

# **A primer on high-dimensional statistics: Part I**

Martin Wainwright

UC Berkeley  
Departments of Statistics, and EECS

Spring School, Westerland, March 2015

# Introduction

- classical asymptotic theory: sample size  $n \rightarrow +\infty$  with number of parameters  $d$  fixed
  - ▶ law of large numbers, central limit theory
  - ▶ consistency of maximum likelihood estimation

# Introduction

- classical asymptotic theory: sample size  $n \rightarrow +\infty$  with number of parameters  $d$  fixed
  - ▶ law of large numbers, central limit theory
  - ▶ consistency of maximum likelihood estimation
  
- modern applications in science and engineering:
  - ▶ large-scale problems: both  $d$  and  $n$  may be large (possibly  $d \gg n$ )
  - ▶ need for **high-dimensional theory** that provides non-asymptotic results for  $(n, d)$

# Introduction

- classical asymptotic theory: sample size  $n \rightarrow +\infty$  with number of parameters  $d$  fixed
  - ▶ law of large numbers, central limit theory
  - ▶ consistency of maximum likelihood estimation
  
- modern applications in science and engineering:
  - ▶ large-scale problems: both  $d$  and  $n$  may be large (possibly  $d \gg n$ )
  - ▶ need for **high-dimensional theory** that provides non-asymptotic results for  $(n, d)$
  
- **curses** and **blessings** of high dimensionality
  - ▶ **exponential explosions in computational complexity**
  - ▶ **statistical curses (sample complexity)**
  - ▶ **concentration of measure**

# Introduction

- modern applications in science and engineering:
  - ▶ large-scale problems: both  $d$  and  $n$  may be large (possibly  $d \gg n$ )
  - ▶ need for **high-dimensional theory** that provides non-asymptotic results for  $(n, d)$
- **curse** and **blessings** of high dimensionality
  - ▶ **exponential explosions in computational complexity**
  - ▶ **statistical curses (sample complexity)**
  - ▶ **concentration of measure**

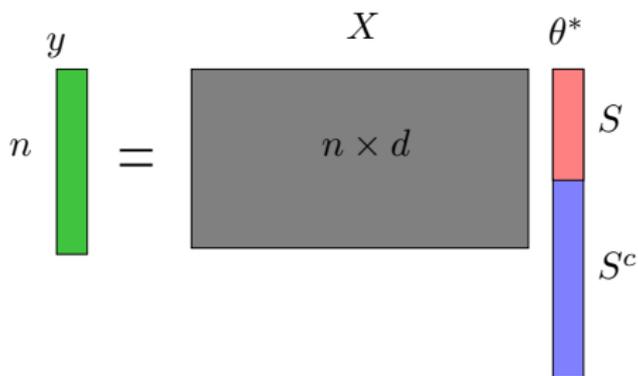
## Key questions:

- What **embedded low-dimensional structures** are present in data?
- How can they can be **exploited algorithmically**?

# Outline

- ① Lecture 1—2: Basics of sparse linear models
  - ▶ Sparse linear systems:  $\ell_0/\ell_1$  equivalence
  - ▶ Noisy case: Lasso,  $\ell_2$ -bounds, prediction error and variable selection
- ② Lectures 2—3: More general theory

# Noiseless linear models and basis pursuit



- under-determined linear system: unidentifiable without constraints
- say  $\theta^* \in \mathbb{R}^d$  is sparse: supported on  $S \subset \{1, 2, \dots, d\}$ .

$\ell_0$ -optimization

$$\theta^* = \arg \min_{\theta \in \mathbb{R}^d} \|\theta\|_0$$
$$X\theta = y$$

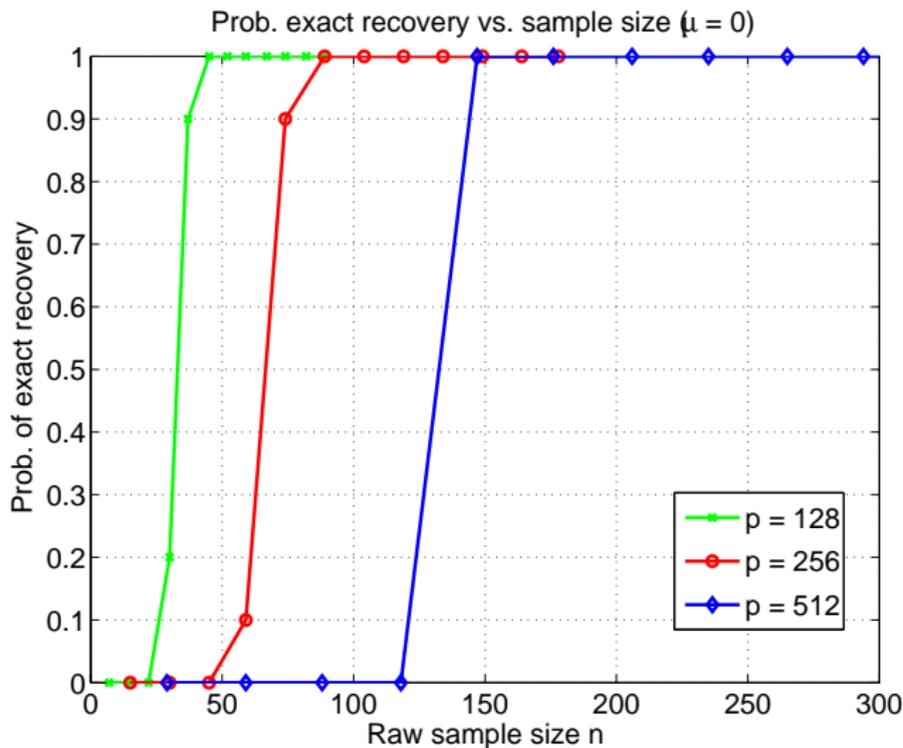
Computationally intractable  
NP-hard

$\ell_1$ -relaxation

$$\hat{\theta} \in \arg \min_{\theta \in \mathbb{R}^d} \|\theta\|_1$$
$$X\theta = y$$

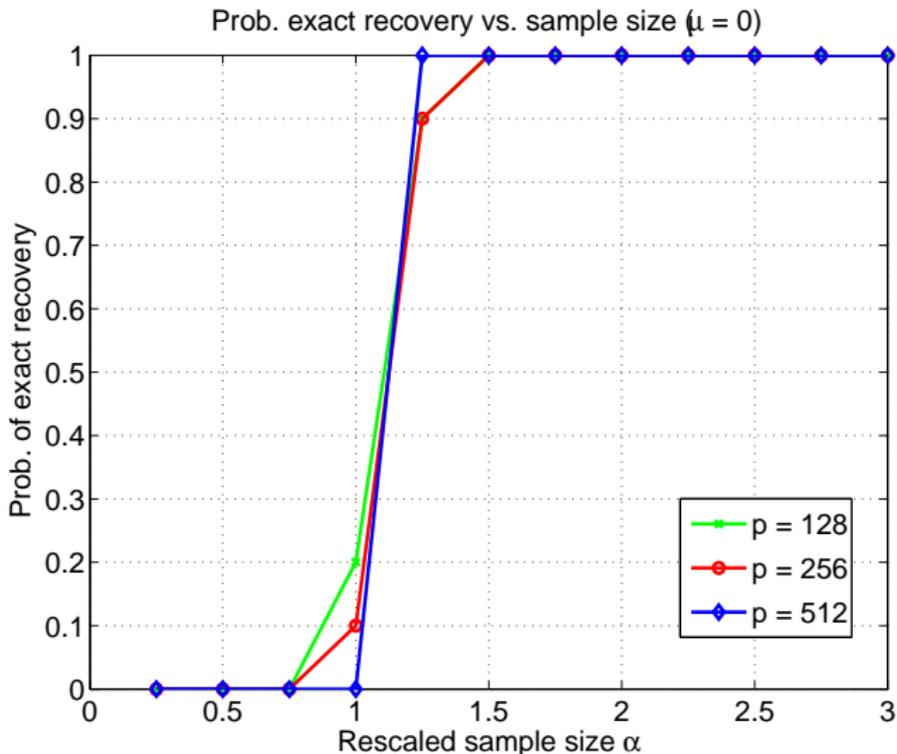
Linear program (easy to solve)  
Basis pursuit relaxation

# Noiseless $\ell_1$ recovery: Unrescaled sample size



Probability of recovery versus sample size  $n$ .

# Noiseless $\ell_1$ recovery: Rescaled



Probabability of recovery versus **rescaled sample size**  $\alpha := \frac{n}{s \log(d/s)}$ .

# Restricted nullspace: necessary and sufficient

## Definition

For a fixed  $S \subset \{1, 2, \dots, d\}$ , the matrix  $X \in \mathbb{R}^{n \times d}$  satisfies the restricted nullspace property w.r.t.  $S$ , or  $\text{RN}(S)$  for short, if

$$\underbrace{\{\Delta \in \mathbb{R}^d \mid X\Delta = 0\}}_{\text{N}(X)} \cap \underbrace{\{\Delta \in \mathbb{R}^d \mid \|\Delta_{S^c}\|_1 \leq \|\Delta_S\|_1\}}_{\text{C}(S)} = \{0\}.$$

(Donoho & Xu, 2001; Feuer & Nemirovski, 2003; Cohen et al, 2009)

# Restricted nullspace: necessary and sufficient

## Definition

For a fixed  $S \subset \{1, 2, \dots, d\}$ , the matrix  $X \in \mathbb{R}^{n \times d}$  satisfies the restricted nullspace property w.r.t.  $S$ , or  $\text{RN}(S)$  for short, if

$$\underbrace{\{\Delta \in \mathbb{R}^d \mid X\Delta = 0\}}_{\text{N}(X)} \cap \underbrace{\{\Delta \in \mathbb{R}^d \mid \|\Delta_{S^c}\|_1 \leq \|\Delta_S\|_1\}}_{\text{C}(S)} = \{0\}.$$

(Donoho & Xu, 2001; Feuer & Nemirovski, 2003; Cohen et al, 2009)

## Proposition

Basis pursuit  $\ell_1$ -relaxation is exact for all  $S$ -sparse vectors  $\iff X$  satisfies  $\text{RN}(S)$ .

# Restricted nullspace: necessary and sufficient

## Definition

For a fixed  $S \subset \{1, 2, \dots, d\}$ , the matrix  $X \in \mathbb{R}^{n \times d}$  satisfies the restricted nullspace property w.r.t.  $S$ , or  $\text{RN}(S)$  for short, if

$$\underbrace{\{\Delta \in \mathbb{R}^d \mid X\Delta = 0\}}_{\text{N}(X)} \cap \underbrace{\{\Delta \in \mathbb{R}^d \mid \|\Delta_{S^c}\|_1 \leq \|\Delta_S\|_1\}}_{\text{C}(S)} = \{0\}.$$

(Donoho & Xu, 2001; Feuer & Nemirovski, 2003; Cohen et al, 2009)

## Proof (sufficiency):

(1) Error vector  $\hat{\Delta} = \theta^* - \hat{\theta}$  satisfies  $X\hat{\Delta} = 0$ , and hence  $\hat{\Delta} \in \text{N}(X)$ .

(2) Show that  $\hat{\Delta} \in \text{C}(S)$

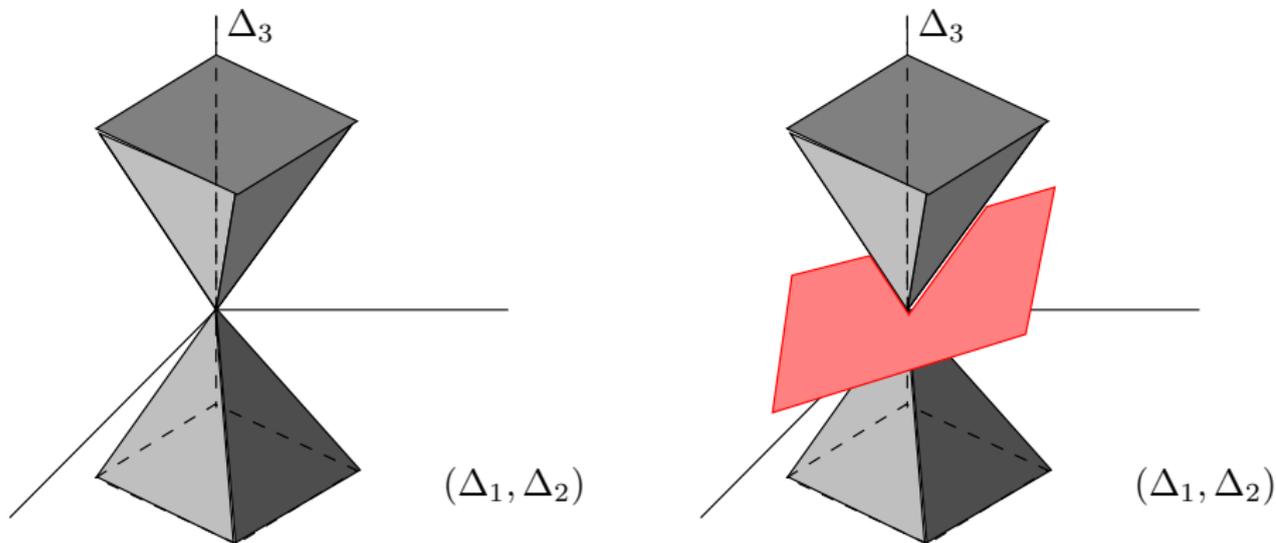
Optimality of  $\hat{\theta}$ :  $\|\hat{\theta}\|_1 \leq \|\theta^*\|_1 = \|\theta_S^*\|_1.$

Sparsity of  $\theta^*$ :  $\|\hat{\theta}\|_1 = \|\theta^* + \hat{\Delta}\|_1 = \|\theta_S^* + \hat{\Delta}_S\|_1 + \|\hat{\Delta}_{S^c}\|_1.$

Triangle inequality:  $\|\theta_S^* + \hat{\Delta}_S\|_1 + \|\hat{\Delta}_{S^c}\|_1 \geq \|\theta_S^*\|_1 - \|\hat{\Delta}_S\|_1 + \|\hat{\Delta}_{S^c}\|_1.$

(3) Hence,  $\hat{\Delta} \in \text{N}(X) \cap \text{C}(S)$ , and  $(\text{RN}) \implies \hat{\Delta} = 0.$

# Illustration of restricted nullspace property



- consider  $\theta^* = (0, 0, \theta_3^*)$ , so that  $S = \{3\}$ .
- error vector  $\widehat{\Delta} = \widehat{\theta} - \theta^*$  belongs to the set

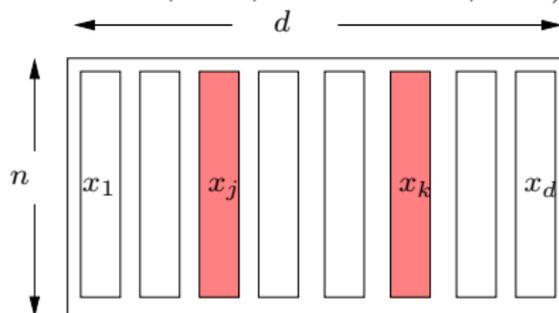
$$\mathbb{C}(S; 1) := \{(\Delta_1, \Delta_2, \Delta_3) \in \mathbb{R}^3 \mid |\Delta_1| + |\Delta_2| \leq |\Delta_3|\}.$$

# Some sufficient conditions

How to verify RN property for a given sparsity  $s$ ?

- ① **Elementwise incoherence condition** (Donoho & Xuo, 2001; Feuer & Nem., 2003)

$$\max_{j,k=1,\dots,d} \left| \left( \frac{X^T X}{n} - I_{d \times d} \right)_{jk} \right| \leq \frac{\delta_1}{s}$$

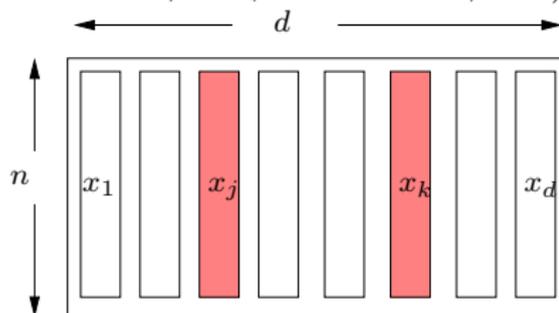


# Some sufficient conditions

How to verify RN property for a given sparsity  $s$ ?

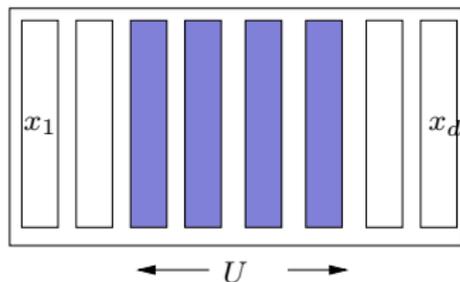
- ① **Elementwise incoherence condition** (Donoho & Xuo, 2001; Feuer & Nem., 2003)

$$\max_{j,k=1,\dots,d} \left| \left( \frac{X^T X}{n} - I_{d \times d} \right)_{jk} \right| \leq \frac{\delta_1}{s}$$



- ② **Restricted isometry**, or submatrix incoherence (Candes & Tao, 2005)

$$\max_{|U| \leq 2s} \left\| \left( \frac{X^T X}{n} - I_{d \times d} \right)_{UU} \right\|_{\text{op}} \leq \delta_{2s}.$$

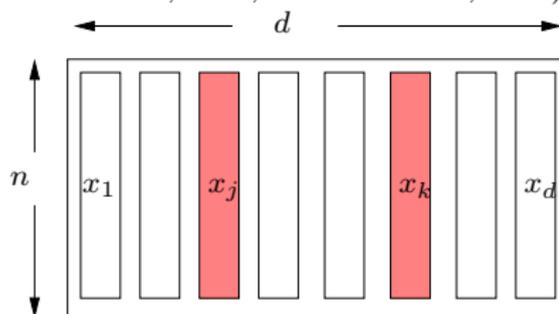


# Some sufficient conditions

How to verify RN property for a given sparsity  $s$ ?

- ① **Elementwise incoherence condition** (Donoho & Xuo, 2001; Feuer & Nem., 2003)

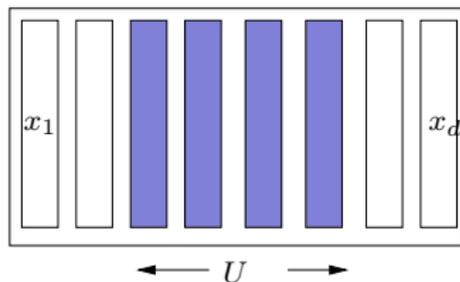
$$\max_{j,k=1,\dots,d} \left| \left( \frac{X^T X}{n} - I_{d \times d} \right)_{jk} \right| \leq \frac{\delta_1}{s}$$



Matrices with i.i.d. sub-Gaussian entries: holds w.h.p. for  $n = \Omega(s^2 \log d)$

- ② **Restricted isometry**, or submatrix incoherence (Candes & Tao, 2005)

$$\max_{|U| \leq 2s} \left\| \left( \frac{X^T X}{n} - I_{d \times d} \right)_{UU} \right\|_{\text{op}} \leq \delta_{2s}.$$

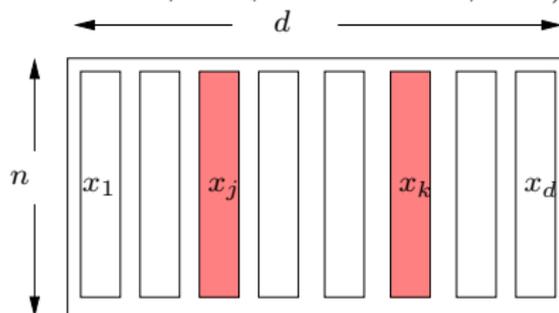


# Some sufficient conditions

How to verify RN property for a given sparsity  $s$ ?

- ① **Elementwise incoherence condition** (Donoho & Xuo, 2001; Feuer & Nem., 2003)

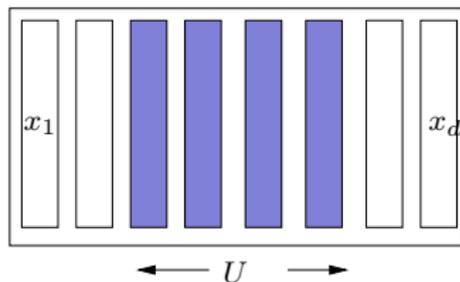
$$\max_{j,k=1,\dots,d} \left| \left( \frac{X^T X}{n} - I_{d \times d} \right)_{jk} \right| \leq \frac{\delta_1}{s}$$



Matrices with i.i.d. sub-Gaussian entries: holds w.h.p. for  $n = \Omega(s^2 \log d)$

- ② **Restricted isometry**, or submatrix incoherence (Candes & Tao, 2005)

$$\max_{|U| \leq 2s} \left\| \left( \frac{X^T X}{n} - I_{d \times d} \right)_{UU} \right\|_{\text{op}} \leq \delta_{2s}.$$



Matrices with i.i.d. sub-Gaussian entries: holds w.h.p. for  $n = \Omega(s \log \frac{d}{s})$

# Violating matrix incoherence (elementwise/RIP)

## Important:

Incoherence/RIP conditions imply RN, but are far from necessary.

Very easy to violate them.....

# Violating matrix incoherence (elementwise/RIP)

Form random design matrix

$$X = \underbrace{\begin{bmatrix} x_1 & x_2 & \dots & x_d \end{bmatrix}}_{d \text{ columns}} = \underbrace{\begin{bmatrix} X_1^T \\ X_2^T \\ \vdots \\ X_n^T \end{bmatrix}}_{n \text{ rows}} \in \mathbb{R}^{n \times d}, \quad \text{each row } X_i \sim N(0, \Sigma), \text{ i.i.d.}$$

**Example:** For some  $\mu \in (0, 1)$ , consider the covariance matrix

$$\Sigma = (1 - \mu)I_{d \times d} + \mu \mathbf{1}\mathbf{1}^T.$$

# Violating matrix incoherence (elementwise/RIP)

Form random design matrix

$$X = \underbrace{\begin{bmatrix} x_1 & x_2 & \dots & x_d \end{bmatrix}}_{d \text{ columns}} = \underbrace{\begin{bmatrix} X_1^T \\ X_2^T \\ \vdots \\ X_n^T \end{bmatrix}}_{n \text{ rows}} \in \mathbb{R}^{n \times d}, \quad \text{each row } X_i \sim N(0, \Sigma), \text{ i.i.d.}$$

**Example:** For some  $\mu \in (0, 1)$ , consider the covariance matrix

$$\Sigma = (1 - \mu)I_{d \times d} + \mu \mathbf{1}\mathbf{1}^T.$$

- **Elementwise incoherence violated:** for any  $j \neq k$

$$\mathbb{P} \left[ \frac{\langle x_j, x_k \rangle}{n} \geq \mu - \epsilon \right] \geq 1 - c_1 \exp(-c_2 n \epsilon^2).$$

# Violating matrix incoherence (elementwise/RIP)

Form random design matrix

$$X = \underbrace{\begin{bmatrix} x_1 & x_2 & \dots & x_d \end{bmatrix}}_{d \text{ columns}} = \underbrace{\begin{bmatrix} X_1^T \\ X_2^T \\ \vdots \\ X_n^T \end{bmatrix}}_{n \text{ rows}} \in \mathbb{R}^{n \times d}, \quad \text{each row } X_i \sim N(0, \Sigma), \text{ i.i.d.}$$

**Example:** For some  $\mu \in (0, 1)$ , consider the covariance matrix

$$\Sigma = (1 - \mu)I_{d \times d} + \mu \mathbf{1}\mathbf{1}^T.$$

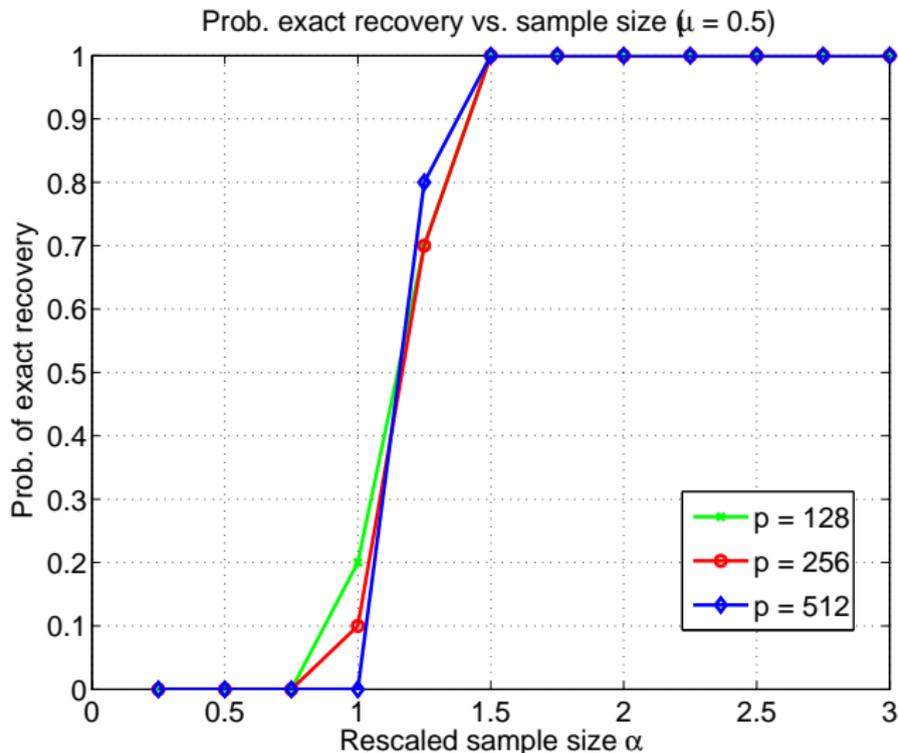
- **Elementwise incoherence violated:** for any  $j \neq k$

$$\mathbb{P} \left[ \frac{\langle x_j, x_k \rangle}{n} \geq \mu - \epsilon \right] \geq 1 - c_1 \exp(-c_2 n \epsilon^2).$$

- **RIP constants tend to infinity** as  $(n, |S|)$  increases:

$$\mathbb{P} \left[ \left\| \frac{X_S^T X_S}{n} - I_{s \times s} \right\|_2 \geq \mu(s-1) - 1 - \epsilon \right] \geq 1 - c_1 \exp(-c_2 n \epsilon^2).$$

# Noiseless $\ell_1$ recovery for $\mu = 0.5$



Probab. versus rescaled sample size  $\alpha := \frac{n}{s \log(d/s)}$ .

# Direct result for restricted nullspace/eigenvalues

**Theorem (Raskutti, W., & Yu, 2010; Rudelson & Zhou, 2012)**

*Random Gaussian/sub-Gaussian matrix  $X \in \mathbb{R}^{n \times d}$  with i.i.d. rows, covariance  $\Sigma$ , and let  $\kappa^2 = \max_j \Sigma_{jj}$  be the maximal variance. Then*

$$\frac{\|X\theta\|_2^2}{n} \geq c_1 \|\Sigma^{1/2}\theta\|_2^2 - c_2 \kappa^2(\Sigma) \frac{\log(e d (\frac{\|\theta\|_2}{\|\theta\|_1})^2)}{n} \|\theta\|_1^2 \quad \text{for all non-zero } \theta \in \mathbb{R}^d$$

*with probability at least  $1 - 2e^{-c_3 n}$ .*

# Direct result for restricted nullspace/eigenvalues

## Theorem (Raskutti, W., & Yu, 2010; Rudelson & Zhou, 2012)

Random Gaussian/sub-Gaussian matrix  $X \in \mathbb{R}^{n \times d}$  with i.i.d. rows, covariance  $\Sigma$ , and let  $\kappa^2 = \max_j \Sigma_{jj}$  be the maximal variance. Then

$$\frac{\|X\theta\|_2^2}{n} \geq c_1 \|\Sigma^{1/2}\theta\|_2^2 - c_2 \kappa^2(\Sigma) \frac{\log(e d (\frac{\|\theta\|_2}{\|\theta\|_1})^2)}{n} \|\theta\|_1^2 \quad \text{for all non-zero } \theta \in \mathbb{R}^d$$

with probability at least  $1 - 2e^{-c_3 n}$ .

- many interesting matrix families are covered
  - ▶ Toeplitz dependency
  - ▶ constant  $\mu$ -correlation (previous example)
  - ▶ covariance matrix  $\Sigma$  can even be degenerate
- related results hold for generalized linear models

## Easy verification of restricted nullspace

- for any  $\Delta \in \mathbb{C}(S)$ , we have

$$\|\Delta\|_1 = \|\Delta_S\|_1 + \|\Delta_{S^c}\|_1 \leq 2\|\Delta_S\| \leq 2\sqrt{s} \|\Delta\|_2$$

- applying previous result:

$$\frac{\|X\Delta\|_2^2}{n} \geq \underbrace{\left\{ c_1 \lambda_{\min}(\Sigma) - 4c_2 \kappa^2(\Sigma) \frac{s \log d}{n} \right\}}_{\gamma(\Sigma)} \|\Delta\|_2^2.$$

## Easy verification of restricted nullspace

- for any  $\Delta \in \mathbb{C}(S)$ , we have

$$\|\Delta\|_1 = \|\Delta_S\|_1 + \|\Delta_{S^c}\|_1 \leq 2\|\Delta_S\| \leq 2\sqrt{s} \|\Delta\|_2$$

- applying previous result:

$$\frac{\|X\Delta\|_2^2}{n} \geq \underbrace{\left\{ c_1 \lambda_{\min}(\Sigma) - 4c_2 \kappa^2(\Sigma) \frac{s \log d}{n} \right\}}_{\gamma(\Sigma)} \|\Delta\|_2^2.$$

- have actually proven much more than restricted nullspace....

## Easy verification of restricted nullspace

- for any  $\Delta \in \mathbb{C}(S)$ , we have

$$\|\Delta\|_1 = \|\Delta_S\|_1 + \|\Delta_{S^c}\|_1 \leq 2\|\Delta_S\| \leq 2\sqrt{s} \|\Delta\|_2$$

- applying previous result:

$$\frac{\|X\Delta\|_2^2}{n} \geq \underbrace{\left\{ c_1 \lambda_{\min}(\Sigma) - 4c_2 \kappa^2(\Sigma) \frac{s \log d}{n} \right\}}_{\gamma(\Sigma)} \|\Delta\|_2^2.$$

- have actually proven much more than restricted nullspace....

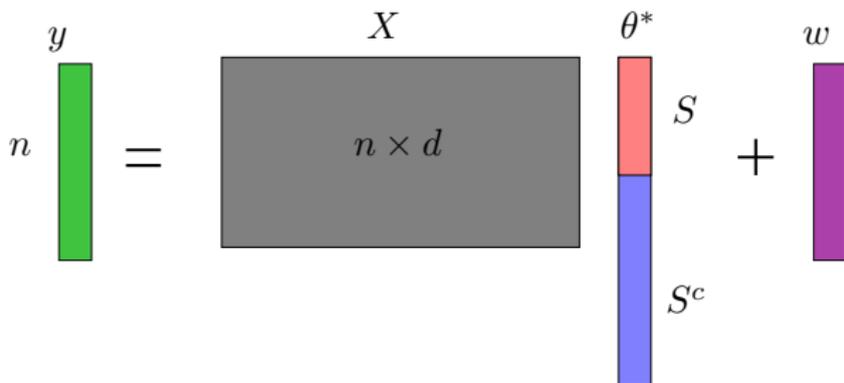
### Definition

A design matrix  $X \in \mathbb{R}^{n \times d}$  satisfies the *restricted eigenvalue* (RE) condition over  $S$  (denote  $\text{RE}(S)$ ) with parameters  $\alpha \geq 1$  and  $\gamma > 0$  if

$$\frac{\|X\Delta\|_2^2}{n} \geq \gamma \|\Delta\|_2^2 \quad \text{for all } \Delta \in \mathbb{R}^d \text{ such that } \|\Delta_{S^c}\|_1 \leq \alpha \|\Delta_S\|_1.$$

# Lasso and restricted eigenvalues

Turning to noisy observations...



**Estimator:** Lasso program

$$\hat{\theta}_{\lambda_n} \in \arg \min_{\theta \in \mathbb{R}^d} \left\{ \frac{1}{2n} \|y - X\theta\|_2^2 + \lambda_n \|\theta\|_1 \right\}.$$

**Goal:** Obtain bounds on { prediction error, parametric error, variable selection }.

# Different error metrics

1 (In-sample) prediction error:  $\|X(\hat{\theta} - \theta^*)\|_2^2/n$

- ▶ “weakest” error measure
- ▶ appropriate when  $\theta^*$  itself not of primary interest
- ▶ strong dependence between columns of  $X$  possible (for slow rate)
- ▶ proof technique: basic inequality

# Different error metrics

- 1 (In-sample) prediction error:  $\|X(\hat{\theta} - \theta^*)\|_2^2/n$ 
  - ▶ “weakest” error measure
  - ▶ appropriate when  $\theta^*$  itself not of primary interest
  - ▶ strong dependence between columns of  $X$  possible (for slow rate)
  - ▶ proof technique: basic inequality
  
- 2 parametric error:  $\|\hat{\theta} - \theta^*\|_r$  for some  $r \in [1, \infty]$ 
  - ▶ appropriate for recovery problems
  - ▶ RE-type conditions appear in both lower/upper bounds
  - ▶ variable selection is not guaranteed
  - ▶ proof technique: basic inequality

# Different error metrics

- 1 (In-sample) prediction error:  $\|X(\hat{\theta} - \theta^*)\|_2^2/n$ 
  - ▶ “weakest” error measure
  - ▶ appropriate when  $\theta^*$  itself not of primary interest
  - ▶ strong dependence between columns of  $X$  possible (for slow rate)
  - ▶ proof technique: basic inequality
- 2 parametric error:  $\|\hat{\theta} - \theta^*\|_r$  for some  $r \in [1, \infty]$ 
  - ▶ appropriate for recovery problems
  - ▶ RE-type conditions appear in both lower/upper bounds
  - ▶ variable selection is not guaranteed
  - ▶ proof technique: basic inequality
- 3 variable selection: is  $\text{supp}(\hat{\theta})$  equal to  $\text{supp}(\theta^*)$ ?
  - ▶ appropriate when non-zero locations are of scientific interest
  - ▶ most stringent of all three criteria
  - ▶ requires incoherence or irrepresentability conditions on  $X$
  - ▶ proof technique: primal-dual witness condition

## Lasso $\ell_2$ -bounds: Four simple steps

Let's analyze constrained version:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2n} \|y - X\theta\|_2^2 \quad \text{such that } \|\theta\|_1 \leq R = \|\theta^*\|_1.$$

---

## Lasso $\ell_2$ -bounds: Four simple steps

Let's analyze constrained version:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2n} \|y - X\theta\|_2^2 \quad \text{such that } \|\theta\|_1 \leq R = \|\theta^*\|_1.$$

---

(1) By **optimality of  $\hat{\theta}$**  and **feasibility of  $\theta^*$** :

$$\frac{1}{2n} \|y - X\hat{\theta}\|_2^2 \leq \frac{1}{2n} \|y - X\theta^*\|_2^2.$$

# Lasso $\ell_2$ -bounds: Four simple steps

Let's analyze constrained version:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2n} \|y - X\theta\|_2^2 \quad \text{such that } \|\theta\|_1 \leq R = \|\theta^*\|_1.$$

---

(1) By **optimality of  $\hat{\theta}$**  and **feasibility of  $\theta^*$** :

$$\frac{1}{2n} \|y - X\hat{\theta}\|_2^2 \leq \frac{1}{2n} \|y - X\theta^*\|_2^2.$$

(2) Derive a **basic inequality**: re-arranging in terms of  $\hat{\Delta} = \hat{\theta} - \theta^*$ :

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle.$$

# Lasso $\ell_2$ -bounds: Four simple steps

Let's analyze constrained version:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2n} \|y - X\theta\|_2^2 \quad \text{such that } \|\theta\|_1 \leq R = \|\theta^*\|_1.$$

---

(1) By **optimality of  $\hat{\theta}$**  and **feasibility of  $\theta^*$** :

$$\frac{1}{2n} \|y - X\hat{\theta}\|_2^2 \leq \frac{1}{2n} \|y - X\theta^*\|_2^2.$$

(2) Derive a **basic inequality**: re-arranging in terms of  $\hat{\Delta} = \hat{\theta} - \theta^*$ :

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle.$$

(3) **Restricted eigenvalue for LHS**; **Hölder's inequality for RHS**

$$\gamma \|\hat{\Delta}\|_2^2 \leq \frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle \leq 2 \|\hat{\Delta}\|_1 \left\| \frac{X^T w}{n} \right\|_\infty.$$

# Lasso $\ell_2$ -bounds: Four simple steps

Let's analyze constrained version:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2n} \|y - X\theta\|_2^2 \quad \text{such that } \|\theta\|_1 \leq R = \|\theta^*\|_1.$$

---

(1) By **optimality of  $\hat{\theta}$**  and **feasibility of  $\theta^*$** :

$$\frac{1}{2n} \|y - X\hat{\theta}\|_2^2 \leq \frac{1}{2n} \|y - X\theta^*\|_2^2.$$

(2) Derive a **basic inequality**: re-arranging in terms of  $\hat{\Delta} = \hat{\theta} - \theta^*$ :

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle.$$

(3) **Restricted eigenvalue for LHS**; **Hölder's inequality for RHS**

$$\gamma \|\hat{\Delta}\|_2^2 \leq \frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle \leq 2 \|\hat{\Delta}\|_1 \left\| \frac{X^T w}{n} \right\|_\infty.$$

(4) As before,  $\hat{\Delta} \in \mathbb{C}(S)$ , so that  $\|\hat{\Delta}\|_1 \leq 2\sqrt{s}\|\hat{\Delta}\|_2$ , and hence

$$\|\hat{\Delta}\|_2 \leq \frac{4}{\gamma} \sqrt{s} \left\| \frac{X^T w}{n} \right\|_\infty.$$

# Lasso error bounds for different models

## Proposition

Suppose that

- vector  $\theta^*$  has support  $S$ , with cardinality  $s$ , and
- design matrix  $X$  satisfies RE( $S$ ) with parameter  $\gamma > 0$ .

For constrained Lasso with  $R = \|\theta^*\|_1$  or regularized Lasso with  $\lambda_n = 2\|X^T w/n\|_\infty$ , any optimal solution  $\hat{\theta}$  satisfies the bound

$$\|\hat{\theta} - \theta^*\|_2 \leq \frac{4\sqrt{s}}{\gamma} \left\| \frac{X^T w}{n} \right\|_\infty.$$

# Lasso error bounds for different models

## Proposition

Suppose that

- vector  $\theta^*$  has support  $S$ , with cardinality  $s$ , and
- design matrix  $X$  satisfies RE( $S$ ) with parameter  $\gamma > 0$ .

For constrained Lasso with  $R = \|\theta^*\|_1$  or regularized Lasso with  $\lambda_n = 2\|X^T w/n\|_\infty$ , any optimal solution  $\hat{\theta}$  satisfies the bound

$$\|\hat{\theta} - \theta^*\|_2 \leq \frac{4\sqrt{s}}{\gamma} \left\| \frac{X^T w}{n} \right\|_\infty.$$

- this is a deterministic result on the set of optimizers
- various corollaries for specific statistical models

# Lasso error bounds for different models

## Proposition

Suppose that

- vector  $\theta^*$  has support  $S$ , with cardinality  $s$ , and
- design matrix  $X$  satisfies RE( $S$ ) with parameter  $\gamma > 0$ .

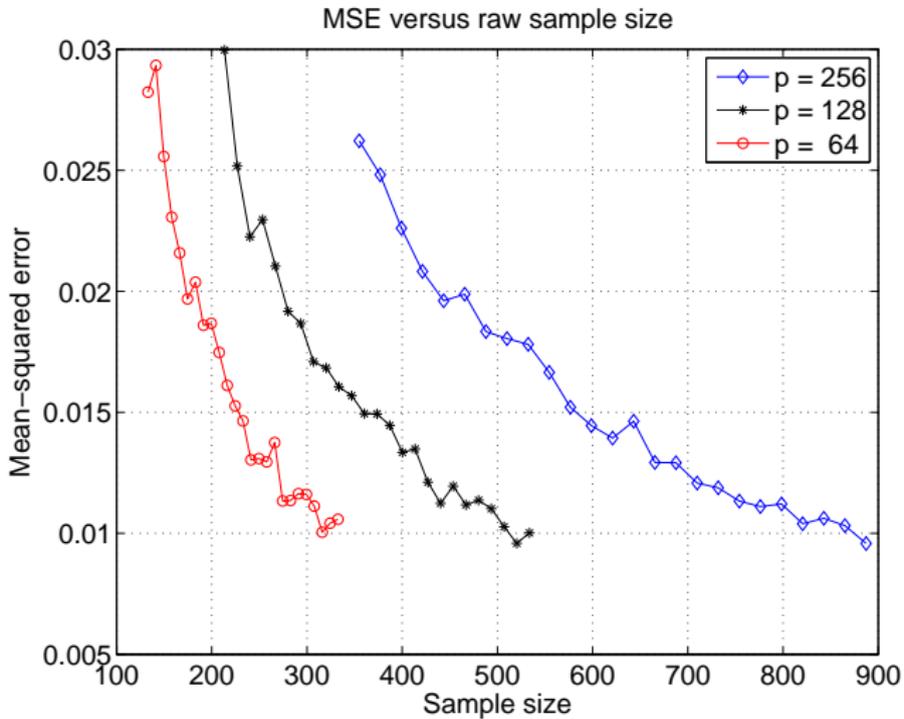
For constrained Lasso with  $R = \|\theta^*\|_1$  or regularized Lasso with  $\lambda_n = 2\|X^T w/n\|_\infty$ , any optimal solution  $\hat{\theta}$  satisfies the bound

$$\|\hat{\theta} - \theta^*\|_2 \leq \frac{4\sqrt{s}}{\gamma} \left\| \frac{X^T w}{n} \right\|_\infty.$$

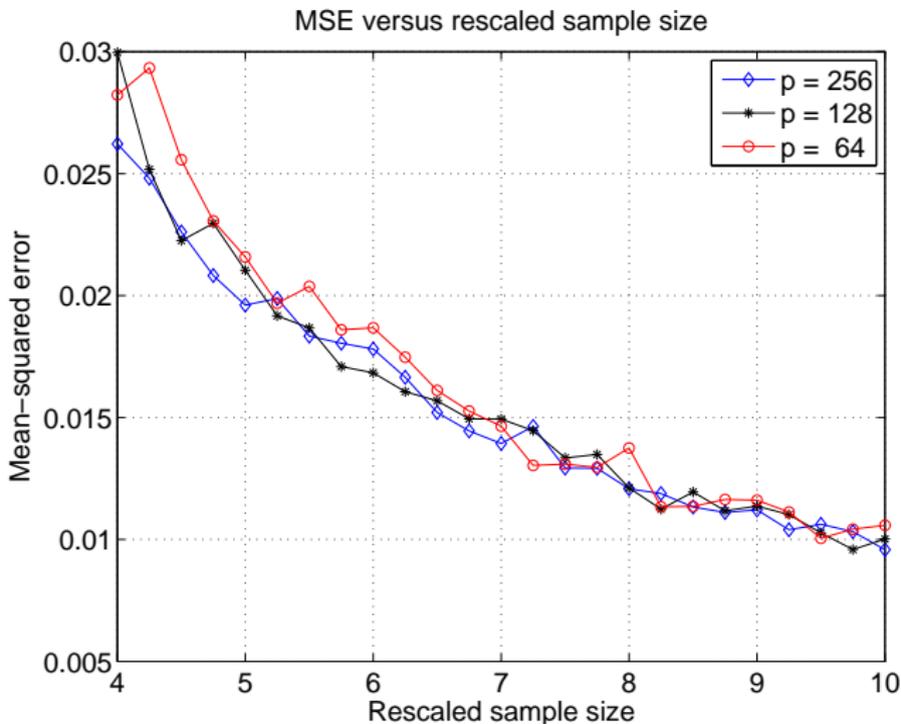
- this is a deterministic result on the set of optimizers
- various corollaries for specific statistical models
  - ▶ Compressed sensing:  $X_{ij} \sim N(0, 1)$  and bounded noise  $\|w\|_2 \leq \sigma\sqrt{n}$
  - ▶ Deterministic design:  $X$  with bounded columns and  $w_i \sim N(0, \sigma^2)$

$$\left\| \frac{X^T w}{n} \right\|_\infty \leq \sqrt{\frac{3\sigma^2 \log d}{n}} \quad \text{w.h.p.} \implies \|\hat{\theta} - \theta^*\|_2 \leq \frac{4\sigma}{\gamma} \sqrt{3 \frac{s \log d}{n}}.$$

# Lasso $\ell_2$ -error: Unrescaled sample size



# Lasso $l_2$ -error: Rescaled sample size



Rescaled sample size  $\frac{n}{s \log p/s}$ .

# Extension to an oracle inequality

Previous theory assumed that  $\theta^*$  was “hard” sparse. Not realistic in practice.

# Extension to an oracle inequality

Previous theory assumed that  $\theta^*$  was “hard” sparse. Not realistic in practice.

## Theorem (An oracle inequality)

Suppose that least-squares loss satisfies  $\gamma$ -RE condition. Then for  $\lambda_n \geq \max\{2\|\frac{X^T w}{n}\|_\infty, \sqrt{\frac{\log d}{n}}\}$ , any optimal Lasso solution satisfies

$$\|\hat{\theta} - \theta^*\|_2^2 \leq \min_{S \subseteq \{1, \dots, d\}} \left\{ \underbrace{\frac{9}{4} \frac{\lambda_n^2}{\gamma^2} |S|}_{\text{estimation error}} + \underbrace{\frac{2\lambda_n}{\gamma} \|\theta_{S^c}^*\|_1}_{\text{approximation error}} \right\}.$$

(cf. Bunea et al., 2007; Buhlmann and van de Geer, 2009; Koltchinski et al., 2011)

# Extension to an oracle inequality

Previous theory assumed that  $\theta^*$  was “hard” sparse. Not realistic in practice.

## Theorem (An oracle inequality)

Suppose that least-squares loss satisfies  $\gamma$ -RE condition. Then for  $\lambda_n \geq \max\{2\|\frac{X^T w}{n}\|_\infty, \sqrt{\frac{\log d}{n}}\}$ , any optimal Lasso solution satisfies

$$\|\hat{\theta} - \theta^*\|_2^2 \leq \min_{S \subseteq \{1, \dots, d\}} \left\{ \underbrace{\frac{9}{4} \frac{\lambda_n^2}{\gamma^2} |S|}_{\text{estimation error}} + \underbrace{\frac{2\lambda_n}{\gamma} \|\theta_{S^c}^*\|_1}_{\text{approximation error}} \right\}.$$

- when  $\theta^*$  is exactly sparse, set  $S = \text{supp}(\theta^*)$  to recover previous result

(cf. Bunea et al., 2007; Buhlmann and van de Geer, 2009; Koltchinski et al., 2011)

# Extension to an oracle inequality

Previous theory assumed that  $\theta^*$  was “hard” sparse. Not realistic in practice.

## Theorem (An oracle inequality)

Suppose that least-squares loss satisfies  $\gamma$ -RE condition. Then for  $\lambda_n \geq \max\{2\|\frac{X^T w}{n}\|_\infty, \sqrt{\frac{\log d}{n}}\}$ , any optimal Lasso solution satisfies

$$\|\hat{\theta} - \theta^*\|_2^2 \leq \min_{S \subseteq \{1, \dots, d\}} \left\{ \underbrace{\frac{9}{4} \frac{\lambda_n^2}{\gamma^2} |S|}_{\text{estimation error}} + \underbrace{\frac{2\lambda_n}{\gamma} \|\theta_{S^c}^*\|_1}_{\text{approximation error}} \right\}.$$

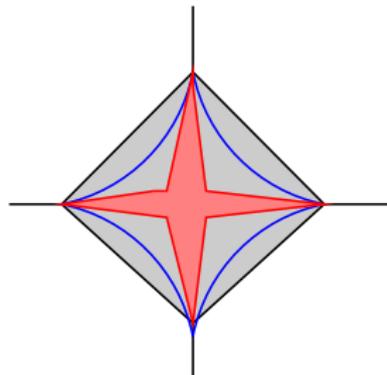
- when  $\theta^*$  is exactly sparse, set  $S = \text{supp}(\theta^*)$  to recover previous result
- more generally, choose  $S$  adaptively to trade-off **estimation error** versus **approximation error**

(cf. Bunea et al., 2007; Buhlmann and van de Geer, 2009; Koltchinski et al., 2011)

# Consequences for $\ell_q$ -“ball” sparsity

- for some  $q \in [0, 1]$ , say  $\theta^*$  belongs to  $\ell_q$ -“ball”

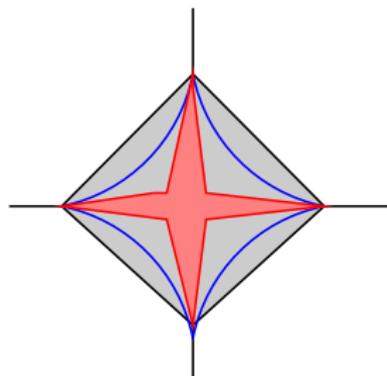
$$\mathbb{B}_q(R_q) := \left\{ \theta \in \mathbb{R}^d \mid \sum_{j=1}^d |\theta_j|^q \leq R_q \right\}.$$



# Consequences for $\ell_q$ -“ball” sparsity

- for some  $q \in [0, 1]$ , say  $\theta^*$  belongs to  $\ell_q$ -“ball”

$$\mathbb{B}_q(R_q) := \left\{ \theta \in \mathbb{R}^d \mid \sum_{j=1}^d |\theta_j|^q \leq R_q \right\}.$$



## Corollary

Consider the linear model  $y = X\theta^* + w$ , where  $X$  satisfies lower RE conditions, and  $w$  has i.i.d  $\sigma$  sub-Gaussian entries. For  $\theta^* \in \mathbb{B}_q(R_q)$ , any Lasso solution satisfies (w.h.p.)

$$\|\hat{\theta} - \theta^*\|_2^2 \lesssim R_q \left( \frac{\sigma^2 \log d}{n} \right)^{1-q/2}.$$

## Are these good results? Minimax theory

- let  $\mathcal{P}$  be a family of probability distributions
- consider a parameter  $\mathbb{P} \mapsto \theta(\mathbb{P})$
- define a metric  $\rho$  on the parameter space

## Are these good results? Minimax theory

- let  $\mathcal{P}$  be a family of probability distributions
- consider a parameter  $\mathbb{P} \mapsto \theta(\mathbb{P})$
- define a metric  $\rho$  on the parameter space

### Definition (Minimax rate)

The minimax rate for  $\theta(\mathcal{P})$  with metric  $\rho$  is given

$$\mathfrak{M}_n(\theta(\mathcal{P}); \rho) := \inf_{\hat{\theta}} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{E}[\rho^2(\hat{\theta}, \theta(\mathbb{P}))],$$

where the infimum ranges over all measurable functions of  $n$  samples.

# Are these good results? Minimax theory

## Definition (Minimax rate)

The minimax rate for  $\theta(\mathcal{P})$  with metric  $\rho$  is given

$$\mathfrak{M}_n(\theta(\mathcal{P}); \rho) := \inf_{\hat{\theta}} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{E}[\rho^2(\hat{\theta}, \theta(\mathbb{P}))],$$

where the infimum ranges over all measurable functions of  $n$  samples.

Concrete example:

- let  $\mathcal{P}$  be family of sparse linear regression problems with  $\theta^* \in \mathbb{B}_q(R_q)$
- consider  $\ell_2$ -error metric  $\rho^2(\hat{\theta}, \theta) = \|\hat{\theta} - \theta\|_2^2$

# Are these good results? Minimax theory

## Definition (Minimax rate)

The minimax rate for  $\theta(\mathcal{P})$  with metric  $\rho$  is given

$$\mathfrak{M}_n(\theta(\mathcal{P}); \rho) := \inf_{\hat{\theta}} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{E}[\rho^2(\hat{\theta}, \theta(\mathbb{P}))],$$

where the infimum ranges over all measurable functions of  $n$  samples.

Concrete example:

- let  $\mathcal{P}$  be family of sparse linear regression problems with  $\theta^* \in \mathbb{B}_q(R_q)$
- consider  $\ell_2$ -error metric  $\rho^2(\hat{\theta}, \theta) = \|\hat{\theta} - \theta\|_2^2$

## Theorem (Raskutti, W. & Yu, 2011)

Under “mild” conditions on design  $X$  and radius  $R_q$ , we have

$$\mathfrak{M}_n(\mathbb{B}_q(R_q); \|\cdot\|_2) \asymp R_q \left( \frac{\sigma^2 \log d}{n} \right)^{1 - \frac{q}{2}}.$$

see Donoho & Johnstone, 1994 for normal sequence model

## Bounds on prediction error

Can predict a new response vector  $y \in \mathbb{R}^n$  via  $\hat{y} = X\hat{\theta}$ . Associated mean-squared error

$$\frac{1}{n} \mathbb{E}[\|y - \hat{y}\|_2^2] = \frac{1}{n} \|X\hat{\theta} - X\theta^*\|_2^2 + \sigma^2.$$

### Theorem

Consider the constrained Lasso with  $R = \|\theta^*\|_1$  or regularized Lasso with  $\lambda_n = 4\sigma\sqrt{\frac{\log d}{n}}$  applied to an  $S$ -sparse problem with  $\sigma$ -sub-Gaussian noise. Then with high probability:

**Slow rate:** If  $X$  has normalized columns ( $\max_j \|X_j\|_2/\sqrt{n} \leq C$ ), then any optimal  $\hat{\theta}$  satisfies the bound

$$\frac{1}{n} \|X\hat{\theta} - X\theta^*\|_2^2 \leq cC R\sigma\sqrt{\frac{\log d}{n}}$$

**Fast rate:** If  $X$  satisfies the  $\gamma$ -RE condition over  $S$ , then

$$\frac{1}{n} \|X\hat{\theta} - X\theta^*\|_2^2 \leq \frac{c\sigma^2}{\gamma} \frac{s \log d}{n}$$

## Prediction error: Proof of slow rate

Let's analyze constrained version:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2n} \|y - X\theta\|_2^2 \quad \text{such that } \|\theta\|_1 \leq R = \|\theta^*\|_1.$$

---

(1) By **optimality of  $\hat{\theta}$**  and **feasibility of  $\theta^*$** , we have the basic inequality for  $\hat{\Delta} = \hat{\theta} - \theta^*$ :

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle.$$

# Prediction error: Proof of slow rate

Let's analyze constrained version:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2n} \|y - X\theta\|_2^2 \quad \text{such that } \|\theta\|_1 \leq R = \|\theta^*\|_1.$$

---

(1) By **optimality of  $\hat{\theta}$**  and **feasibility of  $\theta^*$** , we have the basic inequality for  $\hat{\Delta} = \hat{\theta} - \theta^*$ :

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle.$$

(2) **Hölder's inequality for RHS**

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle \leq 2 \|\hat{\Delta}\|_1 \left\| \frac{X^T w}{n} \right\|_\infty.$$

# Prediction error: Proof of slow rate

Let's analyze constrained version:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2n} \|y - X\theta\|_2^2 \quad \text{such that } \|\theta\|_1 \leq R = \|\theta^*\|_1.$$

---

- (1) By **optimality of  $\hat{\theta}$**  and **feasibility of  $\theta^*$** , we have the basic inequality for  $\hat{\Delta} = \hat{\theta} - \theta^*$ :

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle.$$

- (2) **Hölder's inequality for RHS**

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle \leq 2 \|\hat{\Delta}\|_1 \left\| \frac{X^T w}{n} \right\|_\infty.$$

- (3) Since both  $\hat{\theta}$  and  $\theta^*$  are feasible, we have  $\|\hat{\Delta}\|_1 \leq 2R$  by triangle inequality, and hence

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq 4R \left\| \frac{X^T w}{n} \right\|_\infty.$$

## Prediction error: Proof of fast rate

Let's analyze constrained version:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2n} \|y - X\theta\|_2^2 \quad \text{such that } \|\theta\|_1 \leq R = \|\theta^*\|_1.$$

---

(1) By **optimality of  $\hat{\theta}$**  and **feasibility of  $\theta^*$** , we have the basic inequality for  $\hat{\Delta} = \hat{\theta} - \theta^*$ :

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle.$$

## Prediction error: Proof of fast rate

---

(1) By optimality of  $\hat{\theta}$  and feasibility of  $\theta^*$ , we have the basic inequality for  $\hat{\Delta} = \hat{\theta} - \theta^*$ :

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle.$$

(2) Hölder's inequality for RHS

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle \leq 2 \|\hat{\Delta}\|_1 \left\| \frac{X^T w}{n} \right\|_\infty.$$

## Prediction error: Proof of fast rate

---

(1) By optimality of  $\hat{\theta}$  and feasibility of  $\theta^*$ , we have the basic inequality for  $\hat{\Delta} = \hat{\theta} - \theta^*$ :

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle.$$

(2) Hölder's inequality for RHS

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle \leq 2 \|\hat{\Delta}\|_1 \left\| \frac{X^T w}{n} \right\|_\infty.$$

(3) Since  $\hat{\Delta} \in \mathbb{C}(S)$ , we have  $\|\hat{\Delta}\|_1 \leq 2\sqrt{s}\|\hat{\Delta}\|_2$ , and hence

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq 2\sqrt{s}\|\hat{\Delta}\|_2 \left\| \frac{X^T w}{n} \right\|_\infty.$$

## Prediction error: Proof of fast rate

---

- (1) By optimality of  $\hat{\theta}$  and feasibility of  $\theta^*$ , we have the basic inequality for  $\hat{\Delta} = \hat{\theta} - \theta^*$ :

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle.$$

- (2) Hölder's inequality for RHS

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle \leq 2 \|\hat{\Delta}\|_1 \left\| \frac{X^T w}{n} \right\|_\infty.$$

- (3) Since  $\hat{\Delta} \in \mathbb{C}(S)$ , we have  $\|\hat{\Delta}\|_1 \leq 2\sqrt{s}\|\hat{\Delta}\|_2$ , and hence

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq 2\sqrt{s}\|\hat{\Delta}\|_2 \left\| \frac{X^T w}{n} \right\|_\infty.$$

- (4) Now apply  $\gamma$ -RE condition to RHS

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq 2\sqrt{s}\|\hat{\Delta}\|_2 \left\| \frac{X^T w}{n} \right\|_\infty \leq \frac{1}{\sqrt{\gamma}} \frac{1}{\sqrt{n}} \|X\hat{\Delta}\|_2 \left\| \frac{X^T w}{n} \right\|_\infty.$$

Cancel terms and re-arrange.

# Why RE conditions for fast rate?

## Bothersome issue:

Why should prediction performance depend on an  $RE$ -condition?

- it is not **fundamental**: a method based on  $\ell_0$ -regularization (exponential time) can achieve the fast rate with only column normalization

# Why RE conditions for fast rate?

## Bothersome issue:

Why should prediction performance depend on an *RE*-condition?

- it is not **fundamental**: a method based on  $\ell_0$ -regularization (exponential time) can achieve the fast rate with only column normalization
- some negative evidence: an explicit design matrix and sparse vector ( $k = 2$ ) for which Lasso achieves slow rate Foygel & Srebro (2011)

# Why RE conditions for fast rate?

## Bothersome issue:

Why should prediction performance depend on an *RE*-condition?

- it is not **fundamental**: a method based on  $\ell_0$ -regularization (exponential time) can achieve the fast rate with only column normalization
- some negative evidence: an explicit design matrix and sparse vector ( $k = 2$ ) for which Lasso achieves slow rate Foygel & Srebro (2011)
- ....but adaptive Lasso can achieve the fast rate for this counterexample.

# A computationally-constrained minimax rate

## Complexity classes:

**P**: decision problems solvable in poly. time by a Turing machine

**P/poly**: class **P** plus polynomial-length advice string

## Assumptions:

- standard linear regression model  $y = X\theta^* + w$  where  $w \sim N(0, \sigma^2 I_{n \times n})$
- $\mathbf{NP} \not\subseteq \mathbf{P/poly}$

# A computationally-constrained minimax rate

Complexity classes:

**P**: decision problems solvable in poly. time by a Turing machine

**P/poly**: class **P** plus polynomial-length advice string

Assumptions:

- standard linear regression model  $y = X\theta^* + w$  where  $w \sim N(0, \sigma^2 I_{n \times n})$
- $\text{NP} \not\subseteq \text{P/poly}$

**Theorem (Zhang, W. & Jordan, COLT 2014)**

There is a *fixed “bad” design matrix*  $X \in \mathbb{R}^{n \times d}$  with *RE constant*  $\gamma(X)$  such for any *polynomial-time computable*  $\hat{\theta}$  returning *s-sparse outputs*:

$$\sup_{\theta^* \in \mathbb{B}_0(s)} \mathbb{E} \left[ \frac{\|X(\hat{\theta} - \theta^*)\|_2^2}{n} \right] \gtrsim \frac{\sigma^2}{\gamma^2(X)} \frac{s^{1-\delta} \log d}{n}.$$

# **A primer on high-dimensional statistics: Part II**

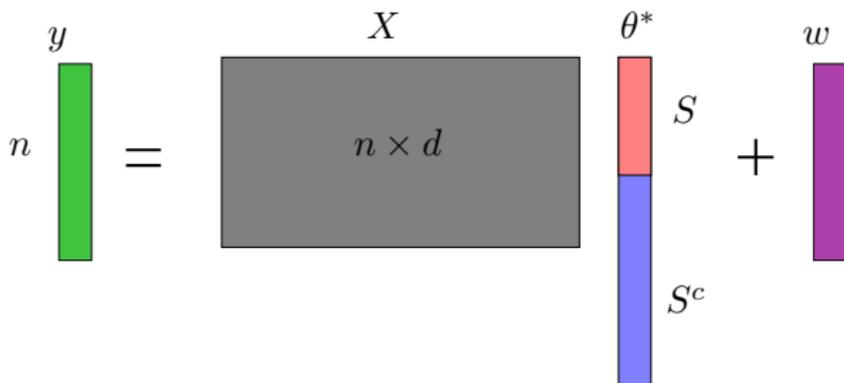
Martin Wainwright

UC Berkeley  
Departments of Statistics, and EECS

Spring School, Westerland, March 2015

# Analysis of Lasso estimator

Turning to noisy observations...



**Estimator:** Lasso program

$$\hat{\theta}_{\lambda_n} \in \arg \min_{\theta \in \mathbb{R}^d} \left\{ \frac{1}{2n} \|y - X\theta\|_2^2 + \lambda_n \|\theta\|_1 \right\}.$$

**Goal:** Obtain bounds on { prediction error, parametric error, variable selection }.

# Different error metrics

1 (In-sample) prediction error:  $\|X(\hat{\theta} - \theta^*)\|_2^2/n$

- ▶ “weakest” error measure
- ▶ appropriate when  $\theta^*$  itself not of primary interest
- ▶ strong dependence between columns of  $X$  possible (for slow rate)
- ▶ proof technique: basic inequality

# Different error metrics

- 1 (In-sample) prediction error:  $\|X(\hat{\theta} - \theta^*)\|_2^2/n$ 
  - ▶ “weakest” error measure
  - ▶ appropriate when  $\theta^*$  itself not of primary interest
  - ▶ strong dependence between columns of  $X$  possible (for slow rate)
  - ▶ proof technique: basic inequality
  
- 2 parametric error:  $\|\hat{\theta} - \theta^*\|_r$  for some  $r \in [1, \infty]$ 
  - ▶ appropriate for recovery problems
  - ▶ RE-type conditions appear in both lower/upper bounds
  - ▶ variable selection is not guaranteed
  - ▶ proof technique: basic inequality

# Different error metrics

- 1 (In-sample) prediction error:  $\|X(\hat{\theta} - \theta^*)\|_2^2/n$ 
  - ▶ “weakest” error measure
  - ▶ appropriate when  $\theta^*$  itself not of primary interest
  - ▶ strong dependence between columns of  $X$  possible (for slow rate)
  - ▶ proof technique: basic inequality
- 2 parametric error:  $\|\hat{\theta} - \theta^*\|_r$  for some  $r \in [1, \infty]$ 
  - ▶ appropriate for recovery problems
  - ▶ RE-type conditions appear in both lower/upper bounds
  - ▶ variable selection is not guaranteed
  - ▶ proof technique: basic inequality
- 3 variable selection: is  $\text{supp}(\hat{\theta})$  equal to  $\text{supp}(\theta^*)$ ?
  - ▶ appropriate when non-zero locations are of scientific interest
  - ▶ most stringent of all three criteria
  - ▶ requires incoherence or irrepresentability conditions on  $X$
  - ▶ proof technique: primal-dual witness condition

# Extension to an oracle inequality

Previous theory assumed that  $\theta^*$  was “hard” sparse. Not realistic in practice.

# Extension to an oracle inequality

Previous theory assumed that  $\theta^*$  was “hard” sparse. Not realistic in practice.

## Theorem (An oracle inequality)

Suppose that least-squares loss satisfies  $\gamma$ -RE condition. Then for  $\lambda_n \geq \max\{2\|\frac{X^T w}{n}\|_\infty, \sqrt{\frac{\log d}{n}}\}$ , any optimal Lasso solution satisfies

$$\|\hat{\theta} - \theta^*\|_2^2 \leq \min_{S \subseteq \{1, \dots, d\}} \left\{ \underbrace{\frac{9}{4} \frac{\lambda_n^2}{\gamma^2} |S|}_{\text{estimation error}} + \underbrace{\frac{2\lambda_n}{\gamma} \|\theta_{S^c}^*\|_1}_{\text{approximation error}} \right\}.$$

(cf. Bunea et al., 2007; Buhlmann and van de Geer, 2009; Koltchinski et al., 2011)

# Extension to an oracle inequality

Previous theory assumed that  $\theta^*$  was “hard” sparse. Not realistic in practice.

## Theorem (An oracle inequality)

Suppose that least-squares loss satisfies  $\gamma$ -RE condition. Then for  $\lambda_n \geq \max\{2\|\frac{X^T w}{n}\|_\infty, \sqrt{\frac{\log d}{n}}\}$ , any optimal Lasso solution satisfies

$$\|\hat{\theta} - \theta^*\|_2^2 \leq \min_{S \subseteq \{1, \dots, d\}} \left\{ \underbrace{\frac{9}{4} \frac{\lambda_n^2}{\gamma^2} |S|}_{\text{estimation error}} + \underbrace{\frac{2\lambda_n}{\gamma} \|\theta_{S^c}^*\|_1}_{\text{approximation error}} \right\}.$$

- when  $\theta^*$  is exactly sparse, set  $S = \text{supp}(\theta^*)$  to recover previous result

(cf. Bunea et al., 2007; Buhlmann and van de Geer, 2009; Koltchinski et al., 2011)

# Extension to an oracle inequality

Previous theory assumed that  $\theta^*$  was “hard” sparse. Not realistic in practice.

## Theorem (An oracle inequality)

Suppose that least-squares loss satisfies  $\gamma$ -RE condition. Then for  $\lambda_n \geq \max\{2\|\frac{X^T w}{n}\|_\infty, \sqrt{\frac{\log d}{n}}\}$ , any optimal Lasso solution satisfies

$$\|\hat{\theta} - \theta^*\|_2^2 \leq \min_{S \subseteq \{1, \dots, d\}} \left\{ \underbrace{\frac{9}{4} \frac{\lambda_n^2}{\gamma^2} |S|}_{\text{estimation error}} + \underbrace{\frac{2\lambda_n}{\gamma} \|\theta_{S^c}^*\|_1}_{\text{approximation error}} \right\}.$$

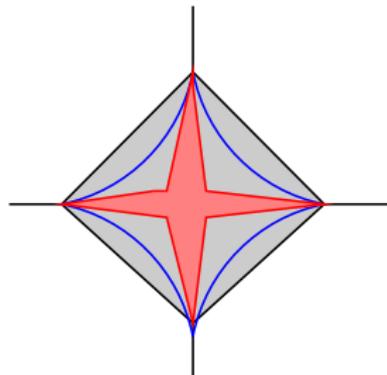
- when  $\theta^*$  is exactly sparse, set  $S = \text{supp}(\theta^*)$  to recover previous result
- more generally, choose  $S$  adaptively to trade-off **estimation error** versus **approximation error**

(cf. Bunea et al., 2007; Buhlmann and van de Geer, 2009; Koltchinski et al., 2011)

# Consequences for $\ell_q$ -“ball” sparsity

- for some  $q \in [0, 1]$ , say  $\theta^*$  belongs to  $\ell_q$ -“ball”

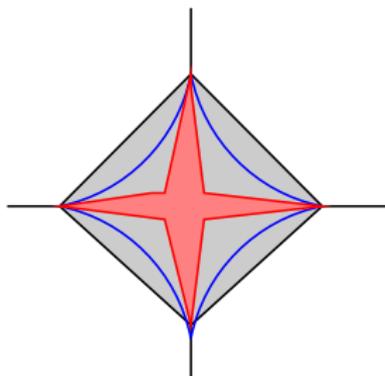
$$\mathbb{B}_q(R_q) := \left\{ \theta \in \mathbb{R}^d \mid \sum_{j=1}^d |\theta_j|^q \leq R_q \right\}.$$



# Consequences for $\ell_q$ -“ball” sparsity

- for some  $q \in [0, 1]$ , say  $\theta^*$  belongs to  $\ell_q$ -“ball”

$$\mathbb{B}_q(R_q) := \left\{ \theta \in \mathbb{R}^d \mid \sum_{j=1}^d |\theta_j|^q \leq R_q \right\}.$$



## Corollary

Consider the linear model  $y = X\theta^* + w$ , where  $X$  satisfies lower RE conditions, and  $w$  has i.i.d  $\sigma$  sub-Gaussian entries. For  $\theta^* \in \mathbb{B}_q(R_q)$ , any Lasso solution satisfies (w.h.p.)

$$\|\hat{\theta} - \theta^*\|_2^2 \lesssim R_q \left( \frac{\sigma^2 \log d}{n} \right)^{1-q/2}.$$

## Are these good results? Minimax theory

- let  $\mathcal{P}$  be a family of probability distributions
- consider a parameter  $\mathbb{P} \mapsto \theta(\mathbb{P})$
- define a metric  $\rho$  on the parameter space

# Are these good results? Minimax theory

- let  $\mathcal{P}$  be a family of probability distributions
- consider a parameter  $\mathbb{P} \mapsto \theta(\mathbb{P})$
- define a metric  $\rho$  on the parameter space

## Definition (Minimax rate)

The minimax rate for  $\theta(\mathcal{P})$  with metric  $\rho$  is given

$$\mathfrak{M}_n(\theta(\mathcal{P}); \rho) := \inf_{\hat{\theta}} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{E}[\rho^2(\hat{\theta}, \theta(\mathbb{P}))],$$

where the infimum ranges over all measurable functions of  $n$  samples.

# Are these good results? Minimax theory

## Definition (Minimax rate)

The minimax rate for  $\theta(\mathcal{P})$  with metric  $\rho$  is given

$$\mathfrak{M}_n(\theta(\mathcal{P}); \rho) := \inf_{\hat{\theta}} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{E}[\rho^2(\hat{\theta}, \theta(\mathbb{P}))],$$

where the infimum ranges over all measurable functions of  $n$  samples.

Concrete example:

- let  $\mathcal{P}$  be family of sparse linear regression problems with  $\theta^* \in \mathbb{B}_q(R_q)$
- consider  $\ell_2$ -error metric  $\rho^2(\hat{\theta}, \theta) = \|\hat{\theta} - \theta\|_2^2$

# Are these good results? Minimax theory

## Definition (Minimax rate)

The minimax rate for  $\theta(\mathcal{P})$  with metric  $\rho$  is given

$$\mathfrak{M}_n(\theta(\mathcal{P}); \rho) := \inf_{\hat{\theta}} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{E}[\rho^2(\hat{\theta}, \theta(\mathbb{P}))],$$

where the infimum ranges over all measurable functions of  $n$  samples.

Concrete example:

- let  $\mathcal{P}$  be family of sparse linear regression problems with  $\theta^* \in \mathbb{B}_q(R_q)$
- consider  $\ell_2$ -error metric  $\rho^2(\hat{\theta}, \theta) = \|\hat{\theta} - \theta\|_2^2$

## Theorem (Raskutti, W. & Yu, 2011)

Under “mild” conditions on design  $X$  and radius  $R_q$ , we have

$$\mathfrak{M}_n(\mathbb{B}_q(R_q); \|\cdot\|_2) \asymp R_q \left( \frac{\sigma^2 \log d}{n} \right)^{1 - \frac{q}{2}}.$$

see Donoho & Johnstone, 1994 for normal sequence model

# Variable selection consistency

## Question

When is Lasso solution unique with  $\text{support}(\hat{\theta}) = \text{support}(\theta^*)$ ?

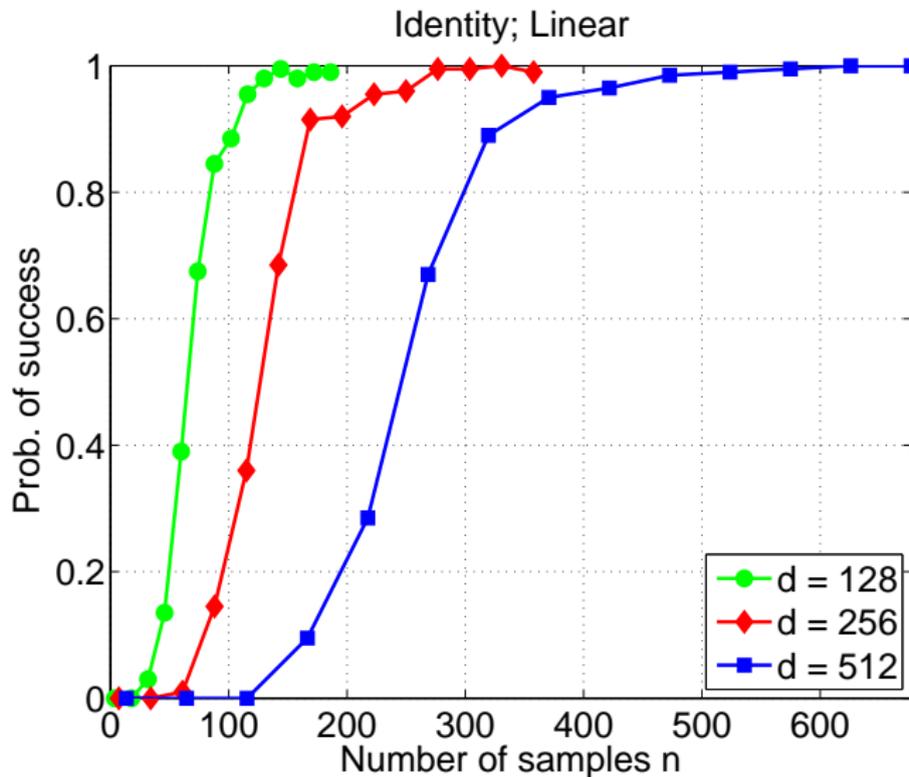
# Variable selection consistency

## Question

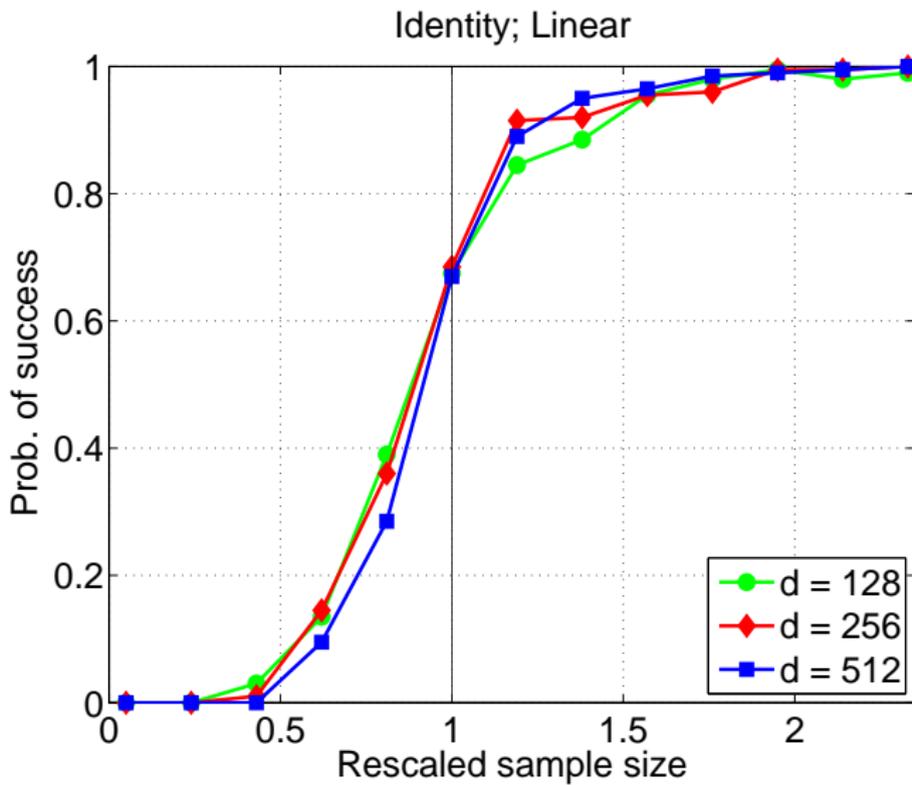
When is Lasso solution unique with  $\text{support}(\hat{\theta}) = \text{support}(\theta^*)$ ?

- Requires a different proof technique, known as a **primal-dual witness method**.
- A procedure that attempts to construct a pair  $(\hat{\theta}, \hat{z}) \in \mathbb{R}^d \times \mathbb{R}^d$  that satisfy the KKT conditions for convex optimality
- When procedure succeeds, it **certifies** the uniqueness and optimality of  $\hat{\theta}$  as a Lasso solution.

# Variable selection performance: unrescaled plots



# Variable selection performance: rescaled plots



Rescaled sample size:  $\frac{n}{s \log(d-s)}$

## Primal-dual witness construction

Consider blocks  $\begin{bmatrix} \hat{\theta}_S & \hat{\theta}_{S^c} \end{bmatrix}$  and  $\begin{bmatrix} \hat{z}_S & \hat{z}_{S^c} \end{bmatrix}$ .

- 1 Set  $\hat{\theta}_{S^c} = 0$ .

## Primal-dual witness construction

Consider blocks  $\begin{bmatrix} \hat{\theta}_S & \hat{\theta}_{S^c} \end{bmatrix}$  and  $\begin{bmatrix} \hat{z}_S & \hat{z}_{S^c} \end{bmatrix}$ .

- 1 Set  $\hat{\theta}_{S^c} = 0$ .
- 2 Solve oracle sub-problem

$$\hat{\theta}_S = \arg \min_{\theta_S \in \mathbb{R}^{|S|}} \left\{ \frac{1}{2n} \|y - X_S \theta_S\|_2^2 + \lambda_n \|\theta_S\|_1 \right\},$$

and choose  $\hat{z} \in \partial \|\theta_S\|_1 \Big|_{\theta_S = \hat{\theta}_S}$  such that  $\frac{1}{n} X_S^T (X_S \hat{\theta}_S - y) + \lambda_n \hat{z}_S = 0$ .  
Require  $\frac{1}{n} \lambda_{\min}(X_S^T X_S) > 0$  in this step.

# Primal-dual witness construction

Consider blocks  $[\hat{\theta}_S \quad \hat{\theta}_{S^c}]$  and  $[\hat{z}_S \quad \hat{z}_{S^c}]$ .

- 1 Set  $\hat{\theta}_{S^c} = 0$ .
- 2 Solve oracle sub-problem

$$\hat{\theta}_S = \arg \min_{\theta_S \in \mathbb{R}^{|S|}} \left\{ \frac{1}{2n} \|y - X_S \theta_S\|_2^2 + \lambda_n \|\theta_S\|_1 \right\},$$

and choose  $\hat{z} \in \partial \|\theta_S\|_1 \Big|_{\theta_S = \hat{\theta}_S}$  such that  $\frac{1}{n} X_S^T (X_S \hat{\theta}_S - y) + \lambda_n \hat{z}_S = 0$ .

Require  $\frac{1}{n} \lambda_{\min}(X_S^T X_S) > 0$  in this step.

- 3 Choose  $\hat{z}_{S^c} \in \mathbb{R}^{d-s}$  to satisfy the zero-subgradient conditions, and such that  $\|\hat{z}_{S^c}\|_\infty < 1$ .

# Primal-dual witness construction

Consider blocks  $[\hat{\theta}_S \quad \hat{\theta}_{S^c}]$  and  $[\hat{z}_S \quad \hat{z}_{S^c}]$ .

- 1 Set  $\hat{\theta}_{S^c} = 0$ .
- 2 Solve oracle sub-problem

$$\hat{\theta}_S = \arg \min_{\theta_S \in \mathbb{R}^{|S|}} \left\{ \frac{1}{2n} \|y - X_S \theta_S\|_2^2 + \lambda_n \|\theta_S\|_1 \right\},$$

and choose  $\hat{z} \in \partial \|\theta_S\|_1 \Big|_{\theta_S = \hat{\theta}_S}$  such that  $\frac{1}{n} X_S^T (X_S \hat{\theta}_S - y) + \lambda_n \hat{z}_S = 0$ .

Require  $\frac{1}{n} \lambda_{\min}(X_S^T X_S) > 0$  in this step.

- 3 Choose  $\hat{z}_{S^c} \in \mathbb{R}^{d-s}$  to satisfy the zero-subgradient conditions, and such that  $\|\hat{z}_{S^c}\|_\infty < 1$ .

## Lemma

If the PDW succeeds, then  $\hat{\theta}$  is the *unique optimal solution* of the Lasso and satisfies  $\text{support}(\hat{\theta}) \subseteq \text{support}(\theta^*)$ .

# Proof sketch

① Form zero-subgradient conditions:

$$\frac{1}{n} \begin{bmatrix} X_S^T X_S & X_{S^c}^T X_S \\ X_{S^c}^T X_S & X_{S^c}^T X_{S^c} \end{bmatrix} \begin{bmatrix} \widehat{\theta}_S - \theta_S^* \\ 0 \end{bmatrix} - \frac{1}{n} \begin{bmatrix} X_S^T \\ X_{S^c}^T \end{bmatrix} w + \lambda_n \begin{bmatrix} \widehat{z}_S \\ \widehat{z}_{S^c} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$$

# Proof sketch

- 1 Form zero-subgradient conditions:

$$\frac{1}{n} \begin{bmatrix} X_S^T X_S & X_{S^c}^T X_S \\ X_{S^c}^T X_S & X_{S^c}^T X_{S^c} \end{bmatrix} \begin{bmatrix} \hat{\theta}_S - \theta_S^* \\ 0 \end{bmatrix} - \frac{1}{n} \begin{bmatrix} X_S^T \\ X_{S^c}^T \end{bmatrix} w + \lambda_n \begin{bmatrix} \hat{z}_S \\ \hat{z}_{S^c} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$$

- 2 Solve for  $\hat{\theta}_S - \theta_S^*$ :

$$\underbrace{\hat{\theta}_S - \theta_S^*}_{U_S} = -(X_S^T X_S)^{-1} X_S^T \mathbf{w} - \lambda_n n (X_S^T X_S)^{-1} z_S.$$

# Proof sketch

- 1 Form zero-subgradient conditions:

$$\frac{1}{n} \begin{bmatrix} X_S^T X_S & X_{S^c}^T X_S \\ X_{S^c}^T X_S & X_{S^c}^T X_{S^c} \end{bmatrix} \begin{bmatrix} \hat{\theta}_S - \theta_S^* \\ 0 \end{bmatrix} - \frac{1}{n} \begin{bmatrix} X_S^T \\ X_{S^c}^T \end{bmatrix} w + \lambda_n \begin{bmatrix} \hat{z}_S \\ \hat{z}_{S^c} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$$

- 2 Solve for  $\hat{\theta}_S - \theta_S^*$ :

$$\underbrace{\hat{\theta}_S - \theta_S^*}_{U_S} = -(X_S^T X_S)^{-1} X_S^T w - \lambda_n n (X_S^T X_S)^{-1} z_S.$$

- 3 Solve for  $z_{S^c}$ :

$$z_{S^c} = \underbrace{X_{S^c}^T X_S (X_S^T X_S)^{-1} z_S}_{\mu} + \underbrace{X_{S^c}^T \left[ I - X_S (X_S^T X_S)^{-1} X_S^T \right]}_{V_{S^c}} \left( \frac{w}{\lambda_n n} \right).$$

Checking that  $\|z_{S^c}\|_\infty < 1$  requires **irrepresentability condition**

$$\max_{j \in S^c} X_j^T \|X_S (X_S^T X_S)^{-1}\|_1 < \alpha < 1.$$

## Bounds on prediction error

Can predict a new response vector  $y \in \mathbb{R}^n$  via  $\hat{y} = X\hat{\theta}$ . Associated mean-squared error

$$\frac{1}{n} \mathbb{E}[\|y - \hat{y}\|_2^2] = \frac{1}{n} \|X\hat{\theta} - X\theta^*\|_2^2 + \sigma^2.$$

### Theorem

Consider the constrained Lasso with  $R = \|\theta^*\|_1$  or regularized Lasso with  $\lambda_n = 4\sigma\sqrt{\frac{\log d}{n}}$  applied to an  $S$ -sparse problem with  $\sigma$ -sub-Gaussian noise. Then with high probability:

**Slow rate:** If  $X$  has normalized columns ( $\max_j \|X_j\|_2/\sqrt{n} \leq C$ ), then any optimal  $\hat{\theta}$  satisfies the bound

$$\frac{1}{n} \|X\hat{\theta} - X\theta^*\|_2^2 \leq cC R\sigma\sqrt{\frac{\log d}{n}}$$

**Fast rate:** If  $X$  satisfies the  $\gamma$ -RE condition over  $S$ , then

$$\frac{1}{n} \|X\hat{\theta} - X\theta^*\|_2^2 \leq \frac{c\sigma^2}{\gamma} \frac{s \log d}{n}$$

## Prediction error: Proof of slow rate

Let's analyze constrained version:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2n} \|y - X\theta\|_2^2 \quad \text{such that } \|\theta\|_1 \leq R = \|\theta^*\|_1.$$

---

(1) By **optimality of  $\hat{\theta}$**  and **feasibility of  $\theta^*$** , we have the basic inequality for  $\hat{\Delta} = \hat{\theta} - \theta^*$ :

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle.$$

# Prediction error: Proof of slow rate

Let's analyze constrained version:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2n} \|y - X\theta\|_2^2 \quad \text{such that } \|\theta\|_1 \leq R = \|\theta^*\|_1.$$

---

(1) By **optimality of  $\hat{\theta}$**  and **feasibility of  $\theta^*$** , we have the basic inequality for  $\hat{\Delta} = \hat{\theta} - \theta^*$ :

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle.$$

(2) **Hölder's inequality for RHS**

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle \leq 2 \|\hat{\Delta}\|_1 \left\| \frac{X^T w}{n} \right\|_\infty.$$

# Prediction error: Proof of slow rate

Let's analyze constrained version:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2n} \|y - X\theta\|_2^2 \quad \text{such that } \|\theta\|_1 \leq R = \|\theta^*\|_1.$$

---

- (1) By **optimality of  $\hat{\theta}$**  and **feasibility of  $\theta^*$** , we have the basic inequality for  $\hat{\Delta} = \hat{\theta} - \theta^*$ :

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle.$$

- (2) **Hölder's inequality for RHS**

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle \leq 2 \|\hat{\Delta}\|_1 \left\| \frac{X^T w}{n} \right\|_\infty.$$

- (3) Since both  $\hat{\theta}$  and  $\theta^*$  are feasible, we have  $\|\hat{\Delta}\|_1 \leq 2R$  by triangle inequality, and hence

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq 4R \left\| \frac{X^T w}{n} \right\|_\infty.$$

## Prediction error: Proof of fast rate

Let's analyze constrained version:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2n} \|y - X\theta\|_2^2 \quad \text{such that } \|\theta\|_1 \leq R = \|\theta^*\|_1.$$

---

(1) By **optimality of  $\hat{\theta}$**  and **feasibility of  $\theta^*$** , we have the basic inequality for  $\hat{\Delta} = \hat{\theta} - \theta^*$ :

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle.$$

## Prediction error: Proof of fast rate

---

(1) By optimality of  $\hat{\theta}$  and feasibility of  $\theta^*$ , we have the basic inequality for  $\hat{\Delta} = \hat{\theta} - \theta^*$ :

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle.$$

(2) Hölder's inequality for RHS

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle \leq 2 \|\hat{\Delta}\|_1 \left\| \frac{X^T w}{n} \right\|_\infty.$$

## Prediction error: Proof of fast rate

---

(1) By optimality of  $\hat{\theta}$  and feasibility of  $\theta^*$ , we have the basic inequality for  $\hat{\Delta} = \hat{\theta} - \theta^*$ :

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle.$$

(2) Hölder's inequality for RHS

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle \leq 2 \|\hat{\Delta}\|_1 \left\| \frac{X^T w}{n} \right\|_\infty.$$

(3) Since  $\hat{\Delta} \in \mathbb{C}(S)$ , we have  $\|\hat{\Delta}\|_1 \leq 2\sqrt{s}\|\hat{\Delta}\|_2$ , and hence

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq 2\sqrt{s}\|\hat{\Delta}\|_2 \left\| \frac{X^T w}{n} \right\|_\infty.$$

## Prediction error: Proof of fast rate

---

- (1) By optimality of  $\hat{\theta}$  and feasibility of  $\theta^*$ , we have the basic inequality for  $\hat{\Delta} = \hat{\theta} - \theta^*$ :

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle.$$

- (2) Hölder's inequality for RHS

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq \frac{2}{n} \langle \hat{\Delta}, X^T w \rangle \leq 2 \|\hat{\Delta}\|_1 \left\| \frac{X^T w}{n} \right\|_\infty.$$

- (3) Since  $\hat{\Delta} \in \mathbb{C}(S)$ , we have  $\|\hat{\Delta}\|_1 \leq 2\sqrt{s}\|\hat{\Delta}\|_2$ , and hence

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq 2\sqrt{s}\|\hat{\Delta}\|_2 \left\| \frac{X^T w}{n} \right\|_\infty.$$

- (4) Now apply  $\gamma$ -RE condition to RHS

$$\frac{1}{n} \|X\hat{\Delta}\|_2^2 \leq 2\sqrt{s}\|\hat{\Delta}\|_2 \left\| \frac{X^T w}{n} \right\|_\infty \leq \frac{1}{\sqrt{\gamma}} \frac{1}{\sqrt{n}} \|X\hat{\Delta}\|_2 \left\| \frac{X^T w}{n} \right\|_\infty.$$

Cancel terms and re-arrange.

# Why RE conditions for fast rate?

## Bothersome issue:

Why should prediction performance depend on an *RE*-condition?

- it is not **fundamental**: a method based on  $\ell_0$ -regularization (exponential time) can achieve the fast rate with only column normalization

# Why RE conditions for fast rate?

## Bothersome issue:

Why should prediction performance depend on an *RE*-condition?

- it is not **fundamental**: a method based on  $\ell_0$ -regularization (exponential time) can achieve the fast rate with only column normalization
- some negative evidence: an explicit design matrix and sparse vector ( $k = 2$ ) for which Lasso achieves slow rate Foygel & Srebro (2011)

# Why RE conditions for fast rate?

## Bothersome issue:

Why should prediction performance depend on an *RE*-condition?

- it is not **fundamental**: a method based on  $\ell_0$ -regularization (exponential time) can achieve the fast rate with only column normalization
- some negative evidence: an explicit design matrix and sparse vector ( $k = 2$ ) for which Lasso achieves slow rate Foygel & Srebro (2011)
- ....but adaptive Lasso can achieve the fast rate for this counterexample.

# A computationally-constrained minimax rate

## Complexity classes:

**P**: decision problems solvable in poly. time by a Turing machine

**P/poly**: class **P** plus polynomial-length advice string

## Assumptions:

- standard linear regression model  $y = X\theta^* + w$  where  $w \sim N(0, \sigma^2 I_{n \times n})$
- $\mathbf{NP} \not\subseteq \mathbf{P/poly}$

# A computationally-constrained minimax rate

Complexity classes:

**P**: decision problems solvable in poly. time by a Turing machine

**P/poly**: class **P** plus polynomial-length advice string

Assumptions:

- standard linear regression model  $y = X\theta^* + w$  where  $w \sim N(0, \sigma^2 I_{n \times n})$
- $\text{NP} \not\subseteq \text{P/poly}$

**Theorem (Zhang, W. & Jordan, COLT 2014)**

There is a *fixed “bad” design matrix*  $X \in \mathbb{R}^{n \times d}$  with *RE constant*  $\gamma(X)$  such for any *polynomial-time computable*  $\hat{\theta}$  returning *s-sparse outputs*:

$$\sup_{\theta^* \in \mathbb{B}_0(s)} \mathbb{E} \left[ \frac{\|X(\hat{\theta} - \theta^*)\|_2^2}{n} \right] \gtrsim \frac{\sigma^2}{\gamma^2(X)} \frac{s^{1-\delta} \log d}{n}.$$

# High-level overview

## Regularized $M$ -estimators:

Many statistical estimators take the form:

$$\underbrace{\hat{\theta}_{\lambda_n}}_{\text{Estimate}} \in \arg \min_{\theta \in \Omega} \left\{ \underbrace{\mathcal{L}(\theta; Z_1^n)}_{\text{Loss function}} + \lambda_n \underbrace{\mathcal{R}(\theta)}_{\text{Regularizer}} \right\}.$$

# High-level overview

## Regularized $M$ -estimators:

Many statistical estimators take the form:

$$\underbrace{\hat{\theta}_{\lambda_n}}_{\text{Estimate}} \in \arg \min_{\theta \in \Omega} \left\{ \underbrace{\mathcal{L}(\theta; Z_1^n)}_{\text{Loss function}} + \lambda_n \underbrace{\mathcal{R}(\theta)}_{\text{Regularizer}} \right\}.$$

Past years have witnessed an explosion of results (compressed sensing, covariance estimation, block-sparsity, graphical models, matrix completion...)

## Question:

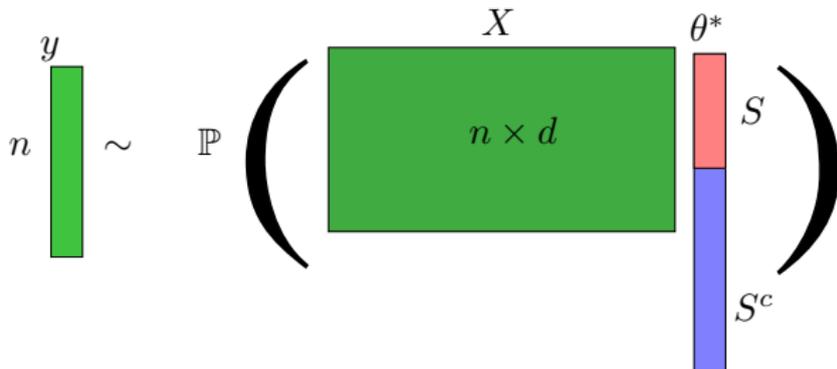
Is there a common set of underlying principles?

# Up until now: Sparse regression

**Set-up:** **Observe**  $(y_i, x_i)$  pairs for  $i = 1, 2, \dots, n$ , where

$$y_i \sim \mathbb{P}(\cdot \mid \langle \theta^*, x_i \rangle),$$

where  $\theta \in \mathbb{R}^d$  is sparse.

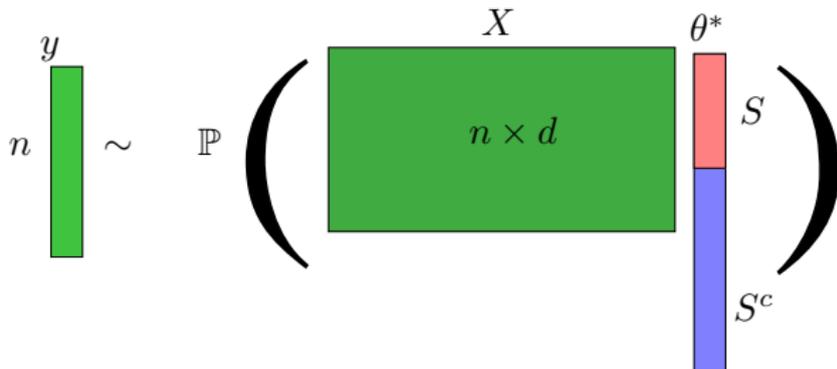


## Up until now: Sparse regression

**Set-up:** **Observe**  $(y_i, x_i)$  pairs for  $i = 1, 2, \dots, n$ , where

$$y_i \sim \mathbb{P}(\cdot \mid \langle \theta^*, x_i \rangle),$$

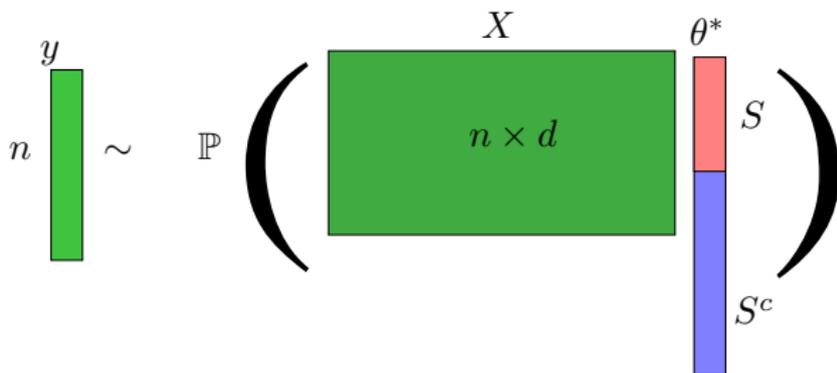
where  $\theta \in \mathbb{R}^d$  is sparse.



**Estimator:**  $\ell_1$ -regularized likelihood

$$\hat{\theta} \in \arg \min_{\theta} \left\{ -\frac{1}{n} \sum_{i=1}^n \log \mathbb{P}(y_i \mid \langle x_i, \theta \rangle) + \lambda \|\theta\|_1 \right\}.$$

## Up until now: Sparse regression



**Example:** Logistic regression for binary responses  $y_i \in \{0, 1\}$ :

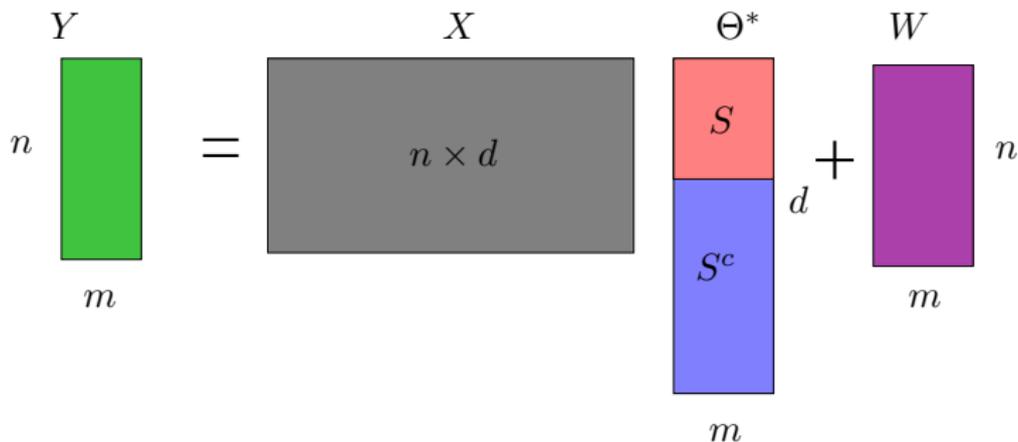
$$\hat{\theta} \in \arg \min_{\theta} \left\{ \frac{1}{n} \sum_{i=1}^n \{ \log(1 + e^{\langle x_i, \theta \rangle}) - y_i \langle x_i, \theta \rangle \} + \lambda \|\theta\|_1 \right\}.$$

## Example: Block sparsity and group Lasso

The diagram illustrates the equation  $Y = X\Theta^* + W$ . Matrix  $Y$  is a green vertical rectangle with height  $n$  and width  $m$ . Matrix  $X$  is a gray horizontal rectangle with height  $n$  and width  $d$ , labeled  $n \times d$ . Matrix  $\Theta^*$  is a vertical rectangle of height  $m$  and width  $d$ , partitioned into a red top block  $S$  and a blue bottom block  $S^c$ . Matrix  $W$  is a purple vertical rectangle with height  $n$  and width  $m$ . The equation is shown as  $Y = X\Theta^* + W$ .

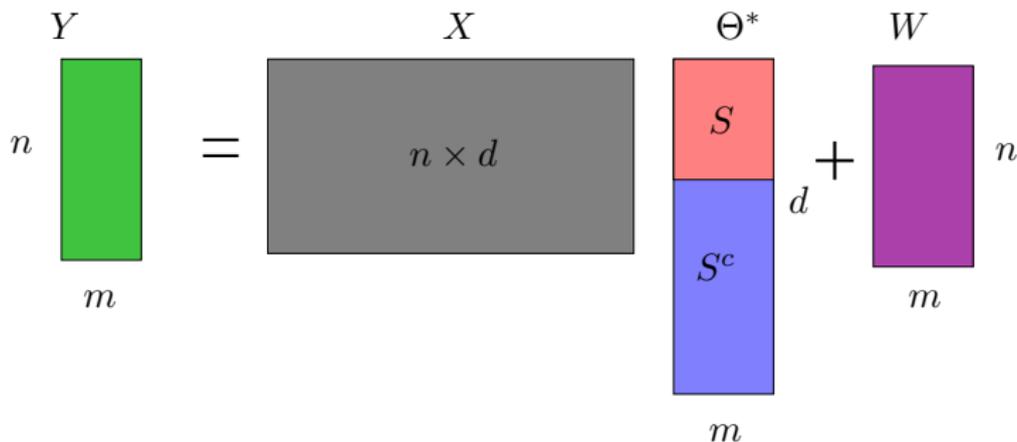
- Matrix  $\Theta^*$  partitioned into **non-zero rows**  $S$  and **zero rows**  $S^c$
- Various applications: multiple-view imaging, gene array prediction, graphical model fitting.

## Example: Block sparsity and group Lasso



- Matrix  $\Theta^*$  partitioned into **non-zero rows**  $S$  and **zero rows**  $S^c$
- Various applications: multiple-view imaging, gene array prediction, graphical model fitting.
- Row-wise  $\ell_1/\ell_2$ -norm  $\|\Theta\|_{1,2} = \sum_{j=1}^d \|\Theta_j\|_2$

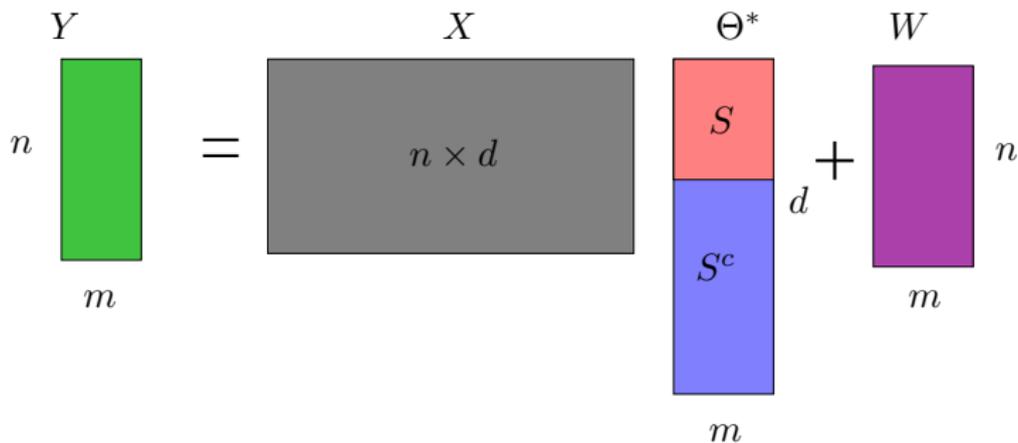
## Example: Block sparsity and group Lasso



- Row-wise  $\ell_1/\ell_2$ -norm  $\|\Theta\|_{1,2} = \sum_{j=1}^d \|\Theta_j\|_2$
- Weighted  $r$ -group Lasso: (Wright et al., 2005; Tropp et al., 2006; Yuan & Lin, 2006)

$$\|\Theta^*\|_{\mathcal{G},r} = \sum_{g \in \mathcal{G}} \omega_g \|\Theta_g\|_r \quad \text{for some } r \in [2, \infty].$$

## Example: Block sparsity and group Lasso

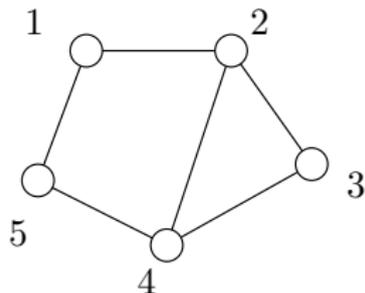
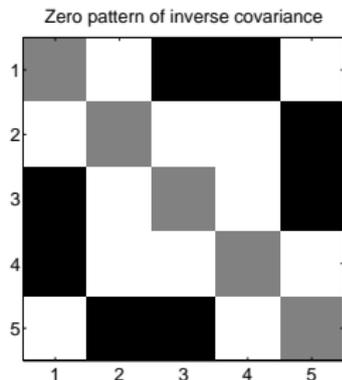


- Row-wise  $\ell_1/\ell_2$ -norm  $\|\Theta\|_{1,2} = \sum_{j=1}^d \|\Theta_j\|_2$
- Weighted  $r$ -group Lasso: (Wright et al., 2005; Tropp et al., 2006; Yuan & Lin, 2006)

$$\|\Theta^*\|_{\mathcal{G},r} = \sum_{g \in \mathcal{G}} \omega_g \|\Theta_g\|_r \quad \text{for some } r \in [2, \infty].$$

- Extensions to { hierarchical, graph-based } groups (e.g., Zhao et al., 2006; Bach et al., 2009; Baraniuk et al., 2009)

# Example: Structured (inverse) covariance matrices



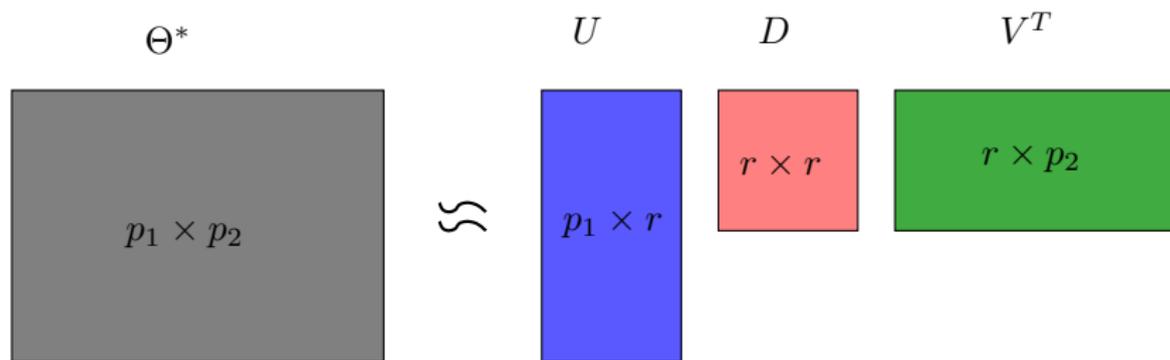
**Set-up:** Samples from random vector with sparse covariance  $\Sigma$  or sparse inverse covariance  $\Theta^* \in \mathbb{R}^{d \times d}$ .

**Estimator** (for inverse covariance)

$$\hat{\Theta} \in \arg \min_{\Theta} \left\{ \left\langle \frac{1}{n} \sum_{i=1}^n x_i x_i^T, \Theta \right\rangle - \log \det(\Theta) + \lambda_n \sum_{j \neq k} |\Theta_{jk}| \right\}$$

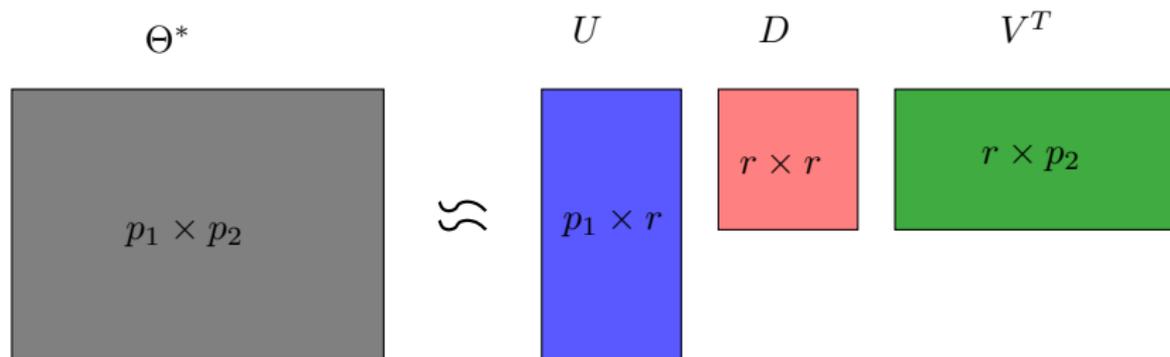
Some past work: Yuan & Lin, 2006; d'Aspremont et al., 2007; Bickel & Levina, 2007; El Karoui, 2007; d'Aspremont et al., 2007; Rothman et al., 2007; Zhou et al., 2007; Friedman et al., 2008; Lam & Fan, 2008; Ravikumar et al., 2008; Zhou, Cai & Huang, 2009; Guo et

## Example: Low-rank matrix approximation



**Set-up:** Matrix  $\Theta^* \in \mathbb{R}^{p_1 \times p_2}$  with rank  $r \ll \min\{p_1, p_2\}$ .

## Example: Low-rank matrix approximation



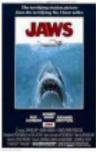
**Set-up:** Matrix  $\Theta^* \in \mathbb{R}^{p_1 \times p_2}$  with rank  $r \ll \min\{p_1, p_2\}$ .

**Least-squares matrix regression:** Given observations  $y_i = \langle X_i, \Theta^* \rangle + w_i$ , solve:

$$\hat{\Theta} \in \arg \min_{\Theta} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - \langle X_i, \Theta \rangle)^2 + \lambda_n \sum_{j=1}^{\min\{p_1, p_2\}} \gamma_j(\Theta) \right\}$$

Some past work: Fazel, 2001; Srebro et al., 2004; Recht, Fazel & Parillo, 2007; Bach, 2008; Candes & Recht, 2008; Keshavan et al., 2009; Rohde & Tsybakov, 2010; Recht, 2009; Negahban & W., 2010; Koltchinski et al., 2011

# Application: Collaborative filtering

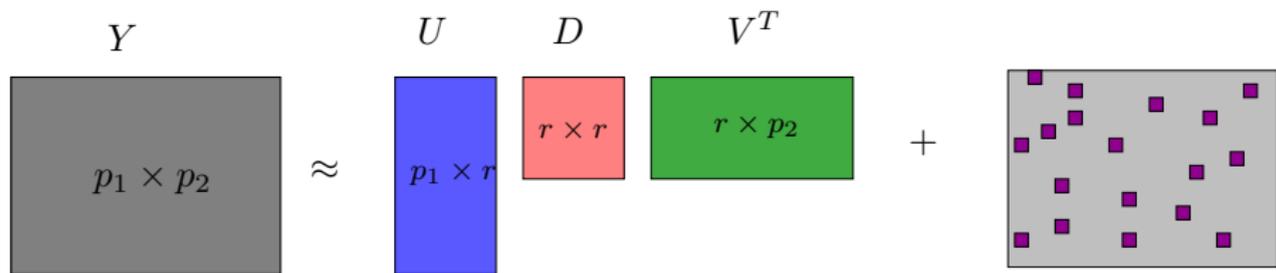
				...	...	
	4	*	3	...	...	*
	3	5	*	...	...	2
	5	4	3	...	...	3
	2	*	*	...	...	1

Universe of  $p_1$  individuals and  $p_2$  films    Observe  $n \ll p_1 p_2$  ratings

(e.g., Srebro, Alon & Jaakkola, 2004; Candes & Recht, 2008)

# Example: Additive matrix decomposition

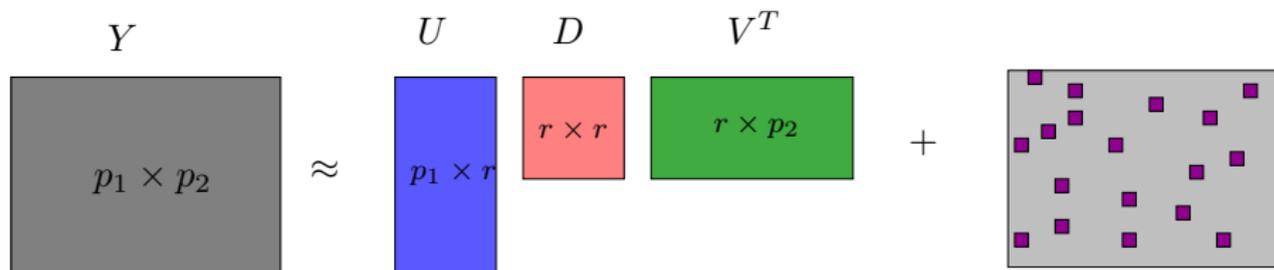
Matrix  $Y$  can be (approximately) decomposed into sum:



$$Y = \underbrace{\Theta^*}_{\text{Low-rank component}} + \underbrace{\Gamma^*}_{\text{Sparse component}}$$

# Example: Additive matrix decomposition

Matrix  $Y$  can be (approximately) decomposed into sum:

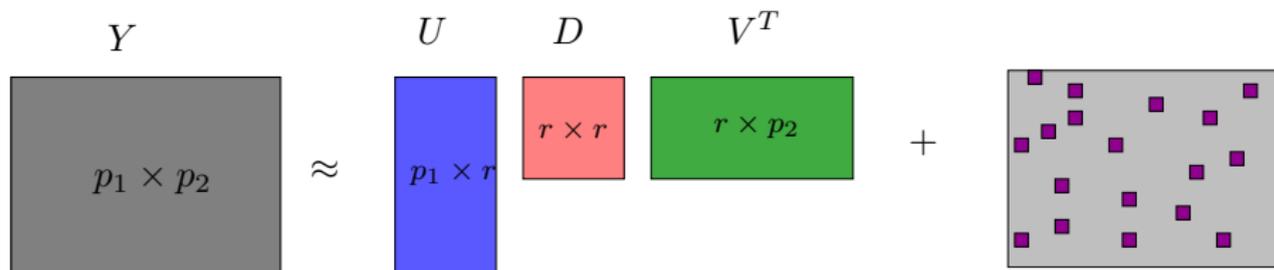


$$Y = \underbrace{\Theta^*}_{\text{Low-rank component}} + \underbrace{\Gamma^*}_{\text{Sparse component}}$$

- Initially proposed by Chandrasekaran, Sanghavi, Parillo & Willsky, 2009
- Various applications:
  - ▶ robust collaborative filtering
  - ▶ robust PCA
  - ▶ graphical model selection with hidden variables

# Example: Additive matrix decomposition

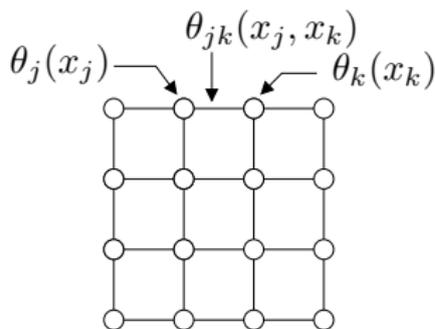
Matrix  $Y$  can be (approximately) decomposed into sum:



$$Y = \underbrace{\Theta^*}_{\text{Low-rank component}} + \underbrace{\Gamma^*}_{\text{Sparse component}}$$

- Initially proposed by Chandrasekaran, Sanghavi, Parillo & Willsky, 2009
- Various applications:
  - ▶ robust collaborative filtering
  - ▶ robust PCA
  - ▶ graphical model selection with hidden variables
- subsequent work: Candes et al., 2010; Xu et al., 2010; Hsu et al., 2010; Agarwal et al., 2011

## Example: Discrete Markov random fields



**Set-up:** Samples from discrete MRF (e.g., Ising or Potts model):

$$\mathbb{P}_\theta(x_1, \dots, x_d) = \frac{1}{Z(\theta)} \exp \left\{ \sum_{j \in V} \theta_j(x_j) + \sum_{(j,k) \in E} \theta_{jk}(x_j, x_k) \right\}.$$

**Estimator:** Given empirical marginal distributions  $\{\hat{\mu}_j, \hat{\mu}_{jk}\}$ :

$$\hat{\Theta} \in \arg \min_{\Theta} \left\{ \sum_{s \in V} \mathbb{E}_{\hat{\mu}_j} [\theta_j(x_j)] + \sum_{(j,k)} \mathbb{E}_{\hat{\mu}_{jk}} [\theta_{jk}(x_j, x_k)] - \log Z(\theta) + \lambda_n \sum_{(j,k)} \|\theta_{jk}\|_F \right\}$$

Some past work: Spirtes et al., 2001; Abbeel et al., 2005; Csiszar & Telata, 2005; Ravikumar et al., 2007; Schneidman et al., 2007; Santhanam & Wainwright, 2008; Sly et al., 2008; Montanari and Pereira, 2009; Anandkumar et al., 2010

# Non-parametric problems: Sparse additive models

- non-parametric regression: **severe** curse of dimensionality!
- many structured classes of non-parametric models are possible:

# Non-parametric problems: Sparse additive models

- non-parametric regression: **severe** curse of dimensionality!
- many structured classes of non-parametric models are possible:
  - ▶ additive models  $f^*(x) = \sum_{j=1}^d f_j^*(x_j)$  (Stone, 1985)
  - ▶ multiple-index models  $f^*(x) = g(B^*x)$

# Non-parametric problems: Sparse additive models

- non-parametric regression: **severe** curse of dimensionality!
- many structured classes of non-parametric models are possible:
  - ▶ additive models  $f^*(x) = \sum_{j=1}^d f_j^*(x_j)$  (Stone, 1985)
  - ▶ multiple-index models  $f^*(x) = g(B^*x)$
  - ▶ sparse additive models:

$$f^*(x) = \sum_{j \in S} f_j^*(x_j) \quad \text{for unknown subset } S$$

(Lin & Zhang, 2003; Meier et al., 2007; Ravikumar et al. 2007; Koltchinski and Yuan, 2008; Raskutti et al., 2010)

# Non-parametric problems: Sparse additive models

Sparse additive models:

$$f^*(x) = \sum_{j \in S}^d f_j^*(x_j) \quad \text{for unknown subset } S$$

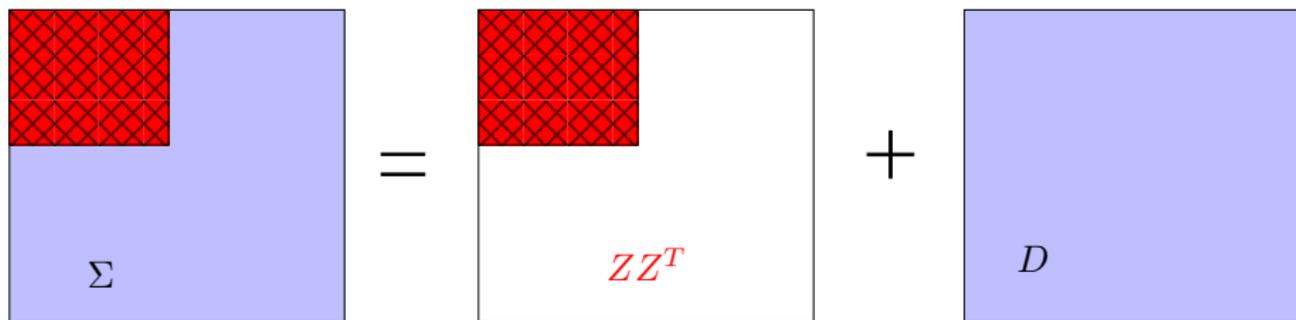
(Lin & Zhang, 2003; Meier et al., 2007; Ravikumar et al. 2007; Koltchinski and Yuan, 2008; Raskutti, W., & Yu, 2010)

Noisy observations  $y_i = f^*(x_i) + w_i$  for  $i = 1, \dots, n$ .

**Estimator:**

$$\hat{f} \in \arg \min_{f = \sum_{j=1}^d f_j} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - \sum_{j=1}^d f_j(x_{ij}))^2 + \lambda \underbrace{\sum_{j=1}^d \|f_j\|_{\mathcal{H}}}_{\|f\|_{1, \mathcal{H}}} + \mu_n \underbrace{\sum_{j=1}^d \|f_j\|_n}_{\|f\|_{1, n}} \right\}.$$

## Example: Sparse principal components analysis



**Set-up:** Covariance matrix  $\Sigma = ZZ^T + D$ , where leading eigenspace  $Z$  has sparse columns.

**Estimator:**

$$\hat{\Theta} \in \arg \min_{\Theta} \left\{ -\langle \Theta, \hat{\Sigma} \rangle + \lambda_n \sum_{(j,k)} |\Theta_{jk}| \right\}$$

Some past work: Johnstone, 2001; Joliffe et al., 2003; Johnstone & Lu, 2004; Zou et al., 2004; d'Asprémont et al., 2007; Johnstone & Paul, 2008; Amini & Wainwright, 2008; Ma, 2012; Berthet & Rigollet, 2012; Nadler et al., 2012

# Motivation and roadmap

- many results on different high-dimensional models
- all based on estimators of the type:

$$\underbrace{\hat{\theta}_{\lambda_n}}_{\text{Estimate}} \in \arg \min_{\theta \in \Omega} \left\{ \underbrace{\mathcal{L}(\theta; Z_1^n)}_{\text{Loss function}} + \lambda_n \underbrace{\mathcal{R}(\theta)}_{\text{Regularizer}} \right\}.$$

# Motivation and roadmap

- many results on different high-dimensional models
- all based on estimators of the type:

$$\underbrace{\hat{\theta}_{\lambda_n}}_{\text{Estimate}} \in \arg \min_{\theta \in \Omega} \left\{ \underbrace{\mathcal{L}(\theta; Z_1^n)}_{\text{Loss function}} + \lambda_n \underbrace{\mathcal{R}(\theta)}_{\text{Regularizer}} \right\}.$$

## Question:

Is there a common set of underlying principles?

# Motivation and roadmap

- many results on different high-dimensional models
- all based on estimators of the type:

$$\underbrace{\hat{\theta}_{\lambda_n}}_{\text{Estimate}} \in \arg \min_{\theta \in \Omega} \left\{ \underbrace{\mathcal{L}(\theta; Z_1^n)}_{\text{Loss function}} + \lambda_n \underbrace{\mathcal{R}(\theta)}_{\text{Regularizer}} \right\}.$$

## Question:

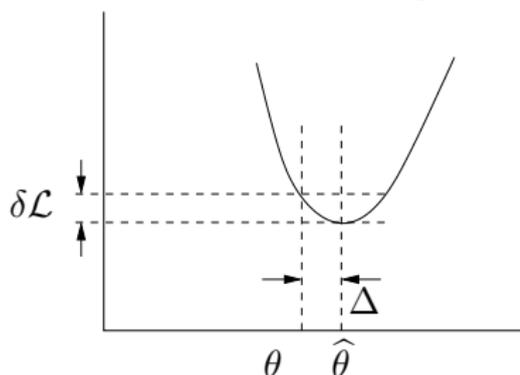
Is there a common set of underlying principles?

**Answer:** Yes, two essential ingredients.

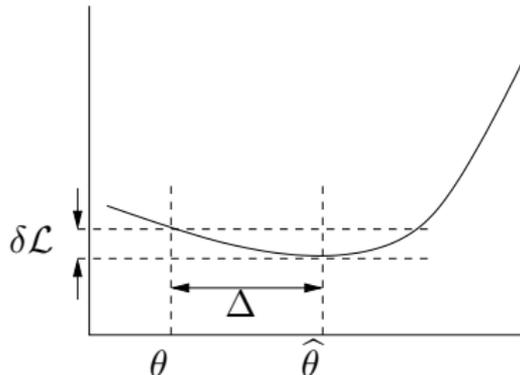
- (I) Restricted strong convexity of loss function
- (II) Decomposability of the regularizer

# (I) Classical role of curvature in statistics

1 Curvature controls difficulty of estimation:



High curvature: easy to estimate



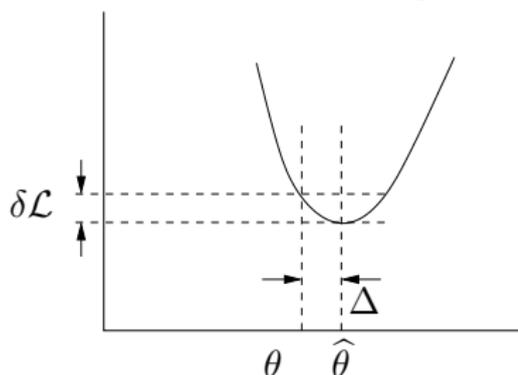
(b) Low curvature: harder

## Canonical example:

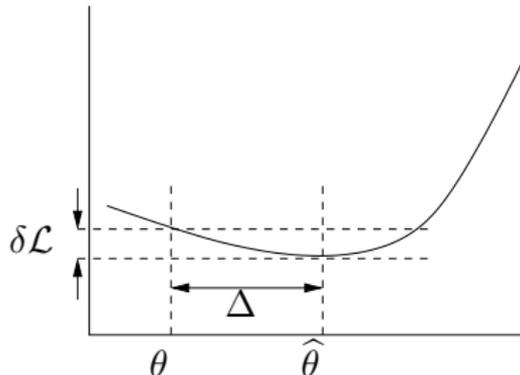
Log likelihood, Fisher information matrix and Cramér-Rao bound.

# (I) Classical role of curvature in statistics

- 1 Curvature controls difficulty of estimation:



High curvature: easy to estimate



(b) Low curvature: harder

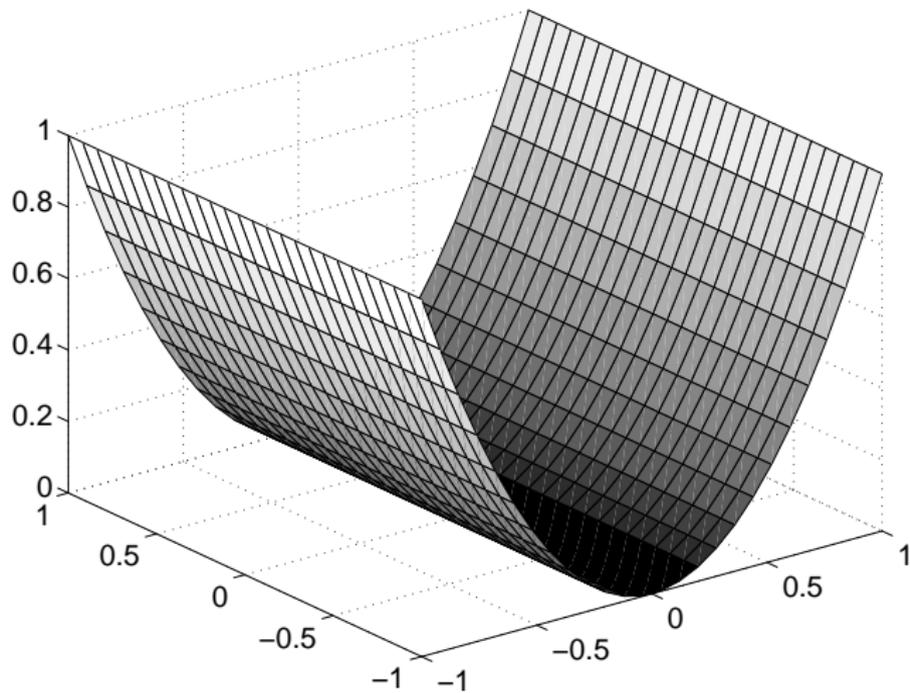
## Canonical example:

Log likelihood, Fisher information matrix and Cramér-Rao bound.

- 2 Formalized by lower bound on Taylor series error  $\mathcal{E}_n(\Delta)$

$$\underbrace{\mathcal{L}(\theta^* + \Delta) - \mathcal{L}(\theta^*) - \langle \nabla \mathcal{L}(\theta^*), \Delta \rangle}_{\mathcal{E}_n(\Delta)} \geq \gamma^2 \|\Delta\|_*^2 \quad \text{for all } \Delta \text{ around } \theta^*.$$

# High dimensions: no strong convexity!



When  $d > n$ , the Hessian  $\nabla^2 \mathcal{L}(\theta; Z_1^n)$  has nullspace of dimension  $d - n$ .

# Restricted strong convexity

## Definition

Loss function  $\mathcal{L}_n$  satisfies restricted strong convexity (RSC) with respect to regularizer  $\mathcal{R}$  if

$$\underbrace{\mathcal{L}_n(\theta^* + \Delta) - \left\{ \mathcal{L}_n(\theta^*) + \langle \nabla \mathcal{L}_n(\theta^*), \Delta \rangle \right\}}_{\text{Taylor error } \mathcal{E}_n(\Delta)} \geq \underbrace{\gamma_\ell \|\Delta\|_2^2}_{\text{Lower curvature}} - \underbrace{\tau_\ell^2 \mathcal{R}^2(\Delta)}_{\text{Tolerance}}$$

for all  $\Delta$  in a suitable neighborhood of  $\theta^*$ .

# Restricted strong convexity

## Definition

Loss function  $\mathcal{L}_n$  satisfies restricted strong convexity (RSC) with respect to regularizer  $\mathcal{R}$  if

$$\underbrace{\mathcal{L}_n(\theta^* + \Delta) - \left\{ \mathcal{L}_n(\theta^*) + \langle \nabla \mathcal{L}_n(\theta^*), \Delta \rangle \right\}}_{\text{Taylor error } \mathcal{E}_n(\Delta)} \geq \underbrace{\gamma_\ell \|\Delta\|_2^2}_{\text{Lower curvature}} - \underbrace{\tau_\ell^2 \mathcal{R}^2(\Delta)}_{\text{Tolerance}}$$

for all  $\Delta$  in a suitable neighborhood of  $\theta^*$ .

- ordinary strong convexity:
  - ▶ special case with tolerance  $\tau_\ell = 0$
  - ▶ does not hold for most loss functions when  $d > n$

# Restricted strong convexity

## Definition

Loss function  $\mathcal{L}_n$  satisfies restricted strong convexity (RSC) with respect to regularizer  $\mathcal{R}$  if

$$\underbrace{\mathcal{L}_n(\theta^* + \Delta) - \left\{ \mathcal{L}_n(\theta^*) + \langle \nabla \mathcal{L}_n(\theta^*), \Delta \rangle \right\}}_{\text{Taylor error } \mathcal{E}_n(\Delta)} \geq \underbrace{\gamma_\ell \|\Delta\|_2^2}_{\text{Lower curvature}} - \underbrace{\tau_\ell^2 \mathcal{R}^2(\Delta)}_{\text{Tolerance}}$$

for all  $\Delta$  in a suitable neighborhood of  $\theta^*$ .

- ordinary strong convexity:
  - ▶ special case with tolerance  $\tau_\ell = 0$
  - ▶ does not hold for most loss functions when  $d > n$
- RSC enforces a lower bound on curvature, but **only** when  $\mathcal{R}^2(\Delta) \ll \|\Delta\|_2^2$

# Restricted strong convexity

## Definition

Loss function  $\mathcal{L}_n$  satisfies restricted strong convexity (RSC) with respect to regularizer  $\mathcal{R}$  if

$$\underbrace{\mathcal{L}_n(\theta^* + \Delta) - \left\{ \mathcal{L}_n(\theta^*) + \langle \nabla \mathcal{L}_n(\theta^*), \Delta \rangle \right\}}_{\text{Taylor error } \mathcal{E}_n(\Delta)} \geq \underbrace{\gamma_\ell \|\Delta\|_2^2}_{\text{Lower curvature}} - \underbrace{\tau_\ell^2 \mathcal{R}^2(\Delta)}_{\text{Tolerance}}$$

for all  $\Delta$  in a suitable neighborhood of  $\theta^*$ .

- ordinary strong convexity:
  - ▶ special case with tolerance  $\tau_\ell = 0$
  - ▶ does not hold for most loss functions when  $d > n$
- RSC enforces a lower bound on curvature, but **only** when  $\mathcal{R}^2(\Delta) \ll \|\Delta\|_2^2$
- a function satisfying RSC can actually be **non-convex**

## Example: RSC $\equiv$ RE for least-squares

- for least-squares loss  $\mathcal{L}(\theta) = \frac{1}{2n} \|y - X\theta\|_2^2$ :

$$\mathcal{E}_n(\Delta) = \mathcal{L}_n(\theta^* + \Delta) - \left\{ \mathcal{L}_n(\theta^*) - \langle \nabla \mathcal{L}_n(\theta^*), \Delta \rangle \right\} = \frac{1}{2n} \|X\Delta\|_2^2.$$

## Example: RSC $\equiv$ RE for least-squares

- for least-squares loss  $\mathcal{L}(\theta) = \frac{1}{2n} \|y - X\theta\|_2^2$ :

$$\mathcal{E}_n(\Delta) = \mathcal{L}_n(\theta^* + \Delta) - \left\{ \mathcal{L}_n(\theta^*) - \langle \nabla \mathcal{L}_n(\theta^*), \Delta \rangle \right\} = \frac{1}{2n} \|X\Delta\|_2^2.$$

- Restricted eigenvalue (RE) condition (van de Geer, 2007; Bickel et al., 2009):

$$\frac{\|X\Delta\|_2^2}{2n} \geq \gamma \|\Delta\|_2^2 \quad \text{for all } \|\Delta_{S^c}\|_1 \leq \|\Delta_S\|_1.$$

## Example: RSC $\equiv$ RE for least-squares

- for least-squares loss  $\mathcal{L}(\theta) = \frac{1}{2n} \|y - X\theta\|_2^2$ :

$$\mathcal{E}_n(\Delta) = \mathcal{L}_n(\theta^* + \Delta) - \left\{ \mathcal{L}_n(\theta^*) - \langle \nabla \mathcal{L}_n(\theta^*), \Delta \rangle \right\} = \frac{1}{2n} \|X\Delta\|_2^2.$$

- Restricted eigenvalue (RE) condition (van de Geer, 2007; Bickel et al., 2009):

$$\frac{\|X\Delta\|_2^2}{2n} \geq \gamma \|\Delta\|_2^2 \quad \text{for all } \Delta \in \mathbb{R}^d \text{ with } \|\Delta\|_1 \leq 2\sqrt{s}\|\Delta\|_2.$$

## Example: Generalized linear models

A broad class of models for relationship between response  $y \in \mathcal{X}$  and predictors  $x \in \mathbb{R}^d$ .

## Example: Generalized linear models

A broad class of models for relationship between response  $y \in \mathcal{X}$  and predictors  $x \in \mathbb{R}^d$ .

Based on families of conditional distributions:

$$\mathbb{P}_\theta(y \mid x, \theta^*) \propto \exp \left\{ \frac{y \langle x, \theta^* \rangle - \Phi(\langle x, \theta^* \rangle)}{c(\sigma)} \right\}.$$

## Example: Generalized linear models

A broad class of models for relationship between response  $y \in \mathcal{X}$  and predictors  $x \in \mathbb{R}^d$ .

Based on families of conditional distributions:

$$\mathbb{P}_\theta(y \mid x, \theta^*) \propto \exp \left\{ \frac{y \langle x, \theta^* \rangle - \Phi(\langle x, \theta^* \rangle)}{c(\sigma)} \right\}.$$

### Examples:

- Linear Gaussian model:  $\Phi(t) = t^2/2$  and  $c(\sigma) = \sigma^2$ .
- Binary response data  $y \in \{0, 1\}$ , Bernoulli model:  $\Phi(t) = \log(1 + e^t)$ .
- Multinomial responses (e.g., ratings)
- Poisson models (count-valued data):  $\Phi(t) = e^t$ .

## GLM-based restricted strong convexity

- let  $\mathcal{R}$  be norm-based regularizer dominating the  $\ell_2$ -norm (e.g.,  $\ell_1$ , group-sparse, nuclear etc.)
- let  $\mathcal{R}^*$  be the associated dual norm
- covariate-Rademacher complexity of norm ball

$$\sup_{\mathcal{R}(u) \leq 1} \left\langle u, \frac{1}{n} \sum_{i=1}^n \varepsilon_i x_i \right\rangle = \mathcal{R}^* \left( \frac{1}{n} \sum_{i=1}^n \varepsilon_i x_i \right)$$

where  $\{\varepsilon_i\}_{i=1}^n$  are i.i.d sign variables

# GLM-based restricted strong convexity

- let  $\mathcal{R}$  be norm-based regularizer dominating the  $\ell_2$ -norm (e.g.,  $\ell_1$ , group-sparse, nuclear etc.)
- let  $\mathcal{R}^*$  be the associated dual norm
- covariate-Rademacher complexity of norm ball

$$\sup_{\mathcal{R}(u) \leq 1} \left\langle u, \frac{1}{n} \sum_{i=1}^n \varepsilon_i x_i \right\rangle = \mathcal{R}^* \left( \frac{1}{n} \sum_{i=1}^n \varepsilon_i x_i \right)$$

where  $\{\varepsilon_i\}_{i=1}^n$  are i.i.d sign variables

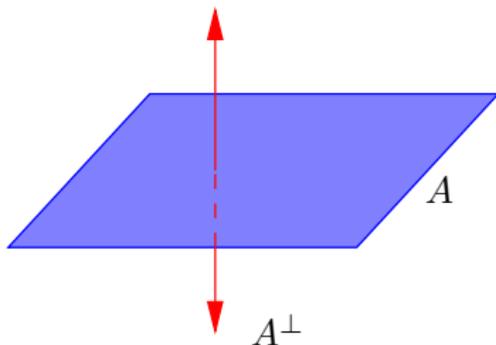
## Theorem (Negahban et al., 2012; W. 2014)

Let the covariates  $\{x_i\}_{i=1}^n$  be sampled i.i.d. Then

$$\underbrace{\mathcal{E}_n(\Delta)}_{\text{Emp. Taylor error}} \geq \underbrace{\bar{\mathcal{E}}(\Delta)}_{\text{Pop. Taylor error}} - c_1 \{t \mathcal{R}(\Delta)\}^2 \quad \text{for all } \|\Delta\|_2 \leq 1$$

with probability at least  $1 - \mathbb{P}[\mathcal{R}^* \left( \frac{1}{n} \sum_{i=1}^n \varepsilon_i x_i \right) \geq t]$ .

## (II) Decomposable regularizers



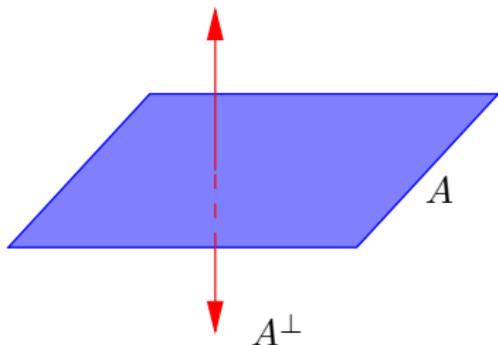
Subspace  $A$ :

Approximation to model parameters

Complementary subspace  $A^\perp$ :

Undesirable deviations.

## (II) Decomposable regularizers

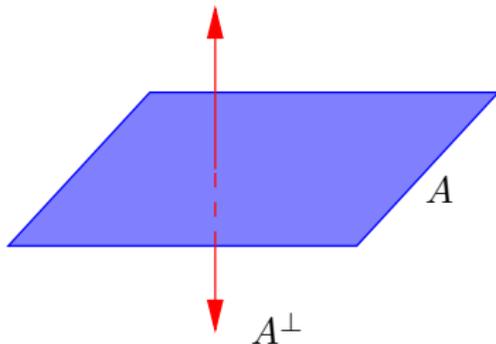


Subspace  $A$ : Approximation to model parameters  
Complementary subspace  $A^\perp$ : Undesirable deviations.

Regularizer  $\mathcal{R}$  decomposes across  $(A, A^\perp)$  if

$$\mathcal{R}(\alpha + \beta) = \mathcal{R}(\alpha) + \mathcal{R}(\beta) \quad \text{for all } \alpha \in A, \text{ and } \beta \in A^\perp.$$

## (II) Decomposable regularizers

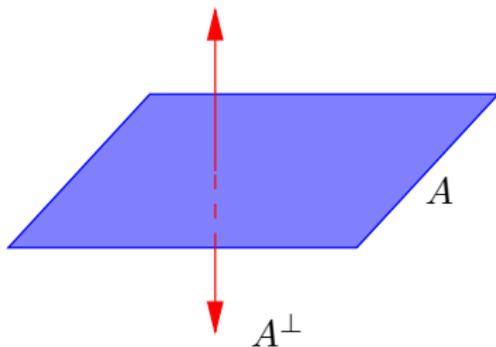


Regularizer  $\mathcal{R}$  decomposes across  $(A, A^\perp)$  if

$$\mathcal{R}(\alpha + \beta) = \mathcal{R}(\alpha) + \mathcal{R}(\beta) \quad \text{for all } \alpha \in A, \text{ and } \beta \in A^\perp.$$

- Includes:
- (weighted)  $\ell_1$ -norms
  - group-sparse norms
  - nuclear norm
  - sums of decomposable norms

## (II) Decomposable regularizers



Regularizer  $\mathcal{R}$  decomposes across  $(A, A^\perp)$  if

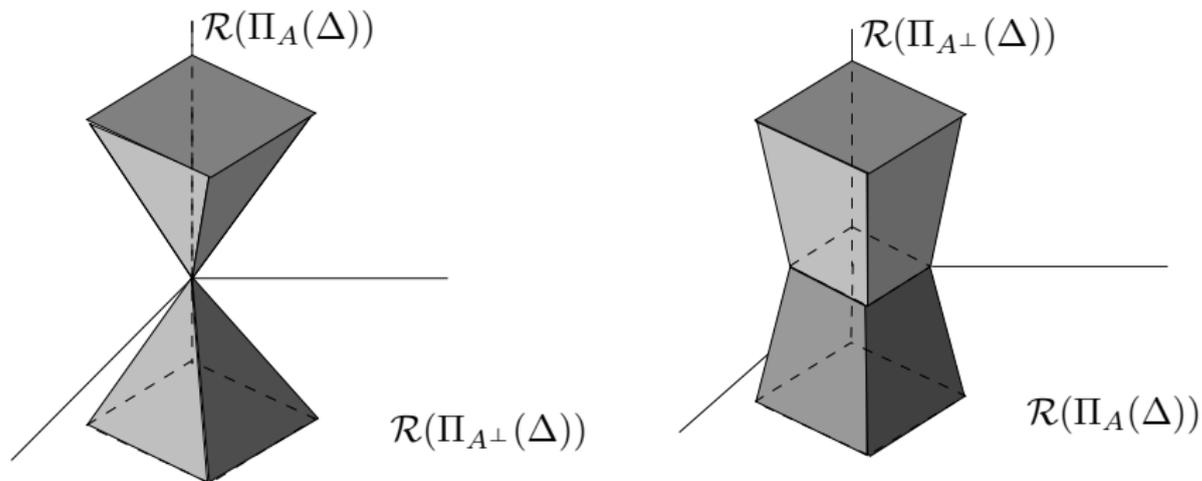
$$\mathcal{R}(\alpha + \beta) = \mathcal{R}(\alpha) + \mathcal{R}(\beta) \quad \text{for all } \alpha \in A, \text{ and } \beta \in A^\perp.$$

### Related definitions:

Geometric decomposability: Candes & Recht, 2012; Chandrasekaran et al., 2012

Weak decomposability: van de Geer, 2012

# Significance of decomposability



(a)  $\mathbb{C}$  for exact model (cone)

(b)  $\mathbb{C}$  for approximate model (star-shaped)

## Lemma

Suppose that  $\mathcal{L}$  is convex, and  $\mathcal{R}$  is decomposable w.r.t.  $A$ . Then as long as  $\lambda_n \geq 2\mathcal{R}^*(\nabla\mathcal{L}(\theta^*; Z_1^n))$ , the error vector  $\widehat{\Delta} = \widehat{\theta}_{\lambda_n} - \theta^*$  belongs to

$$\mathbb{C}(A, B; \theta^*) := \{\Delta \in \Omega \mid \mathcal{R}(\Pi_{A^\perp}(\Delta)) \leq 3\mathcal{R}(\Pi_B(\Delta)) + 4\mathcal{R}(\Pi_{A^\perp}(\theta^*))\}.$$

## Example: Sparse vectors and $\ell_1$ -regularization

- for each subset  $S \subset \{1, \dots, d\}$ , define subspace pairs

$$\begin{aligned}A(S) &:= \{\theta \in \mathbb{R}^d \mid \theta_{S^c} = 0\}, \\B^\perp(S) &:= \{\theta \in \mathbb{R}^d \mid \theta_S = 0\} = A^\perp(S).\end{aligned}$$

- decomposability of  $\ell_1$ -norm:

$$\|\theta_S + \theta_{S^c}\|_1 = \|\theta_S\|_1 + \|\theta_{S^c}\|_1 \quad \text{for all } \theta_S \in A(S) \text{ and } \theta_{S^c} \in B^\perp(S).$$

- natural extension to group Lasso:

- ▶ collection of groups  $\mathcal{G}_j$  that partition  $\{1, \dots, d\}$
- ▶ group norm

$$\|\theta\|_{\mathcal{G}, \alpha} = \sum_j \|\theta_{\mathcal{G}_j}\|_\alpha \quad \text{for some } \alpha \in [1, \infty].$$

## Example: Low-rank matrices and nuclear norm

- for each pair of  $r$ -dimensional subspaces  $U \subseteq \mathbb{R}^{p_1}$  and  $V \subseteq \mathbb{R}^{p_2}$ :

$$A(U, V) := \{\Theta \in \mathbb{R}^{p_1 \times p_2} \mid \text{row}(\Theta) \subseteq V, \text{col}(\Theta) \subseteq U\}$$

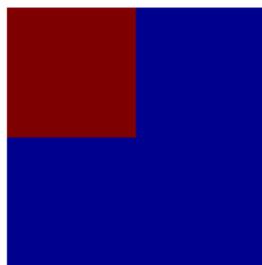
$$B^\perp(U, V) := \{\Gamma \in \mathbb{R}^{p_1 \times p_2} \mid \text{row}(\Gamma) \subseteq V^\perp, \text{col}(\Gamma) \subseteq U^\perp\}.$$

## Example: Low-rank matrices and nuclear norm

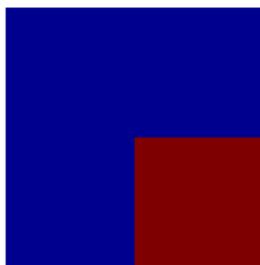
- for each pair of  $r$ -dimensional subspaces  $U \subseteq \mathbb{R}^{p_1}$  and  $V \subseteq \mathbb{R}^{p_2}$ :

$$A(U, V) := \{\Theta \in \mathbb{R}^{p_1 \times p_2} \mid \text{row}(\Theta) \subseteq V, \text{col}(\Theta) \subseteq U\}$$

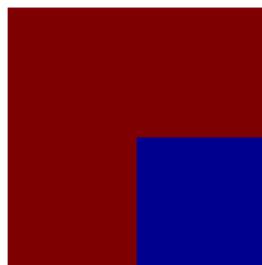
$$B^\perp(U, V) := \{\Gamma \in \mathbb{R}^{p_1 \times p_2} \mid \text{row}(\Gamma) \subseteq V^\perp, \text{col}(\Gamma) \subseteq U^\perp\}.$$



(a)  $\Theta \in A$



(b)  $\Gamma \in B^\perp$

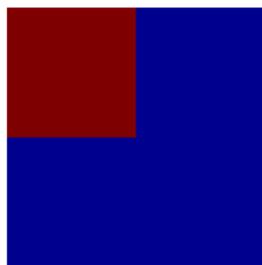


(c)  $\Sigma \in B$

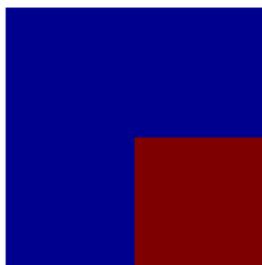
## Example: Low-rank matrices and nuclear norm

- for each pair of  $r$ -dimensional subspaces  $U \subseteq \mathbb{R}^{p_1}$  and  $V \subseteq \mathbb{R}^{p_2}$ :

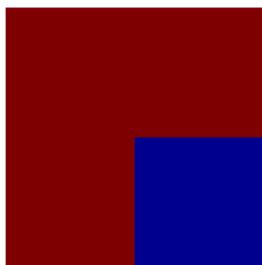
$$A(U, V) := \{\Theta \in \mathbb{R}^{p_1 \times p_2} \mid \text{row}(\Theta) \subseteq V, \text{col}(\Theta) \subseteq U\}$$
$$B^\perp(U, V) := \{\Gamma \in \mathbb{R}^{p_1 \times p_2} \mid \text{row}(\Gamma) \subseteq V^\perp, \text{col}(\Gamma) \subseteq U^\perp\}.$$



(a)  $\Theta \in A$



(b)  $\Gamma \in B^\perp$



(c)  $\Sigma \in B$

- by construction,  $\Theta^T \Gamma = 0$  for all  $\Theta \in A(U, V)$  and  $\Gamma \in B^\perp(U, V)$
- decomposability of nuclear norm  $\|\Theta\|_1 = \sum_{j=1}^{\min\{p_1, p_2\}} \sigma_j(\Theta)$ :

$$\|\Theta + \Gamma\|_1 = \|\Theta\|_1 + \|\Gamma\|_1 \quad \text{for all } \Theta \in A(U, V) \text{ and } \Gamma \in B^\perp(U, V).$$

# Main theorem

Estimator

$$\hat{\theta}_{\lambda_n} \in \arg \min_{\theta \in \mathbb{R}^d} \{ \mathcal{L}_n(\theta; Z_1^n) + \lambda_n \mathcal{R}(\theta) \},$$

where  $\mathcal{L}$  satisfies  $\text{RSC}(\gamma, \tau)$  w.r.t regularizer  $\mathcal{R}$ .

# Main theorem

Estimator

$$\hat{\theta}_{\lambda_n} \in \arg \min_{\theta \in \mathbb{R}^d} \{ \mathcal{L}_n(\theta; Z_1^n) + \lambda_n \mathcal{R}(\theta) \},$$

where  $\mathcal{L}$  satisfies  $\text{RSC}(\gamma, \tau)$  w.r.t regularizer  $\mathcal{R}$ .

## Theorem (Negahban, Ravikumar, W., & Yu, 2012)

Suppose that  $\theta^* \in A$ , and  $\Psi^2(A)\tau_n^2 < 1$ . Then for any regularization parameter  $\lambda_n \geq 2\mathcal{R}^*(\nabla \mathcal{L}(\theta^*; Z_1^n))$ , any solution  $\hat{\theta}_{\lambda_n}$  satisfies

$$\|\hat{\theta}_{\lambda_n} - \theta^*\|_{\star}^2 \lesssim \frac{1}{\gamma^2(\mathcal{L})} \lambda_n^2 \Psi^2(A).$$

Quantities that control rates:

- curvature in RSC:  $\gamma_{\ell}$
- tolerance in RSC:  $\tau$
- dual norm of regularizer:  $\mathcal{R}^*(v) := \sup_{\mathcal{R}(u) \leq 1} \langle v, u \rangle$ .
- optimal subspace const.:  $\Psi(A) = \sup_{\theta \in A \setminus \{0\}} \mathcal{R}(\theta) / \|\theta\|_{\star}$

# Main theorem

Estimator

$$\hat{\theta}_{\lambda_n} \in \arg \min_{\theta \in \mathbb{R}^d} \{ \mathcal{L}_n(\theta; Z_1^n) + \lambda_n \mathcal{R}(\theta) \},$$

## Theorem (Oracle version)

With  $\lambda_n \geq 2\mathcal{R}^*(\nabla \mathcal{L}(\theta^*; Z_1^n))$ , any solution  $\hat{\theta}$  satisfies

$$\|\hat{\theta}_{\lambda_n} - \theta^*\|_{\star}^2 \lesssim \underbrace{\frac{(\lambda'_n)^2}{\gamma^2} \Psi^2(A)}_{\text{Estimation error}} + \underbrace{\frac{\lambda'_n}{\gamma} \mathcal{R}(\Pi_{A^\perp}(\theta^*))}_{\text{Approximation error}}$$

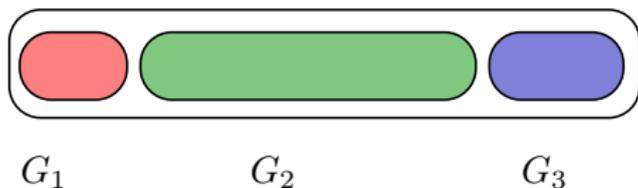
where  $\lambda' = \max\{\lambda, \tau\}$ .

Quantities that control rates:

- curvature in RSC:  $\gamma_\ell$
- tolerance in RSC:  $\tau$
- dual norm of regularizer:  $\mathcal{R}^*(v) := \sup_{\mathcal{R}(u) \leq 1} \langle v, u \rangle$ .
- optimal subspace const.:  $\Psi(A) = \sup_{\theta \in A \setminus \{0\}} \mathcal{R}(\theta) / \|\theta\|_{\star}$

## Example: Group-structured regularizers

Many applications exhibit sparsity with more structure.....



- divide index set  $\{1, 2, \dots, d\}$  into groups  $\mathcal{G} = \{G_1, G_2, \dots, G_T\}$
- for parameters  $\nu_i \in [1, \infty]$ , define block-norm

$$\|\theta\|_{\nu, \mathcal{G}} := \sum_{t=1}^T \|\theta_{G_t}\|_{\nu_t}$$

- group/block Lasso program

$$\hat{\theta}_{\lambda_n} \in \arg \min_{\theta \in \mathbb{R}^d} \left\{ \frac{1}{2n} \|y - X\theta\|_2^2 + \lambda_n \|\theta\|_{\nu, \mathcal{G}} \right\}.$$

- different versions studied by various authors

(Wright et al., 2005; Tropp et al., 2006; Yuan & Li, 2006; Baraniuk, 2008; Obozinski et al., 2008; Zhao et al., 2008; Bach et al., 2009; Lounici et al., 2009)

# Convergence rates for general group Lasso

## Corollary

Say  $\Theta^*$  is supported on group subset  $S_G$ , and  $X$  satisfies *RSC*. Then for regularization parameter

$$\lambda_n \geq 2 \max_{t=1,2,\dots,T} \left\| \frac{X^T w}{n} \right\|_{\nu_t^*}, \quad \text{where } \frac{1}{\nu_t^*} = 1 - \frac{1}{\nu_t},$$

any solution  $\hat{\theta}_{\lambda_n}$  satisfies

$$\|\hat{\theta}_{\lambda_n} - \theta^*\|_2 \leq \frac{2}{\gamma(\mathcal{L})} \Psi_\nu(S_G) \lambda_n, \quad \text{where } \Psi_\nu(S_G) = \sup_{\theta \in A(S_G) \setminus \{0\}} \frac{\|\theta\|_{\nu, S_G}}{\|\theta\|_2}.$$

# Convergence rates for general group Lasso

## Corollary

Say  $\Theta^*$  is supported on group subset  $\mathcal{S}_G$ , and  $X$  satisfies *RSC*. Then for regularization parameter

$$\lambda_n \geq 2 \max_{t=1,2,\dots,T} \left\| \frac{X^T w}{n} \right\|_{\nu_t^*}, \quad \text{where } \frac{1}{\nu_t^*} = 1 - \frac{1}{\nu_t},$$

any solution  $\hat{\theta}_{\lambda_n}$  satisfies

$$\|\hat{\theta}_{\lambda_n} - \theta^*\|_2 \leq \frac{2}{\gamma(\mathcal{L})} \Psi_\nu(\mathcal{S}_G) \lambda_n, \quad \text{where } \Psi_\nu(\mathcal{S}_G) = \sup_{\theta \in A(\mathcal{S}_G) \setminus \{0\}} \frac{\|\theta\|_{\nu, \mathcal{S}_G}}{\|\theta\|_2}.$$

Some special cases with  $m \equiv \max.$  group size

- 1  $\ell_1/\ell_2$  regularization: Group norm with  $\nu = 2$

$$\|\hat{\theta}_{\lambda_n} - \theta^*\|_2^2 = \mathcal{O}\left(\frac{|\mathcal{S}_G| m}{n} + \frac{|\mathcal{S}_G| \log T}{n}\right).$$

# Convergence rates for general group Lasso

## Corollary

Say  $\Theta^*$  is supported on group subset  $\mathcal{S}_G$ , and  $X$  satisfies *RSC*. Then for regularization parameter

$$\lambda_n \geq 2 \max_{t=1,2,\dots,T} \left\| \frac{X^T w}{n} \right\|_{\nu_t^*}, \quad \text{where } \frac{1}{\nu_t^*} = 1 - \frac{1}{\nu_t},$$

any solution  $\hat{\theta}_{\lambda_n}$  satisfies

$$\|\hat{\theta}_{\lambda_n} - \theta^*\|_2 \leq \frac{2}{\gamma(\mathcal{L})} \Psi_\nu(\mathcal{S}_G) \lambda_n, \quad \text{where } \Psi_\nu(\mathcal{S}_G) = \sup_{\theta \in A(\mathcal{S}_G) \setminus \{0\}} \frac{\|\theta\|_{\nu, \mathcal{S}_G}}{\|\theta\|_2}.$$

Some special cases with  $m \equiv \max.$  group size

①  $\ell_1/\ell_\infty$  regularization: group norm with  $\nu = \infty$

$$\|\hat{\theta}_{\lambda_n} - \theta^*\|_2^2 = \mathcal{O}\left(\frac{|\mathcal{S}_G| m^2}{n} + \frac{|\mathcal{S}_G| \log T}{n}\right).$$

## Is adaptive estimation possible?

Consider a group-sparse problem with:

- $T$  groups in total
- each of size  $m$
- $|\mathcal{S}_G|$ -active groups
- $T$  active coefficients per group

Group Lasso will achieve

$$\|\hat{\theta} - \theta^*\|_2^2 \lesssim \frac{|\mathcal{S}_G|m}{n} + \frac{|\mathcal{S}_G| \log |\mathcal{G}|}{n}.$$

Lasso will achieve

$$\|\hat{\theta} - \theta^*\|_2^2 \lesssim \frac{|\mathcal{S}_G| T \log(|\mathcal{G}|m)}{n}.$$

# Is adaptive estimation possible?

Consider a group-sparse problem with:

- $T$  groups in total
- each of size  $m$
- $|\mathcal{S}_G|$ -active groups
- $T$  active coefficients per group

Group Lasso will achieve

$$\|\hat{\theta} - \theta^*\|_2^2 \lesssim \frac{|\mathcal{S}_G|m}{n} + \frac{|\mathcal{S}_G| \log |\mathcal{G}|}{n}.$$

Lasso will achieve

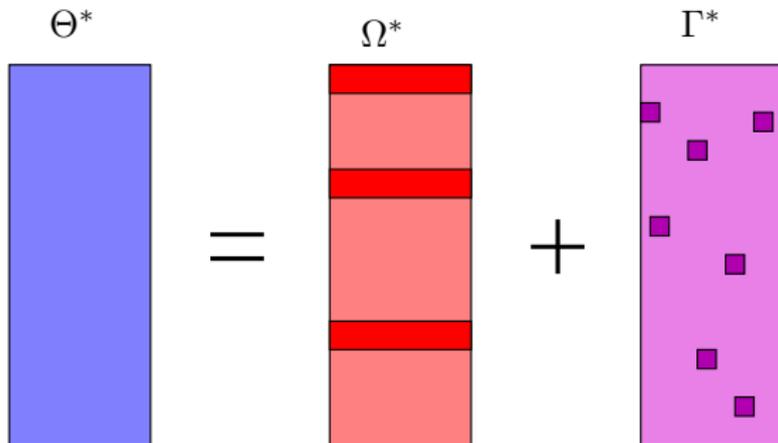
$$\|\hat{\theta} - \theta^*\|_2^2 \lesssim \frac{|\mathcal{S}_G| T \log(|\mathcal{G}|m)}{n}.$$

## Natural question:

Can we design an estimator that optimally adapts to the degree of elementwise versus group sparsity?

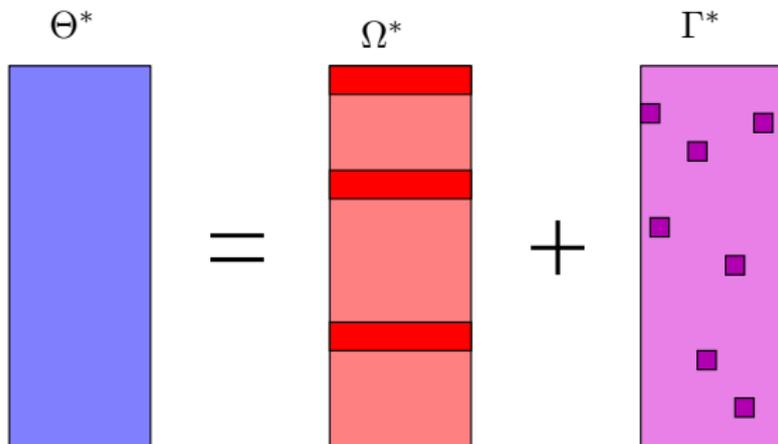
# Answer: Overlap group Lasso

Represent  $\Theta^*$  as a sum of **row-sparse** and **element-wise sparse** matrices.



# Answer: Overlap group Lasso

Represent  $\Theta^*$  as a sum of **row-sparse** and **element-wise sparse** matrices.

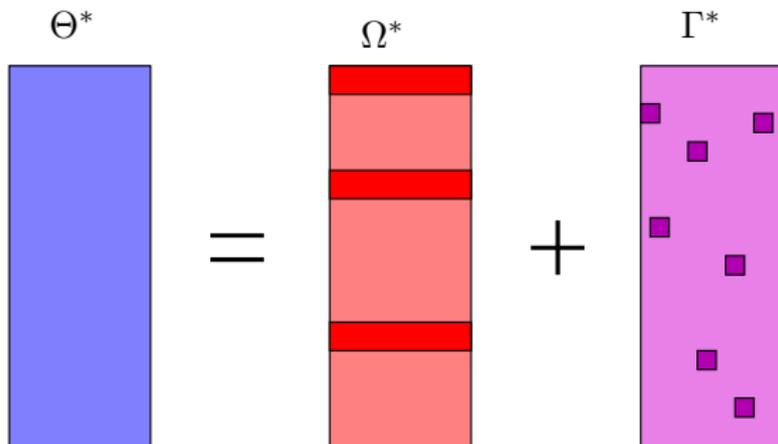


Define new norm on matrix space:

$$\mathcal{R}_\omega(\Theta) = \inf_{\Theta = \Omega + \Gamma} \left\{ \omega \|\Omega\|_{1,2} + \|\Gamma\|_1 \right\}.$$

# Answer: Overlap group Lasso

Represent  $\Theta^*$  as a sum of **row-sparse** and **element-wise sparse** matrices.



Define new norm on matrix space:

$$\mathcal{R}_\omega(\Theta) = \inf_{\Theta = \Omega + \Gamma} \left\{ \omega \|\Omega\|_{1,2} + \|\Gamma\|_1 \right\}.$$

Special case of the overlap group Lasso: (Obozinski et al., 2008; Jalali et al., 2011)

## Example: Adaptivity with overlap group Lasso

Consider regularizer

$$\mathcal{R}_\omega(\Theta) = \inf_{\Theta = \Omega + \Gamma} \left\{ \omega \|\Omega\|_{1,2} + \|\Gamma\|_1 \right\}.$$

with

$$\omega = \frac{\sqrt{m} + \sqrt{\log |\mathcal{G}|}}{\sqrt{\log d}},$$

- $|\mathcal{G}|$  is number of groups
- $m$  is max. group size
- $d$  is number of predictors.

## Example: Adaptivity with overlap group Lasso

Consider regularizer

$$\mathcal{R}_\omega(\Theta) = \inf_{\Theta = \Omega + \Gamma} \left\{ \omega \|\Omega\|_{1,2} + \|\Gamma\|_1 \right\}.$$

with

$$\omega = \frac{\sqrt{m} + \sqrt{\log |\mathcal{G}|}}{\sqrt{\log d}},$$

- $|\mathcal{G}|$  is number of groups
- $m$  is max. group size
- $d$  is number of predictors.

### Corollary

*Under RSC condition on loss function, suppose that  $\Theta^*$  can be decomposed as a sum of an  $|S_{elt}|$ -elementwise sparse matrix and an  $|S_G|$ -groupwise sparse matrix (disjointly). Then for  $\lambda = 4\sigma\sqrt{\frac{\log d}{n}}$ , any optimal solution satisfies (w.h.p.)*

$$\|\hat{\Theta} - \Theta^*\|_F^2 \lesssim \sigma^2 \left\{ \frac{|S_G|m}{n} + \frac{|S_G|\log |\mathcal{G}|}{n} \right\} + \sigma^2 \left\{ \frac{|S_{elt}|\log d}{n} \right\}.$$

## Example: Low-rank matrices and nuclear norm

- low-rank matrix  $\Theta^* \in \mathbb{R}^{p_1 \times p_2}$  that is exactly (or approximately) low-rank
- noisy/partial observations of the form

$$y_i = \langle X_i, \Theta^* \rangle + w_i, \quad i = 1, \dots, n, \quad w_i \text{ i.i.d. noise}$$

- estimate by solving semi-definite program (SDP):

$$\hat{\Theta} \in \arg \min_{\Theta} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - \langle X_i, \Theta \rangle)^2 + \lambda_n \underbrace{\sum_{j=1}^{\min\{p_1, p_2\}} \gamma_j(\Theta)}_{\|\Theta\|_1} \right\}$$

## Example: Low-rank matrices and nuclear norm

- low-rank matrix  $\Theta^* \in \mathbb{R}^{p_1 \times p_2}$  that is exactly (or approximately) low-rank
- noisy/partial observations of the form

$$y_i = \langle X_i, \Theta^* \rangle + w_i, \quad i = 1, \dots, n, \quad w_i \text{ i.i.d. noise}$$

- estimate by solving semi-definite program (SDP):

$$\hat{\Theta} \in \arg \min_{\Theta} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - \langle X_i, \Theta \rangle)^2 + \lambda_n \underbrace{\sum_{j=1}^{\min\{p_1, p_2\}} \gamma_j(\Theta)}_{\|\Theta\|_1} \right\}$$

- various applications:
  - ▶ matrix compressed sensing
  - ▶ matrix completion
  - ▶ rank-reduced multivariate regression (multi-task learning)
  - ▶ time-series modeling (vector autoregressions)
  - ▶ phase-retrieval problems

## Rates for (near) low-rank estimation

For simplicity, consider matrix compressed sensing model:  $X_i$  are random sub-Gaussian projections).

For parameter  $q \in [0, 1]$ , set of near low-rank matrices:

$$\mathbb{B}_q(R_q) = \left\{ \Theta^* \in \mathbb{R}^{p_1 \times p_2} \mid \sum_{j=1}^{\min\{p_1, p_2\}} |\sigma_j(\Theta^*)|^q \leq R_q \right\}.$$

## Rates for (near) low-rank estimation

For simplicity, consider matrix compressed sensing model:  $X_i$  are random sub-Gaussian projections).

For parameter  $q \in [0, 1]$ , set of near low-rank matrices:

$$\mathbb{B}_q(R_q) = \left\{ \Theta^* \in \mathbb{R}^{p_1 \times p_2} \mid \sum_{j=1}^{\min\{p_1, p_2\}} |\sigma_j(\Theta^*)|^q \leq R_q \right\}.$$

### Corollary (Negahban & W., 2011)

With regularization parameter  $\lambda_n \geq 16\sigma \left( \sqrt{\frac{p_1}{n}} + \sqrt{\frac{p_2}{n}} \right)$ , we have w.h.p.

$$\|\hat{\Theta} - \Theta^*\|_F^2 \leq c_0 \frac{R_q}{\gamma(\mathcal{L})^2} \left( \frac{\sigma^2 (p_1 + p_2)}{n} \right)^{1 - \frac{q}{2}}$$

# Rates for (near) low-rank estimation

For parameter  $q \in [0, 1]$ , set of near low-rank matrices:

$$\mathbb{B}_q(R_q) = \left\{ \Theta^* \in \mathbb{R}^{p_1 \times p_2} \mid \sum_{j=1}^{\min\{p_1, p_2\}} |\sigma_j(\Theta^*)|^q \leq R_q \right\}.$$

## Corollary (Negahban & W., 2011)

With regularization parameter  $\lambda_n \geq 16\sigma \left( \sqrt{\frac{p_1}{n}} + \sqrt{\frac{p_2}{n}} \right)$ , we have w.h.p.

$$\|\hat{\Theta} - \Theta^*\|_F^2 \leq c_0 \frac{R_q}{\gamma(\mathcal{L})^2} \left( \frac{\sigma^2(p_1 + p_2)}{n} \right)^{1 - \frac{q}{2}}$$

- for a rank  $r$  matrix  $M$

$$\|M\|_1 = \sum_{j=1}^r \sigma_j(M) \leq \sqrt{r} \sqrt{\sum_{j=1}^r \sigma_j^2(M)} = \sqrt{r} \|M\|_F$$

- solve nuclear norm regularized program with  $\lambda_n \geq \frac{2}{n} \left\| \sum_{i=1}^n w_i X_i \right\|_2$

# Matrix completion

Random operator  $\mathfrak{X} : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^n$  with

$$[\mathfrak{X}(\Theta^*)]_i = d \Theta_{a(i)b(i)}^*$$

where  $(a(i), b(i))$  is a matrix index sampled uniformly at random.

# Matrix completion

Random operator  $\mathfrak{X} : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^n$  with

$$[\mathfrak{X}(\Theta^*)]_i = d \Theta_{a(i)b(i)}^*$$

where  $(a(i), b(i))$  is a matrix index sampled uniformly at random.

Even in noiseless setting, model is **unidentifiable**:

Consider a rank one matrix:

$$\Theta^* = e_1 e_1^T = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & 0 \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix}$$

# Matrix completion

Random operator  $\mathfrak{X} : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^n$  with

$$[\mathfrak{X}(\Theta^*)]_i = d \Theta_{a(i)b(i)}^*$$

where  $(a(i), b(i))$  is a matrix index sampled uniformly at random.

Even in noiseless setting, model is **unidentifiable**:

Consider a rank one matrix:

$$\Theta^* = e_1 e_1^T = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & 0 \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix}$$

Exact recovery based on **eigen-incoherence** involving leverage scores (e.g., Recht & Candes, 2008; Gross, 2009)

## A milder “spikiness” condition

Consider the “poisoned” low-rank matrix:

$$\Theta^* = \Gamma^* + \delta e_1 e_1^T = \Gamma^* + \delta \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & 0 \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix}$$

where  $\Gamma^*$  is rank  $r - 1$ , all eigenvectors perpendicular to  $e_1$ .

Excluded by eigen-incoherence for all  $\delta > 0$ .

## A milder “spikiness” condition

Consider the “poisoned” low-rank matrix:

$$\Theta^* = \Gamma^* + \delta e_1 e_1^T = \Gamma^* + \delta \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & 0 \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix}$$

where  $\Gamma^*$  is rank  $r - 1$ , all eigenvectors perpendicular to  $e_1$ .

Excluded by eigen-incoherence for all  $\delta > 0$ .

Control by **spikiness ratio**:

$$1 \leq \frac{d \|\Theta^*\|_\infty}{\|\Theta^*\|_F} \leq d.$$

Spikiness constraints used in various papers: Oh et al., 2009; Negahban & W. 2010, Koltchinski et al., 2011.

# Uniform law for matrix completion

Let  $\mathfrak{X}_n : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^n$  be **rescaled** matrix completion random operator

$(\mathfrak{X}_n(\Theta))_i \mapsto d \Theta_{a(i), b(i)}$  where index  $(a(i), b(i))$  from uniform distribution.

Define family of zero-mean random variables:

$$Z_n(\Theta) := \frac{\|\mathfrak{X}_n(\Theta)\|_2^2}{n} - \|\Theta\|_F^2, \quad \text{for } \Theta \in \mathbb{R}^{d \times d}.$$

# Uniform law for matrix completion

Let  $\mathfrak{X}_n : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^n$  be **rescaled** matrix completion random operator

$(\mathfrak{X}_n(\Theta))_i \mapsto d \Theta_{a(i), b(i)}$  where index  $(a(i), b(i))$  from uniform distribution.

Define family of zero-mean random variables:

$$Z_n(\Theta) := \frac{\|\mathfrak{X}_n(\Theta)\|_2^2}{n} - \|\Theta\|_F^2, \quad \text{for } \Theta \in \mathbb{R}^{d \times d}.$$

## Theorem (Negahban & W., 2010)

For random matrix completion operator  $\mathfrak{X}_n$ , there are universal positive constants  $(c_1, c_2)$  such that

$$\sup_{\Theta \in \mathbb{R}^{d \times d} \setminus \{0\}} Z_n(\Theta) \leq \underbrace{c_1 d \|\Theta\|_\infty \|\Theta\|_{\text{nuc}} \sqrt{\frac{d \log d}{n}}}_{\text{"low-rank term"}} + \underbrace{c_2 \left( d \|\Theta\|_\infty \sqrt{\frac{d \log d}{n}} \right)^2}_{\text{"spikiness" term}}$$

with probability at least  $1 - \exp(-d \log d)$ .