf{0}$  under RIP.



The key is to decompose  $\mathbf{h}$  into  $\mathbf{h}_{T_0} + \mathbf{h}_{T_1} + \dots$

- $T_0$ : locations of  $k$  largest entries of  $\mathbf{x}$
- $T_1$ : locations of  $k$  largest entries of  $\mathbf{h}$  in  $T_0^c$
- $T_2$ : locations of  $k$  largest entries of  $\mathbf{h}$  in  $(T_0 \cup T_1)^c$
- ...

## Proof of Theorem 6.3

---

The proof proceeds by showing that

1.  $\mathbf{h}_{T_0 \cup T_1}$  dominates  $\mathbf{h}_{(T_0 \cup T_1)^c}$  (by objective function)
2. (converse)  $\mathbf{h}_{(T_0 \cup T_1)^c}$  dominates  $\mathbf{h}_{T_0 \cup T_1}$  (by RIP + feasibility)

These can happen simultaneously only when  $\mathbf{h}$  vanishes

## Proof of Theorem 6.3

**Step 1 (depending only on objective function).** Show that

$$\sum_{j \geq 2} \|\mathbf{h}_{T_j}\| \leq \frac{1}{\sqrt{k}} \|\mathbf{h}_{T_0}\|_1. \quad (6.4)$$

This follows immediately by combining the following 2 observations:

(i) Since  $\mathbf{x} + \mathbf{h}$  is assumed to be a better estimate:

$$\begin{aligned} \|\mathbf{x}\|_1 &\geq \|\mathbf{x} + \mathbf{h}\|_1 = \underbrace{\|\mathbf{x} + \mathbf{h}_{T_0}\|_1 + \|\mathbf{h}_{T_0^c}\|_1}_{\text{since } T_0 \text{ is support of } \mathbf{x}} \geq \underbrace{\|\mathbf{x}\|_1 - \|\mathbf{h}_{T_0}\|_1}_{\text{triangle inequality}} + \|\mathbf{h}_{T_0^c}\|_1 \\ &\implies \|\mathbf{h}_{T_0^c}\|_1 \leq \|\mathbf{h}_{T_0}\|_1 \end{aligned} \quad (6.5)$$

(ii) Since entries of  $\mathbf{h}_{T_{j-1}}$  uniformly dominate those of  $\mathbf{h}_{T_j}$  ( $j \geq 2$ ):

$$\begin{aligned} \|\mathbf{h}_{T_j}\| &\leq \sqrt{k} \|\mathbf{h}_{T_j}\|_\infty \leq \sqrt{k} \frac{\|\mathbf{h}_{T_{j-1}}\|_1}{k} = \frac{1}{\sqrt{k}} \|\mathbf{h}_{T_{j-1}}\|_1 \\ \implies \sum_{j \geq 2} \|\mathbf{h}_{T_j}\| &\leq \frac{1}{\sqrt{k}} \sum_{j \geq 2} \|\mathbf{h}_{T_{j-1}}\|_1 = \frac{1}{\sqrt{k}} \|\mathbf{h}_{T_0^c}\|_1 \end{aligned} \quad (6.6)$$

## Proof of Theorem 6.3

---

**Step 2 (using feasibility + RIP).** Show that  $\exists \rho < 1$  s.t.

$$\|\mathbf{h}_{T_0 \cup T_1}\| \leq \rho \sum_{j \geq 2} \|\mathbf{h}_{T_j}\| \quad (6.7)$$

If this claim holds, then

$$\begin{aligned} \|\mathbf{h}_{T_0 \cup T_1}\| &\leq \rho \sum_{j \geq 2} \|\mathbf{h}_{T_j}\| \stackrel{(6.4)}{\leq} \rho \frac{1}{\sqrt{k}} \|\mathbf{h}_{T_0}\|_1 \\ &\leq \rho \frac{1}{\sqrt{k}} \left( \sqrt{k} \|\mathbf{h}_{T_0}\| \right) = \rho \|\mathbf{h}_{T_0}\| \leq \rho \|\mathbf{h}_{T_0 \cup T_1}\|. \end{aligned} \quad (6.8)$$

Since  $\rho < 1$ , we necessarily have  $\mathbf{h}_{T_0 \cup T_1} = \mathbf{0}$ , which together with (6.5) yields  $\mathbf{h} = \mathbf{0}$



## Proof of Theorem 6.3

---

We now prove (6.7). To connect  $\mathbf{h}_{T_0 \cup T_1}$  with  $\mathbf{h}_{(T_0 \cup T_1)^c}$ , we use feasibility:

$$\mathbf{A}\mathbf{h} = \mathbf{0} \iff \mathbf{A}\mathbf{h}_{T_0 \cup T_1} = -\sum_{j \geq 2} \mathbf{A}\mathbf{h}_{T_j},$$

which taken collectively with RIP yields

$$(1 - \delta_{2k}) \|\mathbf{h}_{T_0 \cup T_1}\|^2 \leq \|\mathbf{A}\mathbf{h}_{T_0 \cup T_1}\|^2 = \left| \langle \mathbf{A}\mathbf{h}_{T_0 \cup T_1}, \sum_{j \geq 2} \mathbf{A}\mathbf{h}_{T_j} \rangle \right|.$$

It follows from (6.3) that for all  $j \geq 2$ ,

$$\begin{aligned} |\langle \mathbf{A}\mathbf{h}_{T_0 \cup T_1}, \mathbf{A}\mathbf{h}_{T_j} \rangle| &\leq |\langle \mathbf{A}\mathbf{h}_{T_0}, \mathbf{A}\mathbf{h}_{T_j} \rangle| + |\langle \mathbf{A}\mathbf{h}_{T_1}, \mathbf{A}\mathbf{h}_{T_j} \rangle| \\ &\stackrel{(6.3)}{\leq} \delta_{2k} (\|\mathbf{h}_{T_0}\| + \|\mathbf{h}_{T_1}\|) \|\mathbf{h}_{T_j}\| \leq \delta_{2k} \sqrt{2} \|\mathbf{h}_{T_0 \cup T_1}\| \cdot \|\mathbf{h}_{T_j}\|, \end{aligned}$$

which gives

$$\begin{aligned} (1 - \delta_{2k}) \|\mathbf{h}_{T_0 \cup T_1}\|^2 &\leq \sum_{j \geq 2} |\langle \mathbf{A}\mathbf{h}_{T_0 \cup T_1}, \mathbf{A}\mathbf{h}_{T_j} \rangle| \\ &\leq \sqrt{2} \delta_{2k} \|\mathbf{h}_{T_0 \cup T_1}\| \sum_{j \geq 2} \|\mathbf{h}_{T_j}\| \end{aligned}$$

This establishes (6.7) if  $\rho := \frac{\sqrt{2}\delta_{2k}}{1-\delta_{2k}} < 1$  (or equivalently,  $\delta_{2k} < \sqrt{2} - 1$ ).

# Robustness for compressible signals

## Theorem 6.4

If  $\delta_{2k} < \sqrt{2} - 1$ , then the solution  $\hat{\mathbf{x}}$  to  $\ell_1$ -minimization obeys

$$\|\hat{\mathbf{x}} - \mathbf{x}\| \lesssim \frac{\|\mathbf{x} - \mathbf{x}_k\|_1}{\sqrt{k}},$$

where  $\mathbf{x}_k$  is best  $k$ -term approximation of  $\mathbf{x}$

- Suppose  $l^{\text{th}}$  largest entry of  $\mathbf{x}$  is  $1/l^\alpha$  for some  $\alpha > 1$ , then

$$\frac{1}{\sqrt{k}} \|\mathbf{x} - \mathbf{x}_k\|_1 \approx \frac{1}{\sqrt{k}} \sum_{l>k} l^{-\alpha} \approx k^{-\alpha+0.5} \ll 1$$

- $\ell_1$ -min works well in recovering compressible signals
- Follows similar arguments as in proof of Theorem 6.3

## Proof of Theorem 6.4

---

**Step 1 (depending only on objective function).** Show that

$$\sum_{j \geq 2} \|\mathbf{h}_{T_j}\| \leq \frac{1}{\sqrt{k}} \|\mathbf{h}_{T_0}\|_1 + \frac{2}{\sqrt{k}} \|\mathbf{x} - \mathbf{x}_{T_0}\|_1. \quad (6.9)$$

This follows immediately by combining the following 2 observations:

(i) Since  $\mathbf{x} + \mathbf{h}$  is assumed to be a better estimate:

$$\begin{aligned} \|\mathbf{x}_{T_0}\|_1 + \|\mathbf{x}_{T_0^c}\|_1 &= \|\mathbf{x}\|_1 \geq \|\mathbf{x} + \mathbf{h}\|_1 = \|\mathbf{x}_{T_0} + \mathbf{h}_{T_0}\|_1 + \|\mathbf{x}_{T_0^c} + \mathbf{h}_{T_0^c}\|_1 \\ &\geq \|\mathbf{x}_{T_0}\|_1 - \|\mathbf{h}_{T_0}\|_1 + \|\mathbf{h}_{T_0^c}\|_1 - \|\mathbf{x}_{T_0^c}\|_1 \\ \implies \|\mathbf{h}_{T_0^c}\|_1 &\leq \|\mathbf{h}_{T_0}\|_1 + 2\|\mathbf{x}_{T_0^c}\|_1 \end{aligned} \quad (6.10)$$

(ii) Recall from (6.6) that  $\sum_{j \geq 2} \|\mathbf{h}_{T_j}\| \leq \frac{1}{\sqrt{k}} \|\mathbf{h}_{T_0^c}\|_1$ .

## Proof of Theorem 6.4

---

**Step 2 (using feasibility + RIP).** Recall from (6.7) that  $\exists \rho < 1$  s.t.

$$\|\mathbf{h}_{T_0 \cup T_1}\| \leq \rho \sum_{j \geq 2} \|\mathbf{h}_{T_j}\| \quad (6.11)$$

If this claim holds, then

$$\begin{aligned} \|\mathbf{h}_{T_0 \cup T_1}\| &\leq \rho \sum_{j \geq 2} \|\mathbf{h}_{T_j}\| \stackrel{\text{(6.10) and (6.6)}}{\leq} \rho \frac{1}{\sqrt{k}} \{ \|\mathbf{h}_{T_0}\|_1 + 2\|\mathbf{x}_{T_0^c}\|_1 \} \\ &\leq \rho \frac{1}{\sqrt{k}} \left( \sqrt{k} \|\mathbf{h}_{T_0}\| + 2\|\mathbf{x}_{T_0^c}\|_1 \right) = \rho \|\mathbf{h}_{T_0}\| + \frac{2\rho}{\sqrt{k}} \|\mathbf{x}_{T_0^c}\|_1 \\ &\leq \rho \|\mathbf{h}_{T_0 \cup T_1}\| + \frac{2\rho}{\sqrt{k}} \|\mathbf{x}_{T_0^c}\|_1. \\ \implies \quad \|\mathbf{h}_{T_0 \cup T_1}\| &\leq \frac{2\rho}{1-\rho} \frac{\|\mathbf{x}_{T_0^c}\|_1}{\sqrt{k}}. \end{aligned} \quad (6.12)$$

---

We highlight in red the part different from the proof of Theorem 6.3.

## Proof of Theorem 6.4

---

Finally, putting the above together yields

$$\begin{aligned}\|\mathbf{h}\| &\leq \|\mathbf{h}_{T_0 \cup T_1}\| + \|\mathbf{h}_{(T_0 \cup T_1)^c}\| \\ &\stackrel{(6.9)}{\leq} \|\mathbf{h}_{T_0 \cup T_1}\| + \frac{1}{\sqrt{k}} \|\mathbf{h}_{T_0}\|_1 + \frac{2}{\sqrt{k}} \|\mathbf{x} - \mathbf{x}_{T_0}\|_1 \\ &\leq \|\mathbf{h}_{T_0 \cup T_1}\| + \|\mathbf{h}_{T_0}\| + \frac{2}{\sqrt{k}} \|\mathbf{x} - \mathbf{x}_{T_0}\|_1 \\ &\leq 2\|\mathbf{h}_{T_0 \cup T_1}\| + \frac{2}{\sqrt{k}} \|\mathbf{x} - \mathbf{x}_{T_0}\|_1 \\ &\stackrel{(6.12)}{\leq} \frac{2(1+\rho)}{1-\rho} \frac{\|\mathbf{x} - \mathbf{x}_{T_0}\|_1}{\sqrt{k}}\end{aligned}$$

## $\ell_1$ recovery in the noisy case

---

In the presence of additive measurement noise,

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w},$$

where  $\|\mathbf{w}\|_2 \leq \epsilon$  is assumed to be bounded.

We can modify the BP algorithm in the following manner:

$$\text{(BP-noisy:)} \quad \hat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x}} \|\mathbf{x}\|_1 \quad \text{subject to} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2 \leq \epsilon.$$

### Theorem 6.5 (Performance of BP via RIP, noisy case)

*If  $\delta_{2k} < \sqrt{2} - 1$ , then for any vector  $\mathbf{x}$ , the solution to basis pursuit (noisy case) satisfies*

$$\|\hat{\mathbf{x}} - \mathbf{x}\|_2 \leq C_0 k^{-1/2} \|\mathbf{x} - \mathbf{x}_k\|_1 + C_1 \epsilon.$$

*where  $\mathbf{x}_k$  is the best  $k$ -term approximation of  $\mathbf{x}$  for some constants  $C_0$  and  $C_1$ .*

## Proof of Theorem 6.5

---

Again let's start by assuming  $\hat{\mathbf{x}} = \mathbf{x} + \mathbf{h}$ . The key difference from the noiseless case is that in Step 2, we now have

$$\begin{aligned}\|\mathbf{A}\mathbf{h}\|_2 &= \|\mathbf{A}(\hat{\mathbf{x}} - \mathbf{x})\|_2 = \|(\mathbf{y} - \mathbf{A}\hat{\mathbf{x}}) - (\mathbf{y} - \mathbf{A}\mathbf{x})\|_2 \\ &\leq \|\mathbf{y} - \mathbf{A}\hat{\mathbf{x}}\|_2 + \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2 \leq 2\epsilon.\end{aligned}$$

Therefore, we need to bound

$$\begin{aligned}\|\mathbf{A}\mathbf{h}_{T_0 \cup T_1}\|_2^2 &= \langle \mathbf{A}\mathbf{h} - \sum_{j \geq 2} \mathbf{A}\mathbf{h}_{T_j}, \mathbf{A}\mathbf{h}_{T_0 \cup T_1} \rangle \\ &\leq \underbrace{\langle \mathbf{A}\mathbf{h}, \mathbf{A}\mathbf{h}_{T_0 \cup T_1} \rangle}_{\leq 2\epsilon\delta_{2k}\|\mathbf{h}_{T_0 \cup T_1}\|_2} - \underbrace{\sum_{j \geq 2} \langle \mathbf{A}\mathbf{h}_{T_j}, \mathbf{A}\mathbf{h}_{T_0 \cup T_1} \rangle}_{\text{bounded as before}}\end{aligned}$$

By plugging in this modification, we show

$$\|\hat{\mathbf{x}} - \mathbf{x}\|_2 = \|\mathbf{h}\|_2 \leq \frac{2(1+\rho)}{1-\rho} \frac{\|\mathbf{x} - \mathbf{x}_k\|_1}{\sqrt{k}} + \frac{2\alpha}{1-\rho}\epsilon,$$

where  $\alpha = \frac{2\sqrt{1+\delta_{2k}}}{1-\delta_{2k}}$ .

# Which design matrix satisfies RIP?

---

**First example:** i.i.d. Gaussian design

## Lemma 6.6

A random matrix  $\mathbf{A} \in \mathbb{R}^{n \times p}$  with i.i.d.  $\mathcal{N}\left(0, \frac{1}{n}\right)$  entries satisfies  $\delta_k < \delta$  with high prob., as long as

$$n \gtrsim \frac{1}{\delta^2} k \log \frac{p}{k}$$

- This is where non-asymptotic random matrix theory enters



# Gaussian random matrices

---

## Lemma 6.7 (See Vershynin '10)

Suppose  $\mathbf{B} \in \mathbb{R}^{n \times k}$  is composed of i.i.d.  $\mathcal{N}(0, 1)$  entries. Then

$$\begin{cases} \mathbb{P}\left(\frac{1}{\sqrt{n}}\sigma_{\max}(\mathbf{B}) > 1 + \sqrt{\frac{k}{n}} + t\right) & \leq e^{-nt^2/2} \\ \mathbb{P}\left(\frac{1}{\sqrt{n}}\sigma_{\min}(\mathbf{B}) < 1 - \sqrt{\frac{k}{n}} - t\right) & \leq e^{-nt^2/2}. \end{cases}$$

- When  $n \gg k$ , one has  $\frac{1}{n}\mathbf{B}^\top \mathbf{B} \approx \mathbf{I}_k$
- Similar results (up to different constants) hold for i.i.d. sub-Gaussian matrix

## Proof of Lemma 6.6

---

1. Fix any index subset  $S \subseteq \{1, \dots, p\}$ ,  $|S| = k$ , then  $\mathbf{A}_S$  (submatrix of  $\mathbf{A}$  consisting of columns at indices from  $S$ ) obeys

$$\left\| \mathbf{A}_S^\top \mathbf{A}_S - \mathbf{I}_k \right\| \leq O\left(\sqrt{k/n}\right) + t$$

with prob. exceeding  $1 - 2e^{-c_1 n t^2}$ , where  $c_1 > 0$  is constant.

2. Taking a union bound over all  $S \subseteq \{1, \dots, p\}$ ,  $|S| = k$  yields

$$\delta_k = \max_{S:|S|=k} \left\| \mathbf{A}_S^\top \mathbf{A}_S - \mathbf{I}_k \right\| \leq O\left(\sqrt{k/n}\right) + t$$

with prob. exceeding  $1 - 2\binom{p}{k} e^{-c_1 n t^2} \geq 1 - 2e^{k \log(ep/k) - c_1 n t^2}$ .

Thus,  $\delta_k < \delta$  with high prob. as long as  $n \gtrsim \delta^{-2} k \log(p/k)$ .

## Other design matrices that satisfy RIP

---

- Random matrices with i.i.d. **sub-Gaussian** entries, as long as

$$n \gtrsim k \log(p/k)$$

- Random partial DFT matrices with

$$n \gtrsim k \log^4 p,$$

where rows of  $\mathbf{A}$  are independently sampled from rows of DFT matrix  $\mathbf{F}$  (Rudelson & Vershynin '08)

- If you have learned entropy method / generic chaining, check out Rudelson & Vershynin '08 and Candes & Plan '11

# Other design matrices that satisfy RIP

---

- Random convolution matrices with

$$n \gtrsim k \log^4 p,$$

where rows of  $\mathbf{A}$  are independently sampled from rows of

$$\mathbf{G} = \begin{bmatrix} g_0 & g_1 & g_2 & \cdots & g_{p-1} \\ g_{p-1} & g_0 & g_1 & \cdots & g_{p-2} \\ g_{p-2} & g_{p-1} & g_0 & \cdots & g_{p-3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ g_1 & g_2 & g_3 & \cdots & g_0 \end{bmatrix}$$

with  $\mathbb{P}(g_i = \pm 1) = 0.5$  (Krahmer, Mendelson, & Rauhut '14)

## **A RIPlless theory**

# Is RIP necessary?

---

- RIP leads to a **universal** result holding simultaneously for all  $k$ -sparse  $x$ 
  - Universality is often not needed as we might only care about a particular  $x$
- There may be a gap between the regime where RIP holds and the regime in which one has minimal measurements
- Certifying RIP is hard

Can we develop a non-universal RIPless theory?

# A standard recipe

---

1. Write out Karush-Kuhn-Tucker (KKT) optimality conditions
  - typically involve certain dual variables
2. Construct dual variables satisfying KKT conditions

# Karush-Kuhn-Tucker (KKT) condition

---

Consider a convex problem

$$\begin{array}{ll} \text{minimize}_x & f(\mathbf{x}) \\ \text{s.t.} & \mathbf{Ax} - \mathbf{y} = \mathbf{0} \end{array}$$

Lagrangian:

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\nu}) := f(\mathbf{x}) + \boldsymbol{\nu}^\top (\mathbf{Ax} - \mathbf{y}) \quad (\boldsymbol{\nu} : \text{Lagrangian multiplier})$$

If  $\mathbf{x}$  is optimizer, then KKT optimality condition reads

$$\begin{cases} \mathbf{0} = \nabla_{\mathbf{v}} \mathcal{L}(\mathbf{x}, \mathbf{v}) \\ \mathbf{0} \in \partial_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{v}) \end{cases}$$



# Karush-Kuhn-Tucker (KKT) condition

---

Consider a convex problem

$$\begin{array}{ll} \text{minimize}_x & f(x) \\ \text{s.t.} & \mathbf{Ax} - \mathbf{y} = \mathbf{0} \end{array}$$

Lagrangian:

$$\mathcal{L}(x, \nu) := f(x) + \nu^\top (\mathbf{Ax} - \mathbf{y}) \quad (\nu : \text{Lagrangian multiplier})$$

If  $x$  is optimizer, then KKT optimality condition reads

$$\begin{cases} \mathbf{Ax} - \mathbf{y} = \mathbf{0} \\ \mathbf{0} \in \partial f(x) + \mathbf{A}^\top \nu \end{cases} \quad (\text{no constraint on } \nu)$$

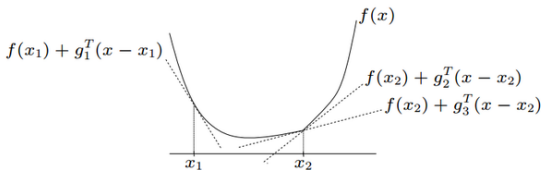
# Subgradient

Consider a convex function  $f(x)$  (possibly nonsmooth).

## Definition 6.8 (Subgradient)

$u \in \partial f(x_0)$  is a subgradient of a convex  $f$  at  $x_0$  if for all  $x$ :

$$f(x) \geq f(x_0) + u^T(x - x_0)$$



**Remark:** if  $f$  is differentiable at  $x_0$ , the only subgradient is the gradient  $\nabla f(x_0)$ .

## Subgradient of $\ell_1$ norm

---

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

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

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

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

# KKT condition for $\ell_1$ minimization

---

$$\begin{aligned} & \text{minimize}_x && \|x\|_1 \\ & \text{s.t.} && Ax - y = 0 \end{aligned}$$

If  $x$  is optimizer, then KKT optimality condition reads

$$\begin{cases} Ax - y = 0, & \text{(naturally satisfied as } x \text{ is truth)} \\ \mathbf{0} \in \partial\|x\|_1 + A^\top \nu & \text{(no constraint on } \nu) \end{cases}$$

$$\iff \exists u \in \text{range}(A^\top) \quad \text{s.t.} \quad \underbrace{\begin{cases} u_i = \text{sign}(x_i), & \text{if } x_i \neq 0 \\ u_i \in [-1, 1], & \text{else} \end{cases}}_{u \text{ is a valid subgradient}}$$

Depends only on signs of  $x_i$ 's irrespective of their magnitudes

# Uniqueness

## Theorem 6.9 (A sufficient—and almost necessary—condition)

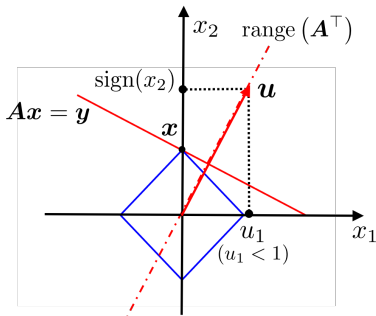
Let  $T := \text{supp}(x)$ . *Suppose  $A_T$  has full column rank. If*

$$\exists \mathbf{u} = \mathbf{A}^\top \boldsymbol{\nu} \text{ for some } \boldsymbol{\nu} \in \mathbb{R}^n \quad \text{s.t.} \quad \begin{cases} u_i = \text{sign}(x_i), & \text{if } x_i \neq 0 \\ u_i \in (-1, 1), & \text{else} \end{cases},$$

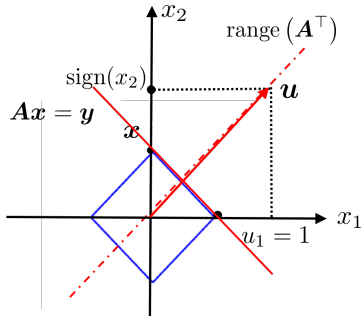
*then  $x$  is unique solution to  $\ell_1$  minimization.*

- Only slightly stronger than KKT!
- $\boldsymbol{\nu}$  is said to be a **dual certificate**
  - recall that  $\boldsymbol{\nu}$  is Lagrangian multiplier
- *Finding  $\boldsymbol{\nu}$  comes down to solving another convex problem*

# Geometric interpretation of dual certificate



When  $|u_1| < 1$ , solution is unique



When  $|u_1| = 1$ , solution is **non-unique**

When we are able to find  $\mathbf{u} \in \text{range}(\mathbf{A}^\top)$  s.t.  $u_2 = \text{sign}(x_2)$  and  $|u_1| < 1$ , then  $\mathbf{x}$  (with  $x_1 = 0$ ) is unique solution to  $\ell_1$  minimization

## Proof of Theorem 6.9

---

Let  $\mathbf{w} \in \partial\|\mathbf{x}\|_1$  be  $\begin{cases} w_i = \text{sign}(x_i), & \text{if } i \in T \text{ (support of } \mathbf{x}\text{);} \\ w_i = \text{sign}(h_i), & \text{else.} \end{cases}$  If  $\mathbf{x} + \mathbf{h}$  is optimizer with  $\mathbf{h}_{T^c} \neq \mathbf{0}$ , then

$$\begin{aligned} \|\mathbf{x}\|_1 \geq \|\mathbf{x} + \mathbf{h}\|_1 &\geq \|\mathbf{x}\|_1 + \langle \mathbf{w}, \mathbf{h} \rangle = \|\mathbf{x}\|_1 + \langle \mathbf{u}, \mathbf{h} \rangle + \langle \mathbf{w} - \mathbf{u}, \mathbf{h} \rangle \\ &= \|\mathbf{x}\|_1 + \langle \underbrace{\mathbf{A}^\top \boldsymbol{\nu}}_{\text{assumption on } \mathbf{u}}, \mathbf{h} \rangle + \sum_{i \notin T} (\text{sign}(h_i)h_i - u_i h_i) \\ &= \|\mathbf{x}\|_1 + \langle \boldsymbol{\nu}, \underbrace{\mathbf{A}\mathbf{h}}_{:=\mathbf{0} \text{ (feasibility)}} \rangle + \sum_{i \notin T} (|h_i| - u_i h_i) \\ &\geq \|\mathbf{x}\|_1 + \sum_{i \notin T} (1 - |u_i|) |h_i| > \|\mathbf{x}\|_1, \end{aligned}$$

resulting in contradiction.

Further, when  $\mathbf{h}_{T^c} = \mathbf{0}$ , one must have  $\mathbf{h}_T = \mathbf{0}$  from left-invertibility of  $\mathbf{A}_T$ , and hence  $\mathbf{h} = \mathbf{h}_T + \mathbf{h}_{T^c} = \mathbf{0}$

# Constructing dual certificates under Gaussian design

---

We illustrate how to construct dual certificates for the following setup

- $\mathbf{x} \in \mathbb{R}^p$  is  $k$ -sparse
- Entries of  $\mathbf{A} \in \mathbb{R}^{n \times p}$  are i.i.d. standard Gaussian
- Sample size  $n$  obeys

$$n \gtrsim k \log p$$



# Constructing dual certificates under Gaussian design

---

$$\begin{aligned} \text{Find} \quad & \boldsymbol{\nu} \in \mathbb{R}^n \\ \text{s.t.} \quad & (\mathbf{A}^\top \boldsymbol{\nu})_T = \text{sign}(\mathbf{x}_T) \end{aligned} \quad (6.13)$$

$$|(\mathbf{A}^\top \boldsymbol{\nu})_i| < 1, \quad i \notin T \quad (6.14)$$

**Step 1:** propose a  $\boldsymbol{\nu}$  compatible with linear constraints (6.13). One candidate is *least squares* solution:

$$\boldsymbol{\nu} = \mathbf{A}_T (\mathbf{A}_T^\top \mathbf{A}_T)^{-1} \text{sign}(\mathbf{x}_T) \quad (\text{explicit expression})$$

- LS solution minimizes  $\|\boldsymbol{\nu}\|$ , which will also be helpful when controlling  $|(\mathbf{A}^\top \boldsymbol{\nu})_i|$
- From Lemma 6.7,  $\mathbf{A}_T^\top \mathbf{A}_T$  is invertible when  $n \gtrsim k \log p$

# Constructing dual certificates under Gaussian design

Step 2: verify (6.14), which amounts to controlling

$$\max_{i \notin T} \left| \left\langle \underbrace{\mathbf{a}_i}_{\text{ith column of } \mathbf{A}}, \underbrace{\mathbf{A}_T(\mathbf{A}_T^\top \mathbf{A}_T)^{-1} \text{sign}(\mathbf{x}_T)}_{\boldsymbol{\nu}} \right\rangle \right|$$

- Since  $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$  and  $\boldsymbol{\nu}$  are **independent** for any  $i \notin T$ ,

$$\max_{i \notin T} |\langle \mathbf{a}_i, \boldsymbol{\nu} \rangle| \lesssim \|\boldsymbol{\nu}\| \sqrt{\log p}$$

- $\|\boldsymbol{\nu}\|$  can be bounded by

$$\begin{aligned} \|\boldsymbol{\nu}\| &\leq \|\mathbf{A}_T(\mathbf{A}_T^\top \mathbf{A}_T)^{-1}\| \cdot \|\text{sgn}(\mathbf{x}_T)\| \\ &= \left\| \left( \underbrace{\mathbf{A}_T^\top \mathbf{A}_T}_{\text{eigenvalues } \asymp n} \right)^{-1/2} \right\| \cdot \sqrt{k} \lesssim \sqrt{k/n} \end{aligned}$$

- When  $n/(k \log p)$  is sufficiently large,  $\max_{i \notin T} |\langle \mathbf{a}_i, \boldsymbol{\nu} \rangle| < 1$

## Details

---

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

$$\mathbb{P}(|\langle \mathbf{a}_i, \boldsymbol{\nu} \rangle| \geq 1 | \boldsymbol{\nu}) \leq 2 \exp\left(-\frac{1}{2\|\boldsymbol{\nu}\|_2^2}\right)$$

- With probability at least  $1 - e^{-cn}$ , we could also bound  $\|\boldsymbol{\nu}\|_2$  as

$$\|\boldsymbol{\nu}\|_2 \leq \sqrt{\frac{2k}{n}}$$

- We have

$$\begin{aligned} \mathbb{P}(\max_{i \in T^c} |\langle \mathbf{a}_i, \boldsymbol{\nu} \rangle| \geq 1) &\leq |T^c| \cdot \mathbb{P}(|\langle \mathbf{a}_i, \boldsymbol{\nu} \rangle| > 1) \quad \text{union bound} \\ &\leq p \int_{\boldsymbol{\nu}} \mathbb{P}(|\langle \mathbf{a}_i, \boldsymbol{\nu} \rangle| \geq 1 | \boldsymbol{\nu}) d\mu(\boldsymbol{\nu}). \end{aligned}$$

## Continued

---

Note that

$$\begin{aligned} & \int_{\boldsymbol{\nu}} \mathbb{P}(|\langle \mathbf{a}_i, \boldsymbol{\nu} \rangle| \geq 1 | \boldsymbol{\nu}) d\mu(\boldsymbol{\nu}) \\ &= \left( \int_{\|\boldsymbol{\nu}\|_2 \leq \sqrt{\frac{2k}{n}}} + \int_{\|\boldsymbol{\nu}\|_2 > \sqrt{\frac{2k}{n}}} \right) \mathbb{P}(|\langle \mathbf{a}_i, \boldsymbol{\nu} \rangle| \geq 1 | \boldsymbol{\nu}) d\mu(\boldsymbol{\nu}) \\ &\leq \int_{\|\boldsymbol{\nu}\|_2 \leq \sqrt{\frac{2k}{n}}} \mathbb{P}(|\langle \mathbf{a}_i, \boldsymbol{\nu} \rangle| \geq 1 | \boldsymbol{\nu}) d\mu(\boldsymbol{\nu}) + \mathbb{P} \left( \|\boldsymbol{\nu}\|_2 > \sqrt{\frac{2k}{n}} \right) \\ &\leq \int_{\|\boldsymbol{\nu}\|_2 \leq \sqrt{\frac{2k}{n}}} 2e^{-\frac{1}{2\|\boldsymbol{\nu}\|_2^2}} d\mu(\boldsymbol{\nu}) + e^{-cn} \\ &\leq 2e^{-\frac{n}{4k}} + e^{-cn} \leq 3e^{-\frac{n}{4k}}, \end{aligned}$$

which gives

$$\mathbb{P}(\max_{i \in T^c} |\langle \mathbf{a}_i, \boldsymbol{\nu} \rangle| \geq 1) \leq 3pe^{-\frac{n}{4k}} \leq p^{-\gamma}$$

by setting  $n = 4(\gamma + 1)k \log p$  for some constant  $\gamma > 0$ .

## More general random sampling

---

Consider a random design: each sampling vector  $\mathbf{a}_i$  is independently drawn from a distribution  $F$

$$\mathbf{a}_i \sim F$$

### Incoherence sampling:

- *Isotropy*:

$$\mathbb{E}[\mathbf{a}\mathbf{a}^\top] = \mathbf{I}, \quad \mathbf{a} \sim F$$

- components of  $\mathbf{a}$ : (i) unit variance; (ii) uncorrelated

- *Incoherence*: let  $\mu(F)$  be the smallest quantity s.t. for  $\mathbf{a} \sim F$ ,

$$\|\mathbf{a}\|_\infty^2 \leq \mu(F) \quad \text{with high probability}$$

# Incoherence

We want  $\mu(F)$  (resp.  $A$ ) to be small (resp. dense)!

What happen if sampling vectors  $a_i$  are sparse?

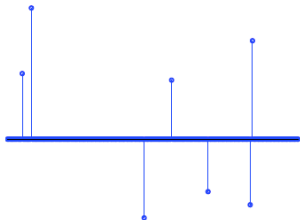
- Example:  $a_i \sim \text{Uniform}(\{\sqrt{p}e_1, \dots, \sqrt{p}e_p\})$

$$\underbrace{\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}}_y \text{ no information} = \sqrt{p} \underbrace{\begin{bmatrix} 1 & & & & & & & \\ & & & 1 & & & & \\ & & & & & & 1 & \\ & & & & & & & 1 \end{bmatrix}}_A \underbrace{\begin{bmatrix} 0 \\ 0 \\ 0 \\ 3 \\ 0 \\ 0 \\ 0 \\ 5 \\ 0 \\ 0 \end{bmatrix}}_x$$

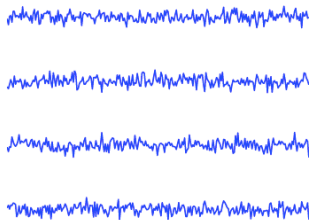
# Incoherent random sampling

---

concentrated vector



incoherent measurements



# A general RIPless theory

---

## Theorem 6.10 (Candes & Plan, '11)

Suppose  $x \in \mathbb{R}^p$  is  $k$ -sparse, and  $a_i \stackrel{\text{ind.}}{\sim} F$  is isotropic. Then  $\ell_1$  minimization is exact and unique with high prob., provided that

$$n \gtrsim \mu(F)k \log p$$

- Near-optimal even for **highly structured** sampling matrices
- Proof idea: produce an (approximate) dual certificate by a clever *golfing scheme* pioneered by David Gross



# Examples of incoherent sampling

---

- **Binary sensing:**  $\mathbb{P}(a[i] = \pm 1) = 0.5$ :

$$\mathbb{E}[\mathbf{a}\mathbf{a}^\top] = \mathbf{I}, \quad \|\mathbf{a}\|_\infty^2 = 1, \quad \mu = 1$$

$$\implies \ell_1\text{-min succeeds if } n \gtrsim k \log p$$

- **Gaussian sensing:**  $\mathbf{a} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ :

$$\mathbb{E}[\mathbf{a}\mathbf{a}^\top] = \mathbf{I}, \quad \|\mathbf{a}\|_\infty^2 \lesssim 2 \log p \quad \Rightarrow \quad \mu \asymp \log p$$

- **Partial Fourier transform:** pick a random frequency  $f \sim \text{Unif}\{0, \frac{1}{p}, \dots, \frac{p-1}{p}\}$  or  $f \sim \text{Unif}[0, 1]$  and set  $a[i] = e^{j2\pi fi}$ .

$$\mathbb{E}[\mathbf{a}\mathbf{a}^\top] = \mathbf{I}, \quad \|\mathbf{a}\|_\infty^2 = 1, \quad \mu = 1$$

$$\implies \ell_1\text{-min succeeds if } n \gtrsim k \log p$$

- Improves upon RIP-based result ( $n \gtrsim k \log^4 p$ )

# Examples of incoherent sampling

---

- Random convolution matrices: rows of  $\mathbf{A}$  are independently sampled from rows of

$$\mathbf{G} = \begin{bmatrix} g_0 & g_1 & g_2 & \cdots & g_{p-1} \\ g_{p-1} & g_0 & g_1 & \cdots & g_{p-2} \\ g_{p-2} & g_{p-1} & g_0 & \cdots & g_{p-3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ g_1 & g_2 & g_3 & \cdots & g_0 \end{bmatrix}$$

with  $\mathbb{P}(g_i = \pm 1) = 0.5$ . One has

$$\begin{aligned} \mathbb{E}[\mathbf{a}\mathbf{a}^\top] &= \mathbf{I}, & \|\mathbf{a}\|_\infty^2 &= 1, & \mu &= 1 \\ \implies & \ell_1\text{-min succeeds if } n \gtrsim k \log p \end{aligned}$$

- Improves upon RIP-based result ( $n \gtrsim k \log^4 p$ )

# A general scheme for dual construction

---

$$\begin{aligned} \text{Find} \quad & \boldsymbol{\nu} \in \mathbb{R}^n \\ \text{s.t.} \quad & \mathbf{A}_T^\top \boldsymbol{\nu} = \text{sign}(\mathbf{x}_T) \end{aligned} \tag{6.15}$$

$$\|\mathbf{A}_{T^c}^\top \boldsymbol{\nu}\|_\infty < 1 \tag{6.16}$$

A candidate: least squares solution w.r.t. (6.19)

$$\boldsymbol{\nu} = \mathbf{A}_T (\mathbf{A}_T^\top \mathbf{A}_T)^{-1} \text{sign}(\mathbf{x}_T) \quad (\text{explicit expression})$$

To verify (6.20), we need to control  $\mathbf{A}_{T^c}^\top \mathbf{A}_T (\mathbf{A}_T^\top \mathbf{A}_T)^{-1} \text{sign}(\mathbf{x}_T)$

- Issue 1: in general,  $\mathbf{A}_{T^c}$  and  $\mathbf{A}_T$  are dependent
- Issue 2:  $(\mathbf{A}_T^\top \mathbf{A}_T)^{-1}$  is hard to deal with

# A general scheme for dual construction

$$\begin{aligned} \text{Find } & \boldsymbol{\nu} \in \mathbb{R}^n \\ \text{s.t. } & \mathbf{A}_T^\top \boldsymbol{\nu} = \text{sign}(\mathbf{x}_T) \end{aligned} \quad (6.17)$$

$$\|\mathbf{A}_{T^c}^\top \boldsymbol{\nu}\|_\infty < 1 \quad (6.18)$$

**Key idea 1:** use iterative scheme to solve  
minimize $_{\boldsymbol{\nu}}$   $\frac{1}{2} \|\mathbf{A}_T^\top \boldsymbol{\nu} - \text{sign}(\mathbf{x}_T)\|^2$

for  $t = 1, 2, \dots$

$$\boldsymbol{\nu}^{(t)} = \boldsymbol{\nu}^{(t-1)} - \underbrace{\mathbf{A}_T \left( \mathbf{A}_T^\top \boldsymbol{\nu}^{(t-1)} - \text{sign}(\mathbf{x}_T) \right)}_{\text{grad of } \frac{1}{2} \|\mathbf{A}_T^\top \boldsymbol{\nu} - \text{sign}(\mathbf{x}_T)\|^2}$$

- Converges to a solution obeying the equality constraint; no inversion involved
- Issue: complicated dependency across iterations

# Golfing scheme (Gross '11)

---

**Key idea 2: sample splitting** — use independent samples for each iteration to decouple statistical dependency

- Partition  $\mathbf{A}$  into  $L$  row blocks  $\underbrace{\mathbf{A}^{(1)} \in \mathbb{R}^{n_1 \times p}, \dots, \mathbf{A}^{(L)} \in \mathbb{R}^{n_L \times p}}_{\text{independent}}$
- for  $t = 1, 2, \dots$  (stochastic gradient)

$$\boldsymbol{\nu}^{(t)} = \boldsymbol{\nu}^{(t-1)} - \underbrace{\mu_t \mathbf{A}_T^{(t)} \left( \mathbf{A}_T^{(t)\top} \boldsymbol{\nu}^{(t-1)} - \text{sign}(\mathbf{x}_T) \right)}_{\in \mathbb{R}^{n_t} \text{ (but we need } \boldsymbol{\nu} \in \mathbb{R}^n)}$$

# Golfing scheme (Gross '11)

**Key idea 2: sample splitting** — use independent samples for each iteration to decouple statistical dependency

- Partition  $\mathbf{A}$  into  $L$  row blocks  $\underbrace{\mathbf{A}^{(1)} \in \mathbb{R}^{n_1 \times p}, \dots, \mathbf{A}^{(L)} \in \mathbb{R}^{n_L \times p}}_{\text{independent}}$
- for  $t = 1, 2, \dots$  (stochastic gradient)

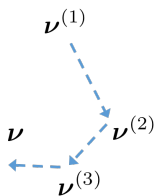
$$\boldsymbol{\nu}^{(t)} = \boldsymbol{\nu}^{(t-1)} - \mu_t \tilde{\mathbf{A}}_T^{(t)} \left( \tilde{\mathbf{A}}_T^{(t)\top} \boldsymbol{\nu}^{(t-1)} - \text{sign}(\mathbf{x}_T) \right)$$

where  $\tilde{\mathbf{A}}^{(t)} = \begin{bmatrix} \mathbf{0} \\ \mathbf{A}^{(t)} \\ \mathbf{0} \end{bmatrix} \in \mathbb{R}^{n \times p}$  is obtained by zero-padding

# Golfing scheme (Gross '11)

---

$$\boldsymbol{\nu}^{(t)} = \boldsymbol{\nu}^{(t-1)} - \mu_t \tilde{\mathbf{A}}_T^{(t)} \left( \underbrace{\tilde{\mathbf{A}}_T^{(t)\top} \boldsymbol{\nu}^{(t-1)} - \text{sign}(\mathbf{x}_T)}_{\text{depends only on } \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(t-1)}} \right)$$



- Statistical independence across iterations
  - By construction,  
 $\mathbf{A}_T^\top \boldsymbol{\nu}^{(t-1)} \in \text{range}(\mathbf{A}_T^{(i)\top}) \cap \dots \cap \text{range}(\mathbf{A}_T^{(L)\top})$
- Each iteration brings us closer to the target (like each golf shot brings us closer to the hole)

# A general scheme for dual construction

$$\begin{aligned} \text{Find} \quad & \boldsymbol{\nu} \in \mathbb{R}^n \\ \text{s.t.} \quad & \mathbf{A}_T^\top \boldsymbol{\nu} = \text{sign}(\mathbf{x}_T) \end{aligned} \tag{6.19}$$

$$\|\mathbf{A}_{T^c}^\top \boldsymbol{\nu}\|_\infty < 1 \tag{6.20}$$

The golfing scheme doesn't yield an exact dual certificate, but an inexact one.

## Theorem 6.11 (Inexact duality)

Let  $T := \text{supp}(\mathbf{x})$ . *Suppose*  $\|(\mathbf{A}_T^\top \mathbf{A}_T)^{-1}\| \leq 2$  *and*  
 $\max_{i \in T^c} \|\mathbf{A}_T^\top \mathbf{a}_i\| \leq 1$ . *If*

$$\exists \mathbf{u} = \mathbf{A}^\top \boldsymbol{\nu} \text{ for some } \boldsymbol{\nu} \in \mathbb{R}^n \quad \text{s.t.} \quad \begin{cases} \|\mathbf{u}_T - \text{sign}(\mathbf{x}_T)\| \leq 1/4, \\ \|\mathbf{u}_{T^c}\|_\infty \leq 1/2, \end{cases} ,$$

*then  $\mathbf{x}$  is unique solution to  $\ell_1$  minimization.*

Proof is similar to Theorem 6.9. The conditions in red is guaranteed with high probability via **concentration inequalities**.



# Reference

---

- [1] “*Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information*,” E. Candes, J. Romberg, and T. Tao, *IEEE Transactions on Information Theory*, 2006.
- [2] “*Compressed sensing*,” D. Donoho, *IEEE Transactions on Information Theory*, 2006
- [3] “*Near-optimal signal recovery from random projections: Universal encoding strategies?*,” E. Candes, and T. Tao, *IEEE Transactions on Information Theory*, 2006.
- [4] “*Lecture notes, Advanced topics in signal processing (ECE 8201)*,” Y. Chi, 2015.
- [5] “*A mathematical introduction to compressive sensing*,” S. Foucart and H. Rauhut, *Springer*, 2013.
- [6] “*On sparse reconstruction from Fourier and Gaussian measurements*,” M. Rudelson and R. Vershynin, *Communications on Pure and Applied Mathematics*, 2008.

# Reference

---

- [7] “*Decoding by linear programming*,” E. Candes, and T. Tao, *IEEE Transactions on Information Theory*, 2006.
- [8] “*The restricted isometry property and its implications for compressed sensing*,” E. Candes, *Compte Rendus de l’Academie des Sciences*, 2008.
- [9] “*Introduction to the non-asymptotic analysis of random matrices*,” R. Vershynin, *Compressed Sensing: Theory and Applications*, 2010.
- [10] “*Iterative hard thresholding for compressed sensing*,” T. Blumensath, M. Davies, *Applied and computational harmonic analysis*, 2009.
- [11] “*Recovering low-rank matrices from few coefficients in any basis*,” D. Gross, *IEEE Transactions on Information Theory*, 2011.
- [12] “*A probabilistic and RIPless theory of compressed sensing*,” E. Candes and Y. Plan, *IEEE Transactions on Information Theory*, 2011.

# Reference

---

- [13] "*Statistical learning with sparsity: the Lasso and generalizations*,"  
T. Hastie, R. Tibshirani, and M. Wainwright, 2015.
- [14] "*Suprema of chaos processes and the restricted isometry property*,"  
F. Krahmer, S. Mendelson, and H. Rauhut, *Communications on Pure and Applied Mathematics*, 2014.