

# Hypothesis Testing for High-Dimensional Regression: Nearly Optimal Sample Size

Adel Javanmard

Stanford University- UC Berkeley

Based on joint work with

Andrea Montanari

January 2015

# Outline

- 1 Problem definition
- 2 Debiasing approach
- 3 Hypothesis testing under nearly optimal sample size

## Problem definition

# Linear model

We focus on linear models:

$$Y = \mathbf{X}\theta_0 + W$$

- $Y \in \mathbb{R}^n$  (response),  $\mathbf{X} \in \mathbb{R}^{n \times p}$  (design matrix),  $\theta_0 \in \mathbb{R}^p$  (parameters)
- Noise vector has independent entries with

$$\mathbb{E}(W_i) = 0, \quad E(W_i^2) = \sigma^2,$$

$$\mathbb{E}(|W_i|^{2+\kappa}) < \infty, \text{ for some } \kappa > 0.$$

# Problem

- **Confidence intervals:** For each  $i \in \{1, \dots, p\}$ ,  $\underline{\theta}_i, \bar{\theta}_i \in \mathbb{R}$  such that

$$\mathbb{P}\left(\theta_{0,i} \in [\underline{\theta}_i, \bar{\theta}_i]\right) \geq 1 - \alpha$$

We would like  $|\underline{\theta}_i - \bar{\theta}_i|$  as small as possible.

- **Hypothesis testing:**

$$H_{0,i} : \theta_{0,i} = 0, \quad H_{A,i} : \theta_{0,i} \neq 0$$

# LASSO

$$\hat{\boldsymbol{\theta}} \equiv \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^p} \left\{ \frac{1}{2n} \|y - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \lambda \|\boldsymbol{\theta}\|_1 \right\}.$$

[Tibshirani 1996, Chen, Donoho 1996]

- Distribution of  $\hat{\boldsymbol{\theta}}$ ?

# LASSO

$$\hat{\boldsymbol{\theta}} \equiv \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^p} \left\{ \frac{1}{2n} \|y - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \lambda \|\boldsymbol{\theta}\|_1 \right\}.$$

[Tibshirani 1996, Chen, Donoho 1996]

- Distribution of  $\hat{\boldsymbol{\theta}}$ ?
- Debiasing approach:  
(LASSO is biased towards small  $\ell_1$  norm.)

# LASSO

$$\hat{\theta} \equiv \arg \min_{\theta \in \mathbb{R}^p} \left\{ \frac{1}{2n} \|y - \mathbf{X}\theta\|_2^2 + \lambda \|\theta\|_1 \right\}.$$

[Tibshirani 1996, Chen, Donoho 1996]

- Distribution of  $\hat{\theta}$ ?
- Debiasing approach:  
(LASSO is biased towards small  $\ell_1$  norm.)

$$\hat{\theta} \xrightarrow{\text{debiasing}} \hat{\theta}^d$$

We characterize distribution of  $\hat{\theta}^d$ .



## Debiasing approach

## Classical setting ( $n \gg p$ )

We know everything about the **least-square** estimator:

$$\hat{\theta}^{\text{LS}} = \frac{1}{n} \hat{\Sigma}^{-1} \mathbf{X}^T Y,$$

where  $\hat{\Sigma} \equiv (\mathbf{X}^T \mathbf{X})/n$  is empirical covariance.

## Classical setting ( $n \gg p$ )

We know everything about the **least-square** estimator:

$$\hat{\boldsymbol{\theta}}^{\text{LS}} = \frac{1}{n} \hat{\boldsymbol{\Sigma}}^{-1} \mathbf{X}^T Y,$$

where  $\hat{\boldsymbol{\Sigma}} \equiv (\mathbf{X}^T \mathbf{X})/n$  is empirical covariance.

- Confidence intervals:

$$[\underline{\theta}_i, \bar{\theta}_i] = [\hat{\theta}_i^{\text{LS}} - c_\alpha \Delta_i, \hat{\theta}_i^{\text{LS}} + c_\alpha \Delta_i], \quad \Delta_i \equiv \sigma \sqrt{\frac{(\hat{\boldsymbol{\Sigma}}^{-1})_{ii}}{n}}$$

## High-dimensional setting ( $n < p$ )

$$\hat{\theta}^{\text{LS}} = \frac{1}{n} \hat{\Sigma}^{-1} \mathbf{X}^T Y$$

**Problem in high dimension:**

$\hat{\Sigma}$  is not invertible!

## High-dimensional setting ( $n < p$ )

$$\hat{\theta}^{\text{LS}} = \frac{1}{n} \hat{\Sigma}^{-1} \mathbf{X}^T \mathbf{Y}$$

Take your favorite  $M \in \mathbb{R}^{p \times p}$ :

$$\begin{aligned} \hat{\theta}^* &= \frac{1}{n} M \mathbf{X}^T \mathbf{Y} \\ &= \frac{1}{n} M \mathbf{X}^T \mathbf{X} \theta_0 + \frac{1}{n} M \mathbf{X}^T \mathbf{W} \\ &= \theta_0 + \underbrace{(M \hat{\Sigma} - \mathbf{I}) \theta_0}_{\text{bias}} + \underbrace{\frac{1}{n} M \mathbf{X}^T \mathbf{W}}_{\text{Gaussian error}} \end{aligned}$$

## Debiased estimator

$$\hat{\theta}^* = \theta_0 + \underbrace{(M\hat{\Sigma} - \mathbf{I})\theta_0}_{\text{bias}} + \underbrace{\frac{1}{n}M\mathbf{X}^T W}_{\text{Gaussian error}}$$

## Debiased estimator

$$\hat{\theta}^* = \theta_0 + \underbrace{(M\hat{\Sigma} - \mathbf{I})\theta_0}_{\text{bias}} + \underbrace{\frac{1}{n}M\mathbf{X}^T W}_{\text{Gaussian error}}$$

Let us (try to) subtract the bias

$$\hat{\theta}^u = \hat{\theta}^* - (M\hat{\Sigma} - \mathbf{I})\hat{\theta}^{\text{Lasso}}$$

## Debiased estimator

$$\hat{\theta}^* = \theta_0 + \underbrace{(M\hat{\Sigma} - \mathbf{I})\theta_0}_{\text{bias}} + \underbrace{\frac{1}{n}M\mathbf{X}^T W}_{\text{Gaussian error}}$$

Let us (try to) subtract the bias

$$\hat{\theta}^u = \hat{\theta}^* - (M\hat{\Sigma} - \mathbf{I})\hat{\theta}^{\text{Lasso}}$$

Debiased estimator ( $\hat{\theta} = \hat{\theta}^{\text{Lasso}}$ )

$$\hat{\theta}^d \equiv \hat{\theta} + \frac{1}{n}M\mathbf{X}^T(Y - \mathbf{X}\hat{\theta})$$



## Debiased estimator: Choosing $M$ ?

$$\hat{\theta}^d \equiv \hat{\theta} + \frac{1}{n} M \mathbf{X}^T (y - \mathbf{X} \hat{\theta})$$

- Low dimensional projection estimator (LDPE)
  - ▶ Start with a linear estimator, debias by a nonlinear estimator
  - ▶  $M$  constructed via nodewise LASSO on  $\mathbf{X}$

[C-H. Zhang, S. S. Zhang]

## Debiased estimator: Choosing $M$ ?

$$\hat{\theta}^d \equiv \hat{\theta} + \frac{1}{n} M \mathbf{X}^\top (y - \mathbf{X} \hat{\theta})$$

- Low dimensional projection estimator (LDPE)
  - ▶ Start with a linear estimator, debias by a nonlinear estimator
  - ▶  $M$  constructed via nodewise LASSO on  $\mathbf{X}$

[C-H. Zhang, S. S. Zhang]

- Approximate inverse of  $\hat{\Sigma}$ : nodewise LASSO on  $\mathbf{X}$   
(under row-sparsity assumption on  $\Sigma^{-1}$ )

[S. van de Geer, P. Bühlmann, Y. Ritov, R. Dezeure]

## Debiased estimator: Choosing $M$ ?

Our approach:

- Optimizing two objectives (bias and variance of  $\hat{\theta}^d$ )

[A. Javanmard, A. Montanari]

$$\sqrt{n}(\hat{\theta}^d - \theta_0) = \underbrace{\sqrt{n}(M\hat{\Sigma} - \mathbf{I})(\theta_0 - \hat{\theta})}_{\text{bias} \downarrow} + Z$$

$$Z|\mathbf{X} \sim \mathbf{N}(0, \underbrace{\sigma^2 M \hat{\Sigma} M^\top}_{\text{noise covariance} \downarrow}), \quad \hat{\Sigma} = \frac{1}{n} \mathbf{X} \mathbf{X}^\top$$

## Debiased estimator: Choosing $M$ ?

Our approach:

- Find  $M$  by solving an optimization problem:

[A. Javanmard, A. Montanari]

$$\begin{aligned} & \underset{M}{\text{minimize}} && \max_{1 \leq i \leq p} (M \hat{\Sigma} M^T)_{i,i} \\ & \text{subject to} && |M \hat{\Sigma} - \mathbf{I}|_{\infty} \leq \xi \end{aligned}$$

## Debiased estimator: Choosing $M$ ?

Our approach:

- Find  $M$  by solving an optimization problem:

[A. Javanmard, A. Montanari]

$$\begin{aligned} & \underset{m_i}{\text{minimize}} && m_i^\top \widehat{\Sigma} m_i \\ & \text{subject to} && \|\widehat{\Sigma} m_i - e_i\|_\infty \leq \xi \end{aligned}$$

The optimization can be decoupled and solved in parallel.

# Main theorems

## Theorem [Javanmard, Montanari 2013] (Deterministic designs)

Let  $\mathbf{X}$  be any deterministic design that satisfies compatibility condition for the set  $S = \text{supp}(\theta_0)$ , ( $|S| \leq s_0$ ), with constant  $\phi_0$ . Further define the coherence parameter

$$\mu_* \equiv \min_{M \in \mathbb{R}^{p \times p}} |M\hat{\Sigma} - \mathbf{I}|_\infty.$$

Let  $K \equiv \max_{i \in [p]} \hat{\Sigma}_{ii}$ . Then, letting  $\lambda = c\sigma\sqrt{\log p/n}$ , we have

$$\sqrt{n}(\hat{\theta}^d - \theta_0) = Z + \Delta, \quad Z \sim \mathbf{N}(0, \sigma^2 M\hat{\Sigma}M^\top)$$

$$\mathbb{P}\left(\|\Delta\|_\infty \geq \frac{4c\mu_*\sigma s_0}{\phi_0^2} \sqrt{\log p}\right) \leq 2p^{-c_0}, \quad c_0 = \frac{c^2}{32K} - 1$$

# Main theorems

## Theorem [Javanmard, Montanari 2013] (Deterministic designs)

Let  $\mathbf{X}$  be any deterministic design that satisfies compatibility condition for the set  $S = \text{supp}(\theta_0)$ , ( $|S| \leq s_0$ ), with constant  $\phi_0$ . Further define the coherence parameter

$$\mu_* \equiv \min_{M \in \mathbb{R}^{p \times p}} |M\hat{\Sigma} - \mathbf{I}|_\infty.$$

Let  $K \equiv \max_{i \in [p]} \hat{\Sigma}_{ii}$ . Then, letting  $\lambda = c\sigma\sqrt{\log p/n}$ , we have

$$\sqrt{n}(\hat{\theta}^d - \theta_0) = Z + \Delta, \quad Z \sim \mathbf{N}(0, \sigma^2 M\hat{\Sigma}M^\top)$$

$$\mathbb{P}\left(\|\Delta\|_\infty \geq \frac{4c\mu_*\sigma s_0}{\phi_0^2} \sqrt{\log p}\right) \leq 2p^{-c_0}, \quad c_0 = \frac{c^2}{32K} - 1$$

### Remark:

$$\mu_* \leq \frac{1}{n} \max_{i \neq j} |\langle \mathbf{X}e_i, \mathbf{X}e_j \rangle|.$$

# Main theorems

## Theorem [Javanmard, Montanari 2013] (Random designs)

Let  $\Sigma$  be such that  $\sigma_{\min}(\Sigma) \geq C_{\min} > 0$  and  $\sigma_{\max}(\Sigma) \leq C_{\max} < \infty$  and  $\max_{i \in [p]} \Sigma_{ii} \leq 1$ . Assume  $\mathbf{X}\Sigma^{-1}$  has independent subgaussian rows with mean zero and subgaussian norm  $K$ . Letting  $\lambda = c\sigma\sqrt{\log p/n}$ , we have

$$\sqrt{n}(\hat{\theta}^d - \theta_0) = Z + \Delta, \quad Z|\mathbf{X} \sim \mathbf{N}(0, \sigma^2 M \hat{\Sigma} M^T),$$
$$\mathbb{P}\left\{ \|\Delta\|_{\infty} \geq \left(\frac{16c\sigma}{C_{\min}}\right) \frac{s_0 \log p}{\sqrt{n}} \right\} \leq 4e^{-c_1 n} + 4p^{-c_2},$$

for some explicit constants  $c_1 = C(K)$ ,  $c_2 = C(c, K, C_{\min}, C_{\max})$ .



# Main theorems

## Theorem [Javanmard, Montanari 2013] (Random designs)

Let  $\Sigma$  be such that  $\sigma_{\min}(\Sigma) \geq C_{\min} > 0$  and  $\sigma_{\max}(\Sigma) \leq C_{\max} < \infty$  and  $\max_{i \in [p]} \Sigma_{ii} \leq 1$ . Assume  $\mathbf{X}\Sigma^{-1}$  has independent subgaussian rows with mean zero and subgaussian norm  $K$ . Letting  $\lambda = c\sigma\sqrt{\log p/n}$ , we have

$$\sqrt{n}(\widehat{\theta}^d - \theta_0) = Z + \Delta, \quad Z|\mathbf{X} \sim \mathbf{N}(0, \sigma^2 M \widehat{\Sigma} M^T),$$
$$\mathbb{P}\left\{ \|\Delta\|_{\infty} \geq \left(\frac{16c\sigma}{C_{\min}}\right) \frac{s_0 \log p}{\sqrt{n}} \right\} \leq 4e^{-c_1 n} + 4p^{-c_2},$$

for some explicit constants  $c_1 = C(K)$ ,  $c_2 = C(c, K, C_{\min}, C_{\max})$ .

### Remark on sample size:

$$\text{If } \frac{n}{(s_0 \log p)^2} \rightarrow \infty \text{ then } \|\Delta\|_{\infty} = o_p(1).$$

# Consequences

- Confidence intervals for single parameters:

$$\lim_{n \rightarrow \infty} \mathbb{P}\left(\theta_{0,i} \in [\underline{\theta}_i, \bar{\theta}_i]\right) \geq 1 - \alpha$$

$$|\underline{\theta}_i - \bar{\theta}_i| \leq (2 + o(1))c_\alpha \sqrt{\frac{\sigma^2}{n}(\Sigma^{-1})_{ii}}$$

(n < p)

# Consequences

- Confidence intervals for single parameters:

$$\lim_{n \rightarrow \infty} \mathbb{P} \left( \theta_{0,i} \in [\underline{\theta}_i, \bar{\theta}_i] \right) \geq 1 - \alpha$$

$$|\underline{\theta}_i - \bar{\theta}_i| \leq (2 + o(1)) c_\alpha \sqrt{\frac{\sigma^2}{n} (\Sigma^{-1})_{ii}}$$

(n < p)

$$|\underline{\theta}_i - \bar{\theta}_i| \leq 2c_\alpha \sqrt{\frac{\sigma^2}{n} (\hat{\Sigma}^{-1})_{ii}}$$

Least square (n > p)

# Consequences

- Confidence intervals for single parameters:

$$\lim_{n \rightarrow \infty} \mathbb{P}\left(\theta_{0,i} \in [\underline{\theta}_i, \bar{\theta}_i]\right) \geq 1 - \alpha$$
$$|\underline{\theta}_i - \bar{\theta}_i| \leq (2 + o(1))c_\alpha \sqrt{\frac{\sigma^2}{n} (\Sigma^{-1})_{ii}}$$

(n < p)

$$|\underline{\theta}_i - \bar{\theta}_i| \leq 2c_\alpha \sqrt{\frac{\sigma^2}{n} (\hat{\Sigma}^{-1})_{ii}}$$

Least square (n > p)

## Remark:

No need for irrepresentability /  $\theta_{\min}$  condition  
(common assumptions for support recovery)

# Consequences

- Confidence intervals for single parameters:

$$\lim_{n \rightarrow \infty} \mathbb{P}\left(\theta_{0,i} \in [\underline{\theta}_i, \bar{\theta}_i]\right) \geq 1 - \alpha$$
$$|\underline{\theta}_i - \bar{\theta}_i| \leq (2 + o(1))c_\alpha \sqrt{\frac{\sigma^2}{n}(\Sigma^{-1})_{ii}}$$

(n < p)

- Hypothesis testing: minimax optimal statistical power

# Framework

- Hypothesis testing

$$H_{0,i} : \theta_{0,i} = 0, \quad H_{A,i} : \theta_{0,i} \neq 0.$$

- Two-sided p-values:

$$P_i = 2 \left( 1 - \Phi \left( \frac{|\hat{\theta}_i^d|}{\tau} \right) \right)$$

with  $\Phi(\cdot)$  cdf of standard normal.

- Decision rule:

$$T_{i,\mathbf{X}}(\mathbf{y}) = \begin{cases} 1 & \text{if } P_i \leq \alpha \quad (\text{reject the null hypothesis } H_{0,i}), \\ 0 & \text{otherwise} \quad (\text{accept the null hypothesis}). \end{cases}$$

## Theorem [Javanmard, Montanari, 2013]

Consider designs with subgaussian rows and let  $S \equiv \text{supp}(\theta_0)$ .

Assume that  $s_0 \equiv |S| = o(\sqrt{n}/\log p)$ .

Then, for any fixed sequence of integers  $i = i(n)$ , we have that for  $i \notin S$

$$\lim_{n \rightarrow \infty} \mathbb{P}_{\theta_0}(T_{i,\mathbf{X}}(y) = 1) \leq \alpha.$$

Further, assuming that for all  $i \in S$ ,  $|\theta_{0,i}| \geq \mu$ , we have

$$\liminf_{p \rightarrow \infty} \frac{1}{1 - \beta_{i,n}(\alpha, \mu)} \mathbb{P}_{\theta_0}(T_{i,\mathbf{X}}(y) = 1) \geq 1,$$

$$1 - \beta_{i,n}(\alpha, \mu) \equiv G\left(\alpha, \frac{\sqrt{n}\mu}{\sigma \sqrt{(\Sigma^{-1})_{ii}}}\right).$$

with  $G(\alpha, u)$  given by ...

## Theorem [Javanmard, Montanari, 2013]

Consider designs with subgaussian rows and let  $S \equiv \text{supp}(\theta_0)$ .

Assume that  $s_0 \equiv |S| = o(\sqrt{n}/\log p)$ .

Then, for any fixed sequence of integers  $i = i(n)$ , we have that for  $i \notin S$

$$\lim_{n \rightarrow \infty} \mathbb{P}_{\theta_0}(T_{i,\mathbf{X}}(y) = 1) \leq \alpha.$$

Further, assuming that for all  $i \in S$ ,  $|\theta_{0,i}| \geq \mu$ , we have

$$\liminf_{p \rightarrow \infty} \frac{1}{1 - \beta_{i,n}(\alpha, \mu)} \mathbb{P}_{\theta_0}(T_{i,\mathbf{X}}(y) = 1) \geq 1,$$

$$1 - \beta_{i,n}(\alpha, \mu) \equiv G\left(\alpha, \frac{\sqrt{n}\mu}{\sigma \sqrt{(\Sigma^{-1})_{ii}}}\right).$$

with  $G(\alpha, u)$  given by ...

Minimax optimal power over the family of  $s_0$ -sparse vectors  $\theta_0$ .



## Related work on bias-correction

- Ridge projection and bias correction [P. Bühlmann]
  - ▶ (Remaining) bias is not negligible.
  - ▶ Conservative tests
  
- Low dimensional projection estimator (LDPE) [C-H. Zhang, S. S. Zhang]
  - ▶ Initial projection based on nodewise LASSO on  $\mathbf{X}$ .
  - ▶ Bias correction via LASSO.

## Hypothesis testing under nearly optimal sample size

## Smaller sample size

- Estimation, prediction:  $n \gtrsim s_0 \log p$ .

[Candés, Tao 2007, Bickel et al. 2009]

- Hypothesis testing, confidence intervals:  $n \gtrsim (s_0 \log p)^2$ .

[This talk]

- ▶ Bias corrected ridge regression [P. Bühlmann]
- ▶ LDPE [C-H. Zhang, S. S. Zhang]
- ▶ Desparsified LASSO [S. van de Geer et. al.]

## Smaller sample size

- Estimation, prediction:  $n \gtrsim s_0 \log p$ .

[Candés, Tao 2007, Bickel et al. 2009]

- Hypothesis testing, confidence intervals:  $n \gtrsim (s_0 \log p)^2$ .

[This talk]

- ▶ Bias corrected ridge regression [P. Bühlmann]
- ▶ LDPE [C-H. Zhang, S. S. Zhang]
- ▶ Desparsified LASSO [S. van de Geer et. al.]

Can we match the optimal sample size,  $n \gtrsim s_0 \log p$  ?

## Theorem [Javanmard, Montanari 2013]

Consider designs with subgaussian rows and assume  $n \gtrsim s_0(\log p)^2$ . Then

$$\limsup_{p \rightarrow \infty} \frac{1}{p - s_0} \sum_{i \in S^c} \mathbb{P}_{\theta_0}(T_{i, \mathbf{X}}(y) = 1) \leq \alpha.$$

Further, assuming that for all  $i \in S$ ,  $|\theta_{0,i}| \geq \mu$ , we have

$$\liminf_{p \rightarrow \infty} \frac{1}{1 - \beta_n^*(\alpha, \theta_0)} \left\{ \frac{1}{s_0} \sum_{i \in S} \mathbb{P}_{\theta_0}(T_{i, \mathbf{X}}(y) = 1) \right\} \geq 1,$$

where

$$1 - \beta_n^*(\alpha, \theta_0) \equiv \frac{1}{s_0} \sum_{i \in S} G\left(\alpha, \frac{\sqrt{n}|\theta_{0,i}|}{\sigma \sqrt{(\Sigma^{-1})_{ii}}}\right).$$

## Theorem [Javanmard, Montanari 2013]

Consider designs with subgaussian rows and assume  $n \gtrsim s_0(\log p)^2$ . Then

$$\limsup_{p \rightarrow \infty} \frac{1}{p - s_0} \sum_{i \in S^c} \mathbb{P}_{\theta_0}(T_{i, \mathbf{X}}(y) = 1) \leq \alpha.$$

Further, assuming that for all  $i \in S$ ,  $|\theta_{0,i}| \geq \mu$ , we have

$$\liminf_{p \rightarrow \infty} \frac{1}{1 - \beta_n^*(\alpha, \theta_0)} \left\{ \frac{1}{s_0} \sum_{i \in S} \mathbb{P}_{\theta_0}(T_{i, \mathbf{X}}(y) = 1) \right\} \geq 1,$$

where

$$1 - \beta_n^*(\alpha, \theta_0) \equiv \frac{1}{s_0} \sum_{i \in S} G\left(\alpha, \frac{\sqrt{n}|\theta_{0,i}|}{\sigma \sqrt{(\Sigma^{-1})_{ii}}}\right).$$

- Controls **average** type I error
- Minimax optimal **average** power

# High-level idea of the proof

Recall

$$\sqrt{n}(\hat{\theta}^d - \theta_0) = Z + \underbrace{\Delta}_{\text{bias}}.$$

$$\|\Delta\|_{\infty} = O_p\left(\frac{s_0 \log p}{\sqrt{n}}\right) \longrightarrow n \gtrsim (s_0 \log p)^2$$

# High-level idea of the proof

Recall

$$\sqrt{n}(\hat{\theta}^d - \theta_0) = Z + \underbrace{\Delta}_{\text{bias}}.$$

$$\|\Delta\|_{\infty} = O_p\left(\frac{s_0 \log p}{\sqrt{n}}\right) \longrightarrow n \gtrsim (s_0 \log p)^2$$

To ensure **average** performance, we do not need to control  $\|\Delta\|_{\infty}$ .



## High-level idea of the proof

A new norm:

$$\|\Delta\|_{(\infty,k)} \equiv \max_{A \subset [p], |A| \geq k} \frac{\|\Delta_A\|_2}{\sqrt{|A|}}.$$

## High-level idea of the proof

A new norm:

$$\|\Delta\|_{(\infty,k)} \equiv \max_{A \subset [p], |A| \geq k} \frac{\|\Delta_A\|_2}{\sqrt{|A|}}.$$

Properties of  $\|\cdot\|_{\infty,k}$ :

- Non-increasing in  $k$
- As  $k$  gets smaller, it gives tighter control on the individual entries of  $\Delta$ .

## High-level idea of the proof

A new norm:

$$\|\Delta\|_{(\infty,k)} \equiv \max_{A \subset [p], |A| \geq k} \frac{\|\Delta_A\|_2}{\sqrt{|A|}}.$$

Properties of  $\|\cdot\|_{\infty,k}$ :

- Non-increasing in  $k$
- As  $k$  gets smaller, it gives tighter control on the individual entries of  $\Delta$ .

### Lemma

$$\|\Delta\|_{(\infty,cs_0)} = O\left(\sqrt{\frac{s_0}{n}} \log p\right).$$

## A few steps of the proof

- Any set  $|A| \geq cs_0$  can be partitioned as

$$A = A_1 \cup A_2 \cup \dots \cup A_L$$

with  $A_i$  disjoint and  $cs_0 \leq |A_i| \leq 2cs_0$ .

## A few steps of the proof

- Any set  $|A| \geq cs_0$  can be partitioned as

$$A = A_1 \cup A_2 \cup \dots \cup A_L$$

with  $A_i$  disjoint and  $cs_0 \leq |A_i| \leq 2cs_0$ .

(WLOG, we can assume  $cs_0 \leq |A| \leq 2cs_0$ .)

## A few steps of the proof

- Any set  $|A| \geq cs_0$  can be partitioned as

$$A = A_1 \cup A_2 \cup \dots \cup A_L$$

with  $A_i$  disjoint and  $cs_0 \leq |A_i| \leq 2cs_0$ .

(WLOG, we can assume  $cs_0 \leq |A| \leq 2cs_0$ .)

- 

$$\Delta \equiv \sqrt{n}(M\hat{\Sigma} - \mathbf{I})(\hat{\theta} - \theta_0).$$

## A few steps of the proof

- Any set  $|A| \geq cs_0$  can be partitioned as

$$A = A_1 \cup A_2 \cup \dots \cup A_L$$

with  $A_i$  disjoint and  $cs_0 \leq |A_i| \leq 2cs_0$ .

(WLOG, we can assume  $cs_0 \leq |A| \leq 2cs_0$ .)



$$\Delta \equiv \sqrt{n}(\Sigma^{-1}\widehat{\Sigma} - \mathbf{I})(\widehat{\theta} - \theta_0).$$

## A few steps of the proof

- Any set  $|A| \geq cs_0$  can be partitioned as

$$A = A_1 \cup A_2 \cup \dots \cup A_L$$

with  $A_i$  disjoint and  $cs_0 \leq |A_i| \leq 2cs_0$ .

(WLOG, we can assume  $cs_0 \leq |A| \leq 2cs_0$ .)

- 

$$\Delta \equiv \sqrt{n}(\Sigma^{-1}\widehat{\Sigma} - \mathbf{I})(\widehat{\theta} - \theta_0).$$

Let  $T = \text{supp}(\theta_0) \cup \text{supp}(\widehat{\theta})$ . We have

$$\|\Delta_A\|_2 \leq \sqrt{n} \|(\Sigma^{-1}\widehat{\Sigma} - \mathbf{I})_{A,T}\|_2 \|(\widehat{\theta} - \theta_0)_T\|_2.$$



## A few steps of the proof (cont'd)

- Applying tail bounds we get

$$\sup_{\substack{|A| \leq cs_0 \\ |T| \leq c's_0}} \|(\Sigma^{-1}\widehat{\Sigma} - \mathbf{I})_{A,T}\|_2 = O\left(\sqrt{\frac{s_0 \log p}{n}}\right).$$

We also know that

$$\|\widehat{\theta} - \theta_0\|_2 = O\left(\sqrt{\frac{s_0 \log p}{n}}\right).$$

- Combining the bounds, we get

$$\frac{\|\Delta_A\|_2}{|A|} \leq O\left(\sqrt{\frac{s_0}{n}} \log p\right).$$

## A few steps of the proof (cont'd)

- Applying tail bounds we get

$$\sup_{\substack{|A| \leq cs_0 \\ |T| \leq c's_0}} \|(\Sigma^{-1}\widehat{\Sigma} - \mathbf{I})_{A,T}\|_2 = O\left(\sqrt{\frac{s_0 \log p}{n}}\right).$$

We also know that

$$\|\widehat{\theta} - \theta_0\|_2 = O\left(\sqrt{\frac{s_0 \log p}{n}}\right).$$

- Combining the bounds, we get

$$\frac{\|\Delta_A\|_2}{|A|} \leq O\left(\sqrt{\frac{s_0}{n}} \log p\right).$$

$$n \gtrsim s_0 (\log p)^2$$

## Standard gaussian design

Suppose  $X_{ij} \sim \mathbf{N}(0, 1)$  independently.

$$\hat{\theta}^d = \hat{\theta} + \frac{1}{n} \mathbf{X}^\top (y - \mathbf{X} \hat{\theta})$$

## Standard gaussian design

Suppose  $X_{ij} \sim N(0, 1)$  independently.

$$\hat{\theta}^d = \hat{\theta} + \frac{1}{n} \mathbf{X}^\top (y - \mathbf{X}\hat{\theta})$$

- SDL test [J., Montanari 2013]

$$\hat{\theta}^d = \hat{\theta} + \frac{\mathbf{d}}{n} \mathbf{X}^\top (y - \mathbf{X}\hat{\theta})$$

## Standard gaussian design

Suppose  $X_{ij} \sim N(0, 1)$  independently.

$$\hat{\theta}^d = \hat{\theta} + \frac{1}{n} \mathbf{X}^\top (y - \mathbf{X}\hat{\theta})$$

- SDL test [J., Montanari 2013]

$$\begin{aligned} \hat{\theta}^d &= \hat{\theta} + \frac{\mathbf{d}}{n} \mathbf{X}^\top (y - \mathbf{X}\hat{\theta}) \\ \mathbf{d} &= \left(1 - \frac{1}{n} \|\hat{\theta}\|_0\right)^{-1} \end{aligned}$$

Based on the analysis of Approximate Message Passing (AMP).

[M. Bayati, D. Donoho, A. Maleki, A. Montanari]

## Exact asymptotic characterization

$$\widehat{\theta}^d \equiv \widehat{\theta} + \frac{d}{n} \mathbf{X}^\top (y - \mathbf{X}\widehat{\theta})$$

### Theorem [M. Bayati, A. Montanari 2012]

Consider the standard gaussian setting where  $n/p \rightarrow \delta$ ,  $s_0/p \rightarrow \varepsilon$  and  $n\sigma^2 \rightarrow \sigma_\infty^2$ . If  $\delta \geq \varepsilon \log(1/\varepsilon)$ , then on finite-dimensional marginals

$$\widehat{\theta}^d = \theta_0 + \tau Z, \quad Z \sim \mathbf{N}(0, \mathbf{I}_{p \times p}),$$

with  $\tau, d$  given by ...

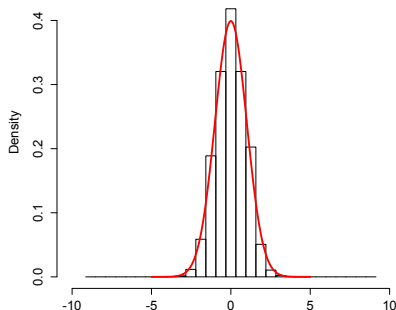
$$\delta \geq \varepsilon \log(1/\varepsilon) \longrightarrow n \geq s_0 \log(p/s_0)$$

## Effect of factor d

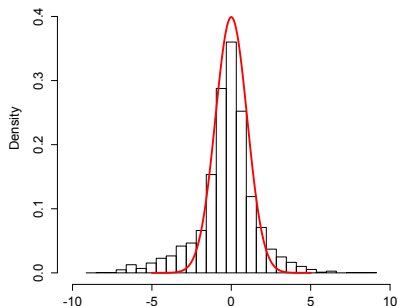
$$d = \left(1 - \frac{\|\hat{\boldsymbol{\theta}}\|_0}{n}\right)^{-1} = 1 + \mathcal{O}\left(\frac{s_0}{n}\right)$$

## Effect of factor d

$$d = \left(1 - \frac{\|\hat{\theta}\|_0}{n}\right)^{-1} = 1 + O\left(\frac{s_0}{n}\right)$$



with factor d



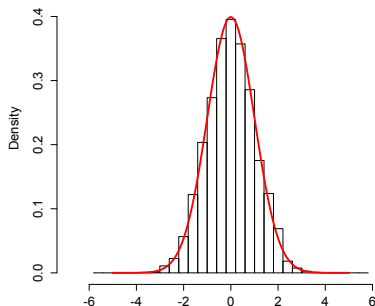
without factor d

Figure: Histogram of  $v = (\hat{\theta}^d - \theta_0)/\tau$  for  $n = 3s_0$  ( $\varepsilon = 0.2$ ,  $\delta = 0.6$ ) and  $p = 3000$ .

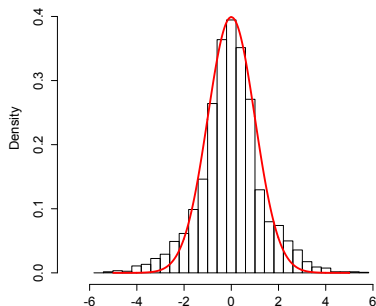


## Effect of factor d

$$d = \left(1 - \frac{\|\hat{\theta}\|_0}{n}\right)^{-1} = 1 + O\left(\frac{s_0}{n}\right)$$



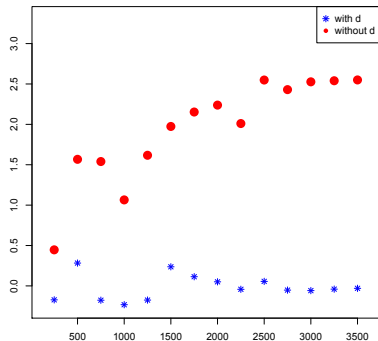
with factor d



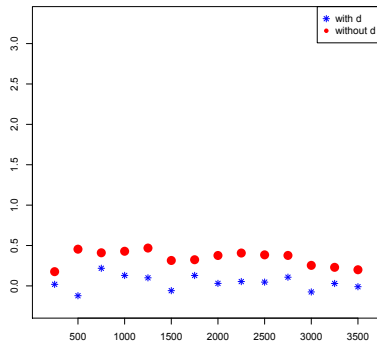
without factor d

**Figure:** Histogram of  $v = (\hat{\theta}^d - \theta_0)/\tau$  for  $n = 30s_0$  ( $\varepsilon = 0.02$ ,  $\delta = 0.6$ ) and  $p = 3000$ .

# Exact asymptotic characterization



$$n = 3s_0$$



$$n = 30s_0$$

Figure: Empirical kurtosis of  $v = (\hat{\theta}^d - \theta_0)/\tau$  with and without normalization factor d.

## Theorem [Javanmard, Montanari, 2013]

Consider the setting where  $n/p \rightarrow \delta$ ,  $s_0/p \rightarrow \varepsilon$  and  $\delta \geq \varepsilon \log(1/\varepsilon)$ . Then, for  $i \notin S$  we have

$$\lim_{p \rightarrow \infty} \mathbb{P}_{\theta_0}(T_{i,\mathbf{X}}(y) = 1) = \alpha.$$

Further, assuming that for all  $i \in S$ ,  $|\theta_{0,i}| \geq \mu$ , we have

$$\lim_{p \rightarrow \infty} \mathbb{P}_{\theta_0}(T_{i,\mathbf{X}}(y) = 1) \geq G\left(\alpha, \frac{\mu}{\tau}\right),$$

with

$$G(\alpha, u) \equiv 2 - \Phi\left(\Phi^{-1}\left(1 - \frac{\alpha}{2}\right) + u\right) - \Phi\left(\Phi^{-1}\left(1 - \frac{\alpha}{2}\right) - u\right).$$

# Summary

- Random designs

- $n \gtrsim s_0(\log p)^2 \checkmark$

- ▶ guarantee on **average** type I error and power
    - ▶ requires good estimate of precision matrix  
(can be done e.g., under sparsity assumption)

- $n \gtrsim s_0 \log p ?$

- Standard gaussian designs

- $n \gtrsim s_0 \log(p/s_0) \checkmark$



[1] A. Javanmard and A. Montanari, *Confidence Intervals and Hypothesis Testing for High-Dimensional Regression*. JMLR, 2014



[2] A. Javanmard and A. Montanari, *Nearly Optimal Sample Size in Hypothesis Testing for High-Dimensional Regression*. Allerton, 2013.



[3] A. Javanmard and A. Montanari, *Hypothesis Testing in High-Dimensional Regression under the Gaussian Random Design Model: Asymptotic Theory*. IEEE transaction on Info. Theory, 2013.

Thanks!