

De-biasing the Lasso: Optimal Sample Size for Gaussian Designs

Adel Javanmard

USC Marshall School of Business
Data Science and Operations department

Based on joint work with

Andrea Montanari

Oct 2015

An example



Kaggle challenge: Identify patients diagnosed with **type-2 diabetes**

Statistical model

Data $(Y_1, X_1), \dots, (Y_n, X_n)$:

- $Y_i =$ Patient i gets type-2 diabetes $\in \{0, 1\}$
- $X_i =$ Features of patient i $\in \mathbb{R}^p$

$$Y_i \sim f_{\theta_0}(\cdot | X_i) \quad \theta_0 \in \mathbb{R}^p$$

$\theta_{0,j}$ = contribution of feature j

Statistical model

Data $(Y_1, X_1), \dots, (Y_n, X_n)$:

- $Y_i =$ Patient i gets type-2 diabetes $\in \{0, 1\}$
- $X_i =$ Features of patient i $\in \mathbb{R}^p$

$$Y_i \sim f_{\theta_0}(\cdot | X_i) \quad \theta_0 \in \mathbb{R}^p$$

$\theta_{0,j}$ = contribution of feature j

Statistical model

Data $(Y_1, X_1), \dots, (Y_n, X_n)$:

- $Y_i =$ Patient i gets type-2 diabetes $\in \{0, 1\}$
- $X_i =$ Features of patient i $\in \mathbb{R}^p$

$$Y_i \sim f_{\theta_0}(\cdot | X_i) \quad \theta_0 \in \mathbb{R}^p$$

$\theta_{0,j}$ = contribution of feature j

Regularized estimator

$$\hat{\theta} \equiv \arg \min_{\theta \in \mathbb{R}^p} \left(\underbrace{\mathcal{L}(\theta)}_{\text{logistic loss}} + \lambda \underbrace{\|\theta\|_1}_{\text{regularizer}} \right).$$

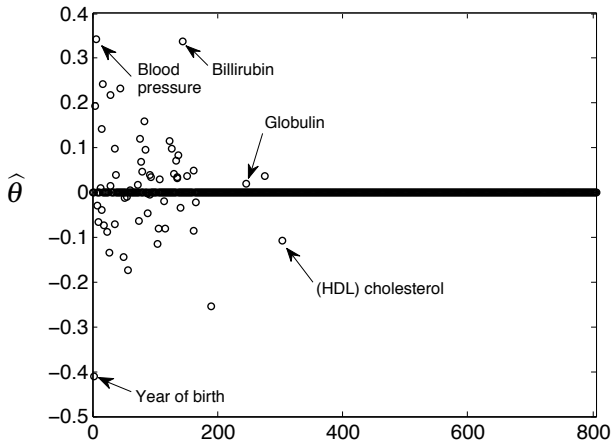
- Convex optimization
- Variable selection

Practice fusion data set (Kaggle)

Database

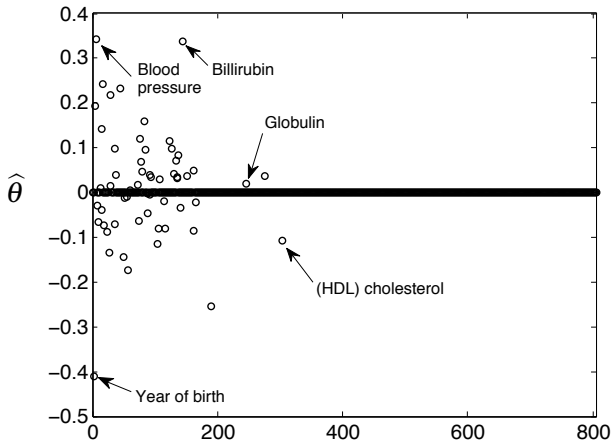


- $n = 500$: patients
- $p = 805$: medical information (meds, lab results, diagnosis, ...)



Regularized logreg selects 62 features
 (λ chosen via cross validation resulting $AUC = 0.75$)

Shall we trust our findings?



Regularized logreg selects 62 features
 (λ chosen via cross validation resulting $AUC = 0.75$)

Shall we **trust** our findings?

In summary

- Will focus on linear model and Lasso
- Compute confidence intervals/p-values

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

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

[Tibshirani 1996, Chen, Donoho 1996]

- Distribution of $\hat{\theta}$?
- Debiasing approach:
(LASSO is biased towards small ℓ_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 ℓ_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:

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

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

Classical setting ($n \gg p$)

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

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

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

- Confidence intervals:

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

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

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

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

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

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 ?

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

- Gaussian design ($x_i \sim \mathcal{N}(0, \Sigma)$)
 - ▶ Assume known Σ (relevant in semi-supervised learning)
 - ▶ $M = \Sigma^{-1}$

[Javanmard, Montanari 2012]

Debiased estimator: Choosing M ?

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

- Gaussian design ($x_i \sim \mathcal{N}(0, \Sigma)$)
 - ▶ Assume known Σ (relevant in semi-supervised learning)
 - ▶ $M = \Sigma^{-1}$

[Javanmard, Montanari 2012]

Does this remind you anything?

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

Debiased estimator: Choosing M ?

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

- Gaussian design ($x_i \sim N(0, \Sigma)$)
 - ▶ Assume known Σ (relevant in semi-supervised learning)
 - ▶ $M = \Sigma^{-1}$

[Javanmard, Montanari 2012]

Does this remind you anything?

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

(pseudo-) Newton method

Debiased estimator: Choosing M ?

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

- Gaussian design ($x_i \sim \mathcal{N}(0, \Sigma)$)
 - ▶ Assume known Σ (relevant in semi-supervised learning)
 - ▶ $M = \Sigma^{-1}$

[Javanmard, Montanari 2012]

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

[S. van de Geer, P. Bühlmann, Y. Ritov, R. Dezeure 2014]

Debiased estimator: Choosing M ?

Our approach:

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

[A. Javanmard, A. Montanari 2014]

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

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

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:

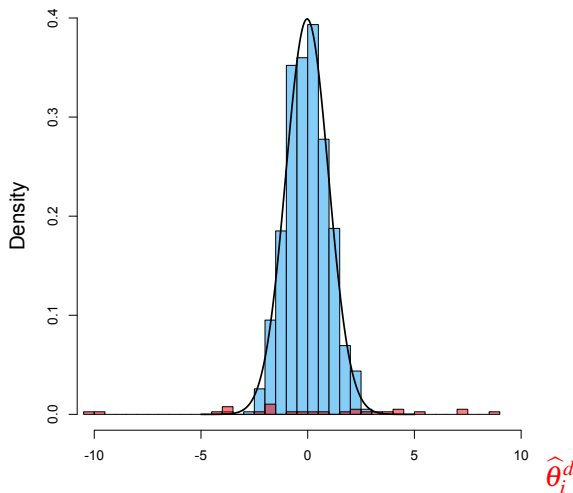
- Find M by solving an optimization problem:

[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.

What does it look like?



$\hat{\theta}_i^d$

Can estimate σ
'Ground truth' from $n_{\text{tot}} = 10,000$ records.

Confidence intervals

Neglecting the bias ($\widehat{\sigma}$ estimator of σ)

$$\widehat{\theta}_i^d \approx N(\theta_{0,i}, \Delta_i^2), \quad \Delta_i^2 \equiv \frac{\widehat{\sigma}^2}{n} (M\widehat{\Sigma}M^T)_{ii}$$

$$[\underline{\theta}_i, \bar{\theta}_i] = [\widehat{\theta}_i^d - c_\alpha \Delta_i, \widehat{\theta}_i^d + c_\alpha \Delta_i]$$

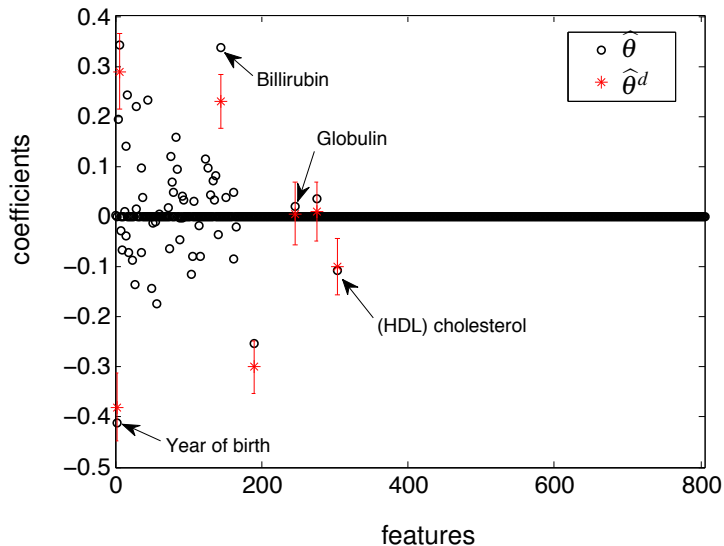
Confidence intervals

Neglecting the bias ($\hat{\sigma}$ estimator of σ)

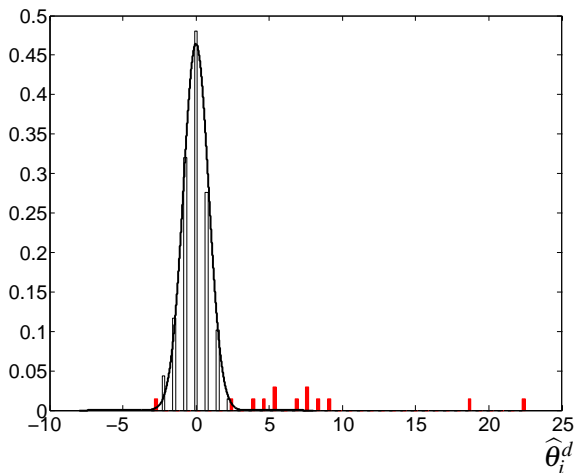
$$\hat{\theta}_i^d \approx N(\theta_{0,i}, \Delta_i^2), \quad \Delta_i^2 \equiv \frac{\hat{\sigma}^2}{n} (M\hat{\Sigma}M^T)_{ii}$$

$$[\underline{\theta}_i, \bar{\theta}_i] = [\hat{\theta}_i^d - c_\alpha \Delta_i, \hat{\theta}_i^d + c_\alpha \Delta_i]$$

What does it look like?



UCI crime dataset



$$n = 84, p = 102, n_{\text{tot}} = 1994.$$

A theorem

Theorem [Javanmard, Montanari 2013] (Deterministic designs)

Let \mathbf{X} be any deterministic design that satisfies compatibility condition. Define the coherence parameter

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

Let $s_0 = |\text{supp}(\theta_0)|$. Then

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

$$\|\Delta\|_\infty \leq c\mu_*\sigma s_0\sqrt{\log p}, \quad \text{w.h.p.}$$

A theorem

Theorem [Javanmard, Montanari 2013] (Deterministic designs)

Let \mathbf{X} be any deterministic design that satisfies compatibility condition. Define the coherence parameter

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

Let $s_0 = |\text{supp}(\theta_0)|$. Then

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

$$\|\Delta\|_\infty \leq c\mu_*\sigma s_0\sqrt{\log p}, \quad \text{w.h.p.}$$

Remark:

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

A theorem

Theorem [Javanmard, Montanari 2013] (Random designs)

Consider population covariance Σ with bounded eigenvalues and assume $\mathbf{X}\Sigma^{-1}$ has independent subgaussian rows. Then

$$\sqrt{n}(\hat{\theta}^d - \theta_0) = \underbrace{Z}_{\text{Gaussian}} + \underbrace{\Delta}_{\text{Bias}}$$
$$\|\Delta\|_{\infty} \leq c\sigma \frac{s_0 \log p}{\sqrt{n}}, \quad \text{w.h.p.}$$

A theorem

Theorem [Javanmard, Montanari 2013] (Random designs)

Consider population covariance Σ with bounded eigenvalues and assume $\mathbf{X}\Sigma^{-1}$ has independent subgaussian rows. Then

$$\sqrt{n}(\hat{\theta}^d - \theta_0) = \underbrace{Z}_{\text{Gaussian}} + \underbrace{\Delta}_{\text{Bias}}$$
$$\|\Delta\|_{\infty} \leq c\sigma \frac{s_0 \log p}{\sqrt{n}}, \quad \text{w.h.p.}$$

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:

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

(n < p)

Consequences

- Confidence intervals for single parameters:

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

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

(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:

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

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

(n < p)

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

Least square (n > p)

Remark:

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

Hypothesis testing (based on de-biased estimator)

- Null/alternative hypothesis:

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

- Two-sided p-values:

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

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

- We provide precise characterization of type I and type II error.
- Test (using de-biased estimator) has minimax optimal statistical power.

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.

Further related work

- Debiasing:
 - Group sparsity [R. Mitra & C.H.Zhang 2014]
 - Confidence interval for inverse covariance estimation [J. Jankova, S.v.d. Geer 2015]
 - Genomics [Q.Zhao et. al. 2015, B. Rakitsch 2015]
 - Econometrics [A. Belloni & V. Chernozhukov 2014, D. Kozbur 2015]
- Other methods for uncertainty assessment
 - Uncertainty quantification under group sparsity [Q.Zhou 2015]
 - Post double selection [Belloni et. al. 2014]

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, $n \gtrsim s_0 \log p$?

Where is the bottleneck?

Where is the bottleneck?

The bias is given by

$$\Delta = \sqrt{n}(\Omega\hat{\Sigma} - \mathbf{I})(\theta_0 - \hat{\theta}^{\text{Lasso}}).$$

Earlier work bound bias a simple $\ell_1 - \ell_\infty$ inequality:

$$\begin{aligned}\|\Delta\|_\infty &\leq \sqrt{n}\|M\hat{\Sigma} - \mathbf{I}\|_\infty\|\theta^* - \hat{\theta}^{\text{Lasso}}\|_1 \\ &\leq \sqrt{n} \times C\sqrt{\frac{\log p}{n}} \times Cs_0\sigma\sqrt{\frac{\log p}{n}} \\ &\leq C^2\sigma\frac{s_0\log p}{\sqrt{n}}.\end{aligned}$$

Plan for this part

- Focus on Gaussian design: $x_i \sim \mathbf{N}(0, \Sigma)$
- Assume that Σ is known. (See paper for unknown covariance.)
- We show that the required sample rate is indeed artifact of the argument!

Plan for this part

- Focus on Gaussian design: $x_i \sim \mathbf{N}(0, \Sigma)$
- Assume that Σ is known. (See paper for unknown covariance.)
- We show that the required sample rate is indeed artifact of the argument!

De-biased estimator is asymptotically Gaussian under condition $n \gtrsim s_0(\log p)^2$.

'Leave-one-out' technique

Fix coordinate i .

- Define

$$\begin{aligned}\hat{\theta}^p &\equiv \arg \min_{\theta} \frac{1}{2n} \|y - X\theta\|^2 + \lambda \|\theta\|_1 \\ &\text{subject to } \hat{\theta}_i^p = \theta_{0,i}\end{aligned}$$

'Leave-one-out' technique

Fix coordinate i .

- Define

$$\begin{aligned}\hat{\theta}^p &\equiv \arg \min_{\theta} \frac{1}{2n} \|y - X\theta\|^2 + \lambda \|\theta\|_1 \\ &\text{subject to } \hat{\theta}_i^p = \theta_{0,i}\end{aligned}$$

We then have

$$y - X\hat{\theta}^p = w + \tilde{x}_i(\theta_{0,i} - \hat{\theta}_i^p) + X_{\sim i}(\theta_{0,\sim i} - \hat{\theta}_{\sim i}^p)$$

'Leave-one-out' technique

Fix coordinate i .

- Define

$$\begin{aligned}\hat{\theta}^p &\equiv \arg \min_{\theta} \frac{1}{2n} \|y - X\theta\|^2 + \lambda \|\theta\|_1 \\ &\text{subject to } \hat{\theta}_i^p = \theta_{0,i}\end{aligned}$$

We then have

$$y - X\hat{\theta}^p = w + \tilde{x}_i(\theta_{0,i} - \hat{\theta}_i^p) + X_{\sim i}(\theta_{0,\sim i} - \hat{\theta}_{\sim i}^p)$$

$\hat{\theta}^p$ is the Lasso estimator when \tilde{x}_i is left out!

'Leave-one-out' technique

Let v be the i th column of $X\Sigma^{-1}$.

The bias is given by

$$\Delta_i = R_1 + R_2 + R_3$$

'Leave-one-out' technique

Let v be the i th column of $X\Sigma^{-1}$.

The bias is given by

$$\Delta_i = R_1 + R_2 + R_3$$

$$R_1 = \sqrt{n} \left(1 - \frac{\langle v, \tilde{x}_i \rangle}{n} \right) (\hat{\theta}_i^{\text{Lasso}} - \theta_i^*)$$

$$R_2 = \frac{v^\top}{\sqrt{n}} X_{\sim i} (\theta_{0, \sim i} - \hat{\theta}_{\sim i}^{\text{p}})$$

$$R_3 = \frac{v^\top}{\sqrt{n}} X_{\sim i} (\hat{\theta}_{\sim i}^{\text{p}} - \hat{\theta}_{\sim i}^{\text{Lasso}})$$

'Leave-one-out' technique

Let v be the i th column of $X\Sigma^{-1}$.

The bias is given by

$$\Delta_i = R_1 + R_2 + R_3$$

$$R_1 = \sqrt{n} \left(1 - \frac{\langle v, \tilde{x}_i \rangle}{n} \right) (\hat{\theta}_i^{\text{Lasso}} - \theta_i^*) \xrightarrow{\text{concentration}}$$

$$R_2 = \frac{v^\top}{\sqrt{n}} X_{\sim i} (\theta_{0, \sim i} - \hat{\theta}_{\sim i}^{\text{p}}) \xrightarrow{\text{independence}}$$

$$R_3 = \frac{v^\top}{\sqrt{n}} X_{\sim i} (\hat{\theta}_{\sim i}^{\text{p}} - \hat{\theta}_{\sim i}^{\text{Lasso}}) \xrightarrow{\text{perturbation}}$$



Summary

Combining the bounds on R_1, R_2, R_3 , we obtain

$$\|\Delta\|_\infty \leq C \sqrt{\frac{s_0}{n}} \log p, \quad \text{w.h.p}$$

Summary

Combining the bounds on R_1, R_2, R_3 , we obtain

$$\|\Delta\|_\infty \leq C \sqrt{\frac{s_0}{n}} \log p, \quad \text{w.h.p}$$

Therefore,

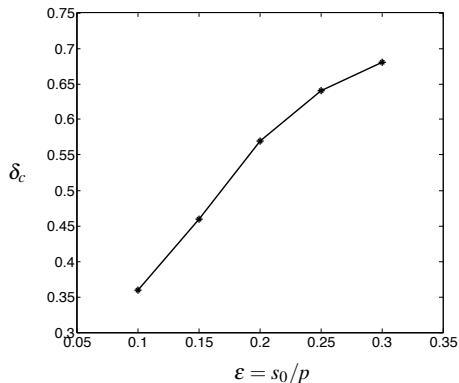
$$\|\Delta\|_\infty \rightarrow 0 \quad \text{provided that} \quad n \geq s_0(\log p)^2.$$

Numerical illustration

- Fix $p = 3000$
- Design matrix X with rows i.i.d. from $N(0, \Sigma)$
- $\Sigma_{ij} = 0.8^{|i-j|}$
- Define $\delta = n/p$ (undersampling rate) and $\varepsilon = s_0/p$ (sparsity proportion)
- δ_c : Critical value above which the de-biased estimator is Gaussian.

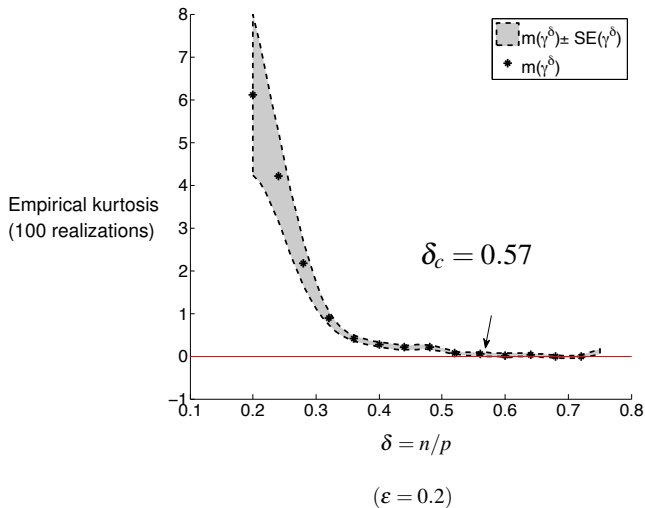
Numerical illustration

- Fix $p = 3000$
- Design matrix X with rows i.i.d. from $N(0, \Sigma)$
- $\Sigma_{ij} = 0.8^{|i-j|}$
- Define $\delta = n/p$ (undersampling rate) and $\varepsilon = s_0/p$ (sparsity proportion)
- δ_c : Critical value above which the de-biased estimator is Gaussian.



How to define δ_c ?

- Fix ε and change $\delta = n/p$.



Conclusion

- De-biasing regularized estimators
- Compute confidence intervals/p-values for high dimensional models
- Optimal sample size for Gaussian designs

Thanks!

