

A statistical mechanics approach to de-biasing and uncertainty estimation in LASSO



← Takashi Takahashi and YK
Tokyo Institute of Technology,
Japan

Outline

- Motivation and setup
- Statistical mechanics approach
 - Replica analysis
 - Cavity/TAP approach
- Experimental validation
- Summary

Underdetermined sparse regression

$$\begin{array}{c} M(<N) \\ \left\{ \begin{array}{c} \text{Data} \\ \mathbf{y} \end{array} \right\} = \begin{array}{c} M \\ \left\{ \begin{array}{c} \mathbf{A} \\ \text{(Sensing/design matrix)} \end{array} \right\} \end{array} \begin{array}{c} N \\ \left\{ \begin{array}{c} \mathbf{x}_0 \\ \text{Parameter / signal} \end{array} \right\} \end{array} + \begin{array}{c} \text{Noise vector} \\ \mathbf{n} \end{array}$$

- Basic example of high-dimensional statistics
- Underdetermined as $M < N$
 - Unique estimator can not be determined by minimizing $\mathcal{L} = \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$ squared loss.
 - There are uncountably many solutions...

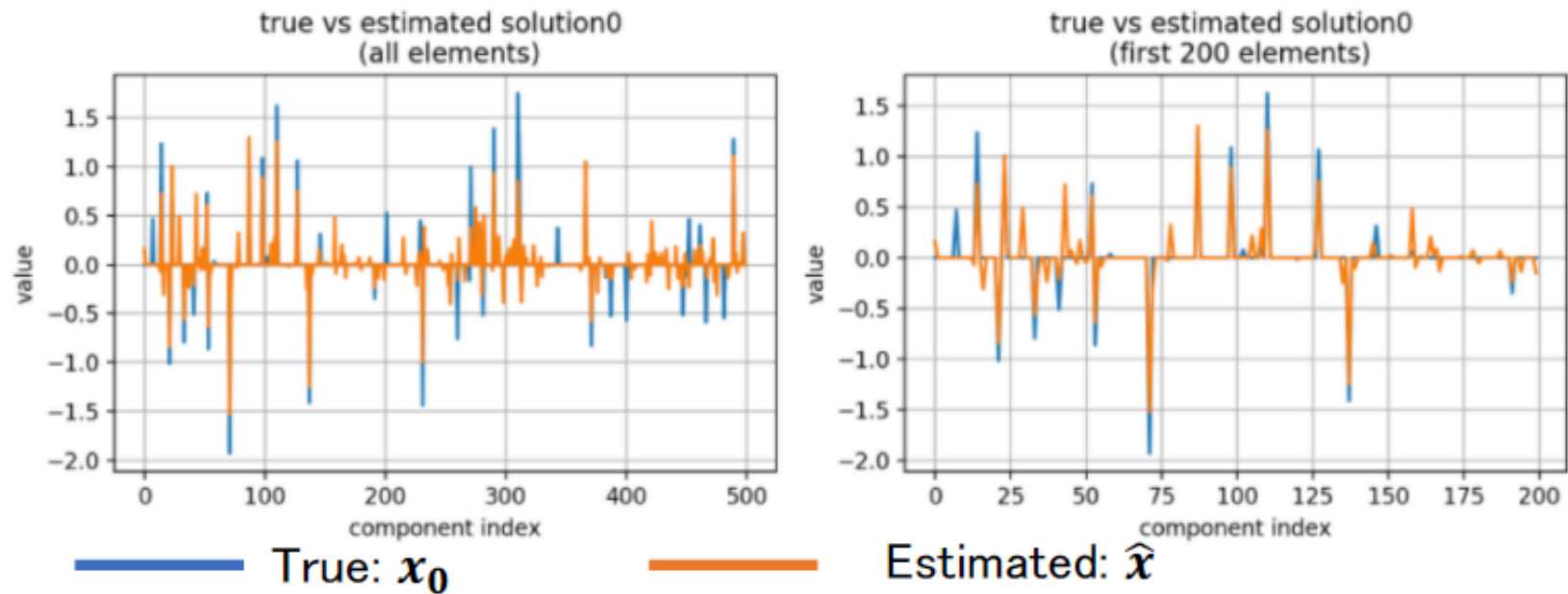
LASSO/L1-norm regularized estimator

- A popular solution:
Least absolute shrinkage and selection operator (LASSO)

$$\hat{\mathbf{x}}^{\text{LASSO}}(\mathbf{y}, A; \lambda) = \arg \min_x \left\{ \underbrace{\frac{1}{2} \|\mathbf{y} - A\mathbf{x}\|_2^2}_{\text{loss}} + \underbrace{\lambda \|\mathbf{x}\|_1}_{\text{regularizer}} \right\}$$

- Formulated as a convex optimization problem
- An estimator is given as a numerical solution
- Obtained estimator is sparse
 - LASSO performs on
regression + *variable selection*

Example of LASSO solution



■ Parameters

$$A_{\mu i} \sim \text{i.i.d. } \mathcal{N}(0, N^{-1/2}), x_{0,i} \sim \text{i.i.d. } 0.9\delta(x) + 0.1\mathcal{N}(0,1), n \sim \text{i.i.d. } \mathcal{N}(0, 0.05)$$

$$N = 500, \alpha = M/N = 0.5$$

Good properties of LASSO

1. Easy to perform

$$Q_1 : \min_x \left\{ \frac{1}{2} \|y - Ax\|_2^2 + \lambda \|x\|_1 \right\}$$

is a convex optimization problem

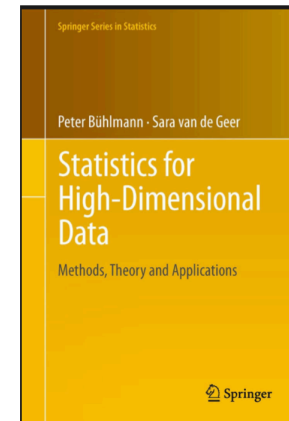
2. Consistent in ideal settings

$$\|x_0 - \hat{x}\|_2 \leq \frac{Cs_0\sigma^2}{M} \log N \quad \text{with high probability}$$

when

- x_0 is s_0 -sparse ($\|x_0\|_0 = s_0$)
- A satisfies certain “good” properties

See “Statistics for High-Dimensional Data” (Bulmann and van de Geer) for details



Unsatisfactory properties of LASSO

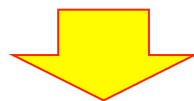
$$\hat{\mathbf{x}}^{\text{LASSO}}(\mathbf{y}, A; \lambda) = \arg \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{y} - A\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1 \right\}$$

1. Biased

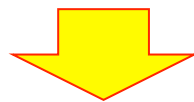
$$\left| \mathbb{E} \left[\hat{x}_i^{\text{LASSO}} \right]_{A, \xi} - x_{o,i} \right| > 0 \text{ for } \lambda > 0$$

2. Dist. of $\hat{\mathbf{x}}^{\text{LASSO}}$ is unknown

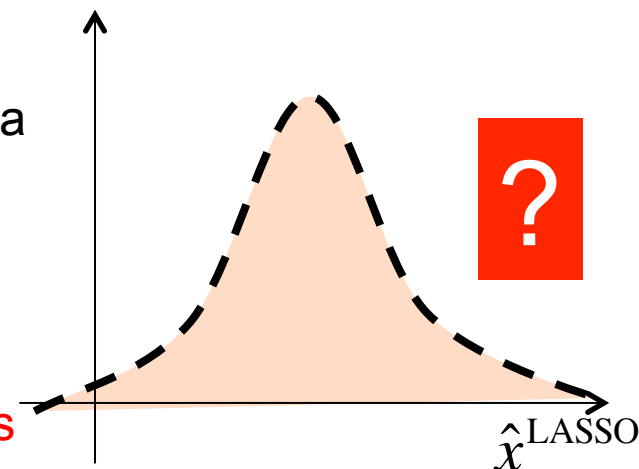
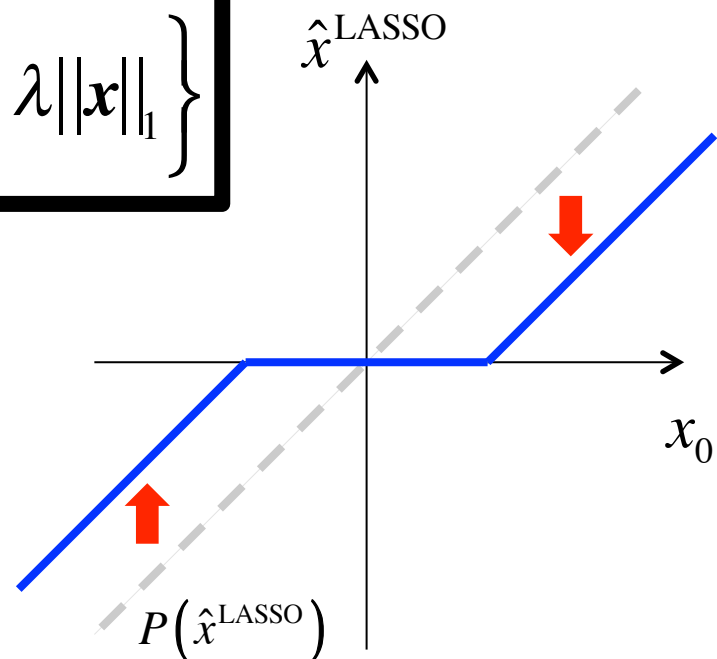
$\hat{\mathbf{x}}^{\text{LASSO}}(\mathbf{y}, A; \lambda)$: No analytical expression
Complicated dependence on data



We cannot evaluate its **confidence interval**



We cannot quantify the **reliability of obtained results**



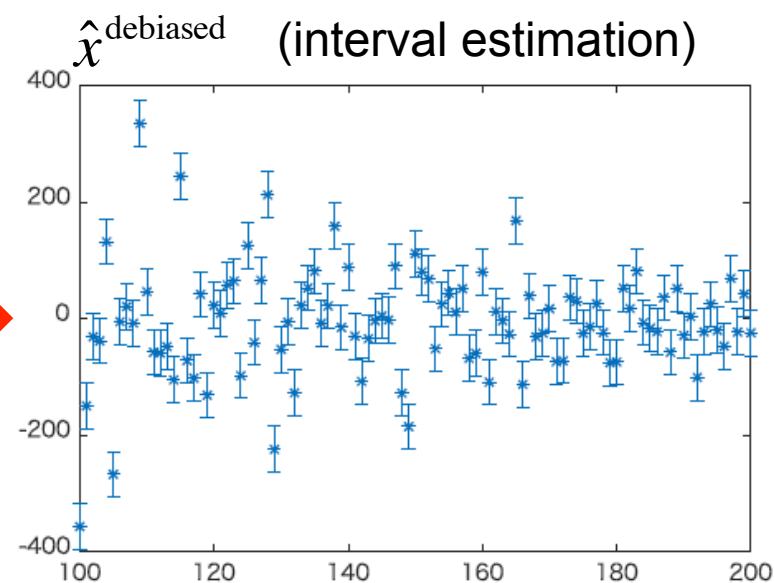
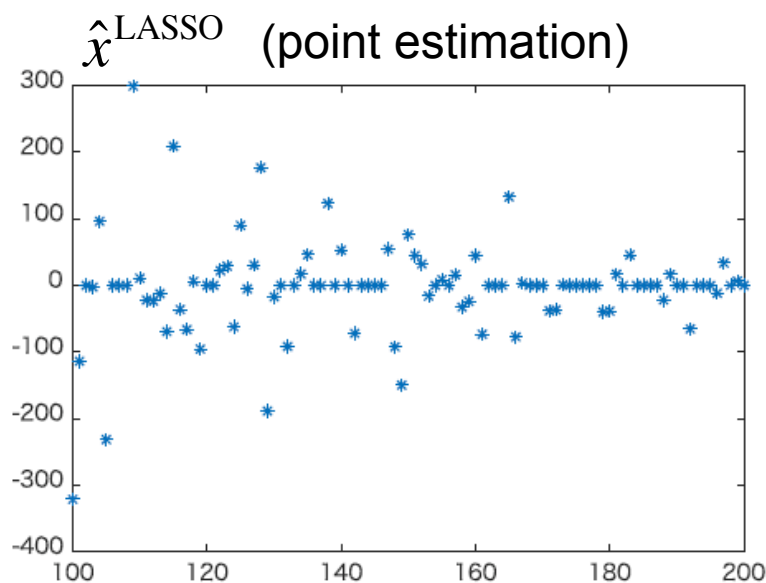
Aim of this talk: Post-processing of LASSO

We aim to construct an estimator $\hat{x}^{\text{debiased}}$ that has the following properties from the LASSO solution.

1. Unbiased

$$\left| \mathbb{E} \left[\hat{x}_i^{\text{debiased}} \right]_{A,n} - x_{o,i} \right| = 0 \text{ for } \forall \lambda > 0$$

2. Confidence intervals are available



Related work

- Javanmard and Montanari (2014)
- van de Geer et al (2014)
- Zhang and Zhang (2014)
 - Fixed \mathbf{x}_0
 - Fixed A
 - Very sparse signals $\rho = \|\mathbf{x}_0\|_0 / N \rightarrow 0$
 - Mathematically rigorous
- Ours
 - Fixed \mathbf{x}_0
 - Random A
 - Finite sparsity $\rho = \|\mathbf{x}_0\|_0 / N = O(1)$
 - Careful review of the existing stat. mech. results

Setup

$$\underbrace{\begin{matrix} M \\ \left\{ \right. \end{matrix} \mathbf{y}}_{\text{Data}} = \underbrace{\begin{matrix} M \\ \left\{ \right. \end{matrix} \underbrace{\mathbf{A}}_{\text{(Sensing/design matrix)}} \underbrace{\begin{matrix} \left. \right\} N \\ \mathbf{x}_0 \end{matrix}}_{\text{Parameter /signal}} + \underbrace{\mathbf{n}}_{\text{Noise vector}}$$

■ Assumptions

- $\mathbf{n} \sim \mathcal{N}(0_M, \sigma^2 I_M)$
- \mathbf{x}_0 : ρN - sparse $(\|\mathbf{x}_0\|_0 = N\rho, \rho = O(1); \|\mathbf{x}_0\|_2^2 = N\rho\sigma_x^2)$
- \mathbf{A} : random matrix satisfying the following property
 $\mathbf{A}^\top \mathbf{A} = \mathbf{O} \mathbf{D} \mathbf{O}^\top$: eigenvalue decomposition
 $\begin{cases} \mathbf{O} \sim \text{uniform dist. of } O(N) \\ \mathbf{D}_{ii} \sim \text{a fixed eigenvalue dist.} \end{cases} \leftarrow \text{Can be relaxed slightly } \mu_{\mathbf{A}^\top \mathbf{A}}(s)$

Statistical mechanics approach

- LASSO cost = Hamiltonian

$$\frac{1}{2} \|\mathbf{y} - A\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1 = H(\mathbf{x}|\mathbf{y}, A; \lambda)$$

- Gibbs-Boltzmann distribution

$$p_\beta(\mathbf{x}|\mathbf{y}, A; \lambda) = \frac{e^{-\beta H(\mathbf{x}|\mathbf{y}, A; \lambda)}}{Z_\beta(\mathbf{y}, A; \lambda)}$$

- LASSO solution = Mean in $\beta \rightarrow \infty$

$$\begin{aligned} \hat{\mathbf{x}}^{\text{LASSO}}(\mathbf{y}, A; \lambda) &= \arg \min_x \left\{ \frac{1}{2} \|\mathbf{y} - A\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1 \right\} \\ &= \lim_{\beta \rightarrow +\infty} \int \mathbf{x} p_\beta(\mathbf{x}|\mathbf{y}, A; \lambda) d\mathbf{x} \end{aligned}$$

One can analyze $\hat{\mathbf{x}}^{\text{LASSO}}$ using standard methods of stat. mech.

Replica + RMT analysis

- Replica symmetric free energy keeping \mathbf{x}_0 fixed

– Parisi and Potters (1995), Takeda et al (2006), ... $\left(\alpha \triangleq \frac{M}{N}, Dz \triangleq \frac{e^{-\frac{z^2}{2}}}{\sqrt{2\pi}} \right)$

$$f = -\lim_{\beta \rightarrow \infty} \lim_{n \rightarrow 0} \lim_{N \rightarrow \infty} \frac{1}{\beta} \frac{\partial}{\partial n} \ln \mathbb{E} \left[\left(Z_\beta(\mathbf{y}, A; \lambda) \right)^n \right]_{A,n}$$

$$= \text{extr}_{\chi, \hat{\chi}, Q, \hat{Q}, m, \hat{m}} \left\{ \begin{aligned} & -G'(-\chi; A^\top A) (Q - 2m + \rho \sigma_X^2 - \chi \sigma^2) - \frac{\alpha}{2} \sigma^2 + \frac{\hat{Q}Q}{2} - \frac{\hat{\chi}\chi}{2} - \hat{m}m \\ & + \lim_{N \rightarrow \infty} \sum_{i \sim 1}^N \int \min_{x_i} \left\{ \frac{\hat{Q}}{2} x_i^2 - \left(\hat{m} x_{o,i} + \sqrt{\hat{\chi}} z_i \right) x_i + \lambda |x_i| \right\} Dz_i \end{aligned} \right\}$$

➔ $\mathbb{E} \left[\left(\hat{x}_i^{\text{LASSO}} \right)^p \right]_{A,n} = \int \left(\arg \min_{x_i} \left\{ \frac{\hat{Q}}{2} x_i^2 - \left(\hat{m} x_{o,i} + \sqrt{\hat{\chi}} z_i \right) x_i + \lambda |x_i| \right\} \right)^p Dz_i$

$$G(x; A^\top A) \triangleq \text{extr}_z \left\{ -\frac{1}{2} \int \rho_{A^\top A}(s) \ln(z - s) ds + \frac{zx}{2} \right\} - \frac{1}{2} \ln|x| - \frac{1}{2}$$

Characterizes the property of the matrix ensemble

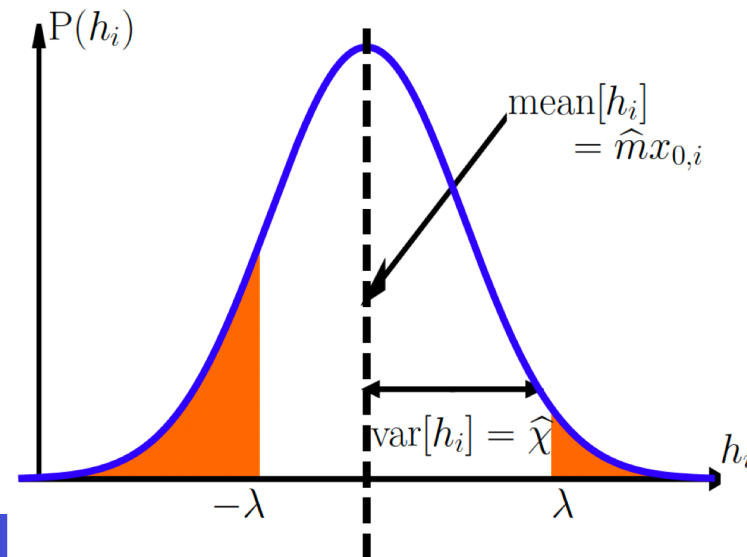
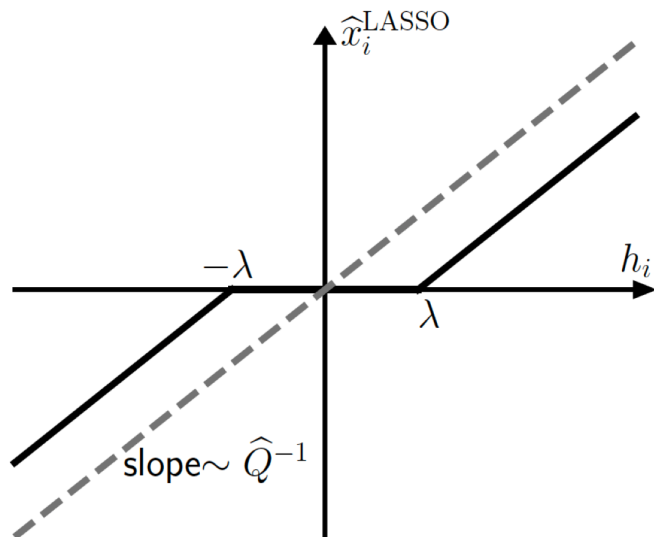
Decoupling principle

- Estimator of each component: Soft-thresholding of an **statistically independent Gaussian field**

$$\mathbb{E} \left[\left(\hat{x}_i^{\text{LASSO}} \right)^p \right]_{A,n} = \int \left(\arg \min_{x_i} \left\{ \frac{\hat{Q}}{2} x_i^2 - \underbrace{\left(\hat{m} x_{0,i} + \sqrt{\hat{\chi}} z_i \right)}_{h_i} x_i + \lambda |x_i| \right\} \right)^p D z_i$$

$$\equiv h_i \sim \mathcal{N}(\hat{m} x_{0,i}, \hat{\chi})$$

$$\hat{x}_i^{\text{LASSO}} = \hat{Q}^{-1} (h_i - \lambda \operatorname{sgn}(h_i)) \Theta(|h_i| - \lambda)$$



Consideration

- Statistical properties of the estimator are completely determined by 3 macroscopic variables $\hat{Q}, \hat{m} (= \hat{Q}), \hat{\chi}$.
- If these variables and the Gaussian field h_i are available, one can construct a debiased estimator $\hat{x}_i^{\text{debiased}}$ as follows:

$$\left\{ \begin{array}{l} \bullet h_i (= \hat{m}x_{0,i} + \sqrt{\hat{\chi}}z_i) \\ \bullet \hat{x}_i^{\text{debiased}} = \frac{h_i}{\hat{Q}} = \frac{\hat{m}}{\hat{Q}}x_{0,i} + \frac{\sqrt{\hat{\chi}}}{\hat{Q}}z_i = x_{0,i} + \frac{\sqrt{\hat{\chi}}}{\hat{Q}}z_i \sim \mathcal{N}\left(x_{0,i}, \frac{\hat{\chi}}{\hat{Q}^2}\right) \\ \bullet 100(1 - \alpha_{\text{sig}})\% \text{ CI: } \left[\hat{x}_i^{\text{debiased}} - \hat{Q}^{-1}\sqrt{\hat{\chi}}u\left(1 - \frac{\alpha_{\text{sig}}}{2}\right), \hat{x}_i^{\text{debiased}} + \hat{Q}^{-1}\sqrt{\hat{\chi}}u\left(1 - \frac{\alpha_{\text{sig}}}{2}\right) \right] \end{array} \right.$$

How can we know $\hat{Q}, \hat{m} (= \hat{Q}), \hat{\chi}$, and $\{h_i\}$ from a single sample of data?



Cavity/TAP approach

Cavity/TAP approach

- Stat. mech. evaluation of the LASSO estimator for a **single sample** of \mathbf{y}, A

Gibbs free energy

$$\begin{aligned}\Phi(\mathbf{m} | \mathbf{y}, A; \lambda) &= \max_{\mathbf{h}} \left[\mathbf{h} \cdot \mathbf{m} - \lim_{\beta \rightarrow +\infty} \frac{1}{\beta} \ln \int e^{-\beta H(\mathbf{x} | \mathbf{y}, A; \lambda) + \beta \mathbf{h} \cdot \mathbf{x}} d\mathbf{x} \right] \\ &= \max_{\mathbf{h}} \left[\mathbf{h} \cdot \mathbf{m} + \min_{\mathbf{x}} \left[\frac{1}{2} \|\mathbf{y} - A\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1 - \mathbf{h} \cdot \mathbf{x} \right] \right] \\ &\quad \left(= \frac{1}{2} \|\mathbf{y} - A\mathbf{m}\|_2^2 + \lambda \|\mathbf{m}\|_1 \right)\end{aligned}$$

$$\hat{\mathbf{x}}^{\text{LASSO}}(\mathbf{y}, A; \lambda) = \arg \min_{\mathbf{m}} \{ \Phi(\mathbf{m} | \mathbf{y}, A; \lambda) \}$$

Cavity/TAP approach

- Adaptive TAP/EC approach
 - Oppor and Winther (2001,2006), YK and Vehkaperä (2014), ...

Generalized free energy

$$\tilde{\Phi}(\mathbf{m} | \mathbf{y}, A; \lambda, l) = \max_{\mathbf{h}} \left[\mathbf{h} \cdot \mathbf{m} + \min_{\mathbf{x}} \left[\frac{l}{2} \|\mathbf{y} - A\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1 - \mathbf{h} \cdot \mathbf{x} \right] \right]$$

$$\Phi = \tilde{\Phi}(l=1) = \int_0^1 \frac{\partial \tilde{\Phi}(l)}{\partial l} dl + \tilde{\Phi}(l=0)$$

$$\approx \max_{\mathbf{h}} \left[\frac{1}{2} \|\mathbf{y} - A\mathbf{m}\|_2^2 - \frac{\Lambda \|\mathbf{m}\|_2^2}{2} + \mathbf{h} \cdot \mathbf{m} - \frac{1}{2\Lambda} \sum_{i=1}^N (|h_i| - \lambda)^2 \Theta(|h_i| - \lambda) \right]$$

Macroscopic moment matching
+ Gaussian approximation

$$\left(\Lambda = 2G'(-\chi; A^\top A), \chi = \frac{1}{N\Lambda} \sum_{i=1}^N \Theta(|h_i| - \lambda) \right)$$

Statistical property of matrix ensemble is summarized in G-func.

Adaptive TAP equation

- Extremum condition of the free energy offers

$$\begin{cases} \mathbf{h} = \Lambda \mathbf{m} + A^\top (\mathbf{y} - A \mathbf{m}) \\ m_i = \frac{h_i - \lambda \operatorname{sgn}(h_i)}{\Lambda} \Theta(|h_i| - \lambda) \\ \Lambda = 2G'(-\chi; A^\top A) \\ \chi = \frac{1}{N\Lambda} \sum_{i=1}^N \Theta(|h_i| - \lambda) \end{cases}$$

Onsager reaction coeff

- We can evaluate the Gaussian fields $\{h_i\}$ from these equations.

But, how can we evaluate macroscopic variables $\hat{Q}, \hat{m}, \hat{\chi}$, which are averaged quantities w.r.t. A, \mathbf{y} from a single sample?

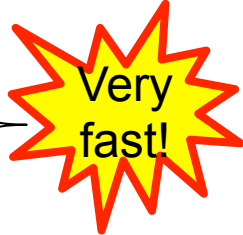
Self-averaging property

- For $N \gg 1$, the self-averaging property indicates that the macroscopic variables $\hat{Q}, \hat{m}, \hat{\chi}$ can be evaluated from a single sample of A, \mathbf{y} .
- Concretely, the replica/cavity equivalence for $N \gg 1$ provides the following correspondence:

Replica (averaged)	Cavity/TAP (single sample)
$\hat{Q} = \hat{m}$	$\Lambda (= 2G'(-\chi))$
$\hat{\chi}$	$\frac{G''(-\chi) \ \mathbf{y} - A\hat{\mathbf{x}}^{\text{LASSO}}\ _2^2}{M(G'(-\chi) - \chi G''(-\chi))} + \frac{-G''(-\chi) + 2G'(-\chi)}{G'(-\chi) - \chi G''(-\chi)} \sigma^2$

- We need to estimate the noise variance σ^2 from data.

Main result

1. $\hat{\mathbf{x}}^{\text{LASSO}} = \arg \min_x \left[\frac{1}{2} \|\mathbf{y} - A\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\| \right]$ Any algorithm:
VAMP, LARS, CD, ...
 2. Solve $\Lambda = 2G'(-\chi; A^\top A)$, $\chi = \frac{1}{\Lambda N} \# \{i \mid x_i^{\text{LASSO}} \neq 0\}$
 3. $\mathbf{h} = \Lambda \hat{\mathbf{x}}^{\text{LASSO}} + A^\top (\mathbf{y} - A \hat{\mathbf{x}}^{\text{LASSO}})$
 4. $\hat{\chi} = \frac{G''(-\chi) \|\mathbf{y} - A \hat{\mathbf{x}}^{\text{LASSO}}\|_2^2}{M (G'(-\chi) - \chi G''(-\chi))} + \frac{-G''(-\chi) + 2G'(-\chi)}{G'(-\chi) - \chi G''(-\chi)} \sigma^2$
 5. $\hat{\mathbf{x}}^{\text{debiased}} = \frac{\mathbf{h}}{\Lambda} \sim \mathcal{N}\left(\mathbf{x}_0, \frac{\hat{\chi}}{\Lambda^2}\right)$
- 

As for the noise variance σ^2 , a cross-validation-based estimator (Reid et al (2013)) empirically exhibits good performance.

Experimental set up

We examined the utility of the developed debiasing method by application to the following two matrix ensembles of $M = \alpha N$ ($\alpha = 0.5$).

1. i.i.d. Gaussian ensemble

- $A_{ij} \sim \mathcal{N}(0, N^{-1})$
- AED: Marchenko-Pastur distribution

$$\mu_{A^\top A}(s) = (1 - \alpha)\delta(s) + \frac{\alpha}{2\pi} \frac{\sqrt{\left[\left(1 + \sqrt{\alpha}\right)^2 - s \right] \left[s - \left(1 - \sqrt{\alpha}\right)^2 \right]}}{s}$$

2. Random DCT ensemble

- Random sampling of $M = \alpha N$ rows from $N \times N$ DCT matrix
- AED: $\mu_{A^\top A}(s) = (1 - \alpha)\delta(s) + \alpha\delta(s - 1)$
- Not rotationally invariant!!! But, the developed technique is applicable (Cakmak and Oppor (2018))

Experimental set up

Other settings:

- Sparse signal: Bernoulli-Gaussian dist.

$$x_{0,i} \sim_{\text{i.i.d.}} (1 - \rho)\delta(x) + \rho\mathcal{N}(0,1)$$

- Noise estimator: CV-based (Reid et al (2013))

$$\hat{\sigma}^2 \equiv \frac{1}{M - \#\{i | \hat{x}_i^{\text{LASSO}} \neq 0\}} \left\| \mathbf{y} - \mathbf{A} \hat{\mathbf{x}}^{\text{LASSO}}(\mathbf{y}, \mathbf{A}; \hat{\lambda}) \right\|_2^2$$

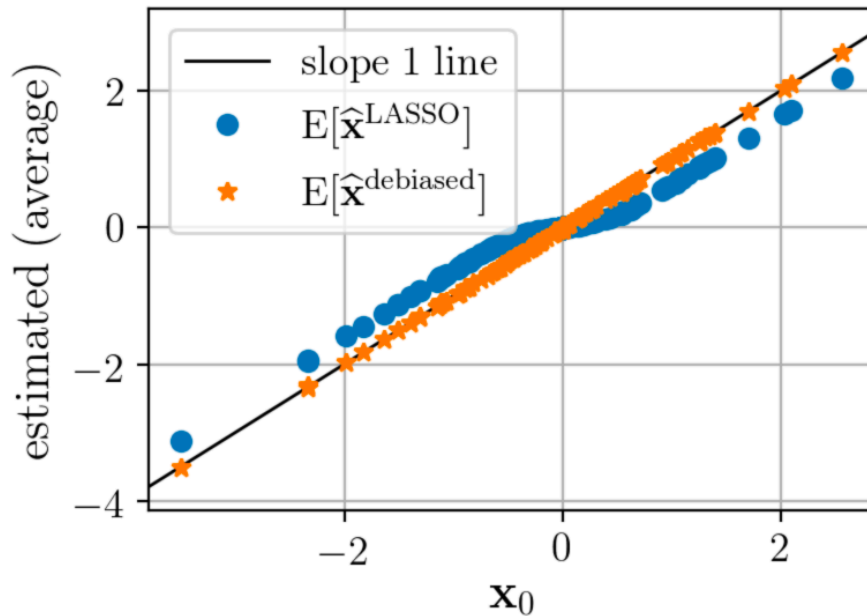
$\hat{\lambda}$: Determined so that CV error is minimized
Approximate CV formula (Obuchi and YK (2016))

- System parameters:

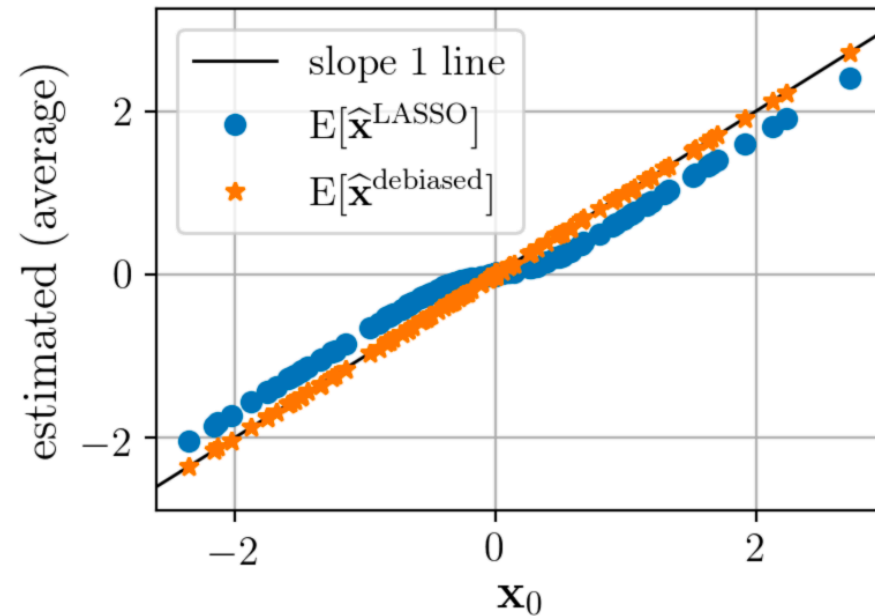
$$N = 1000, \alpha = M/N = 0.5, \sigma^2 = 0.2, \rho = 0.1$$

Debiasing

i.i.d. Gaussian

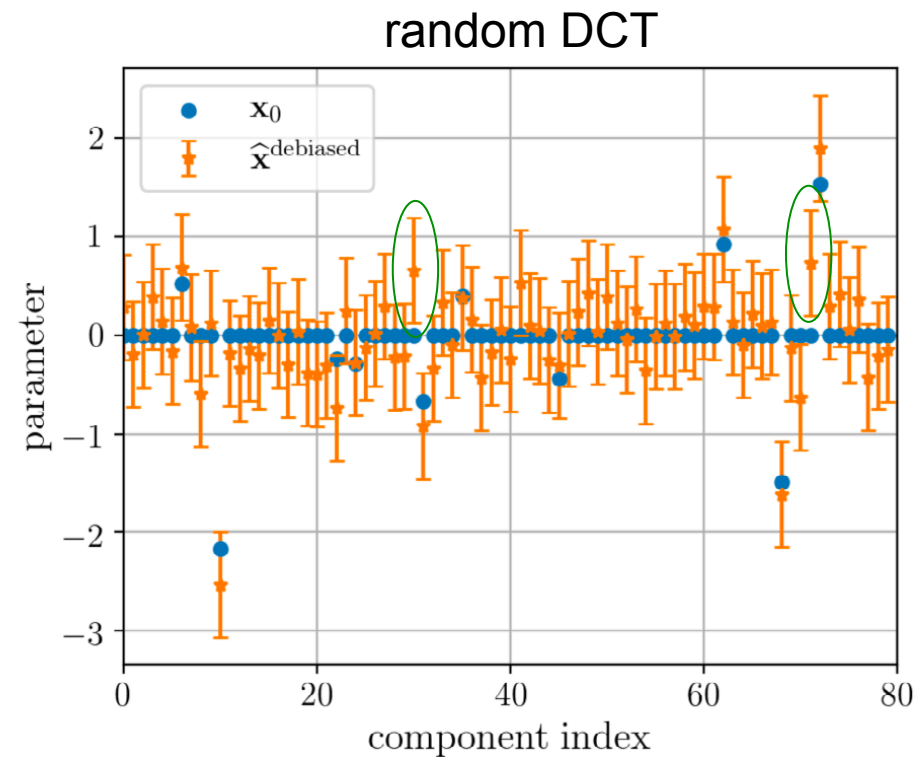
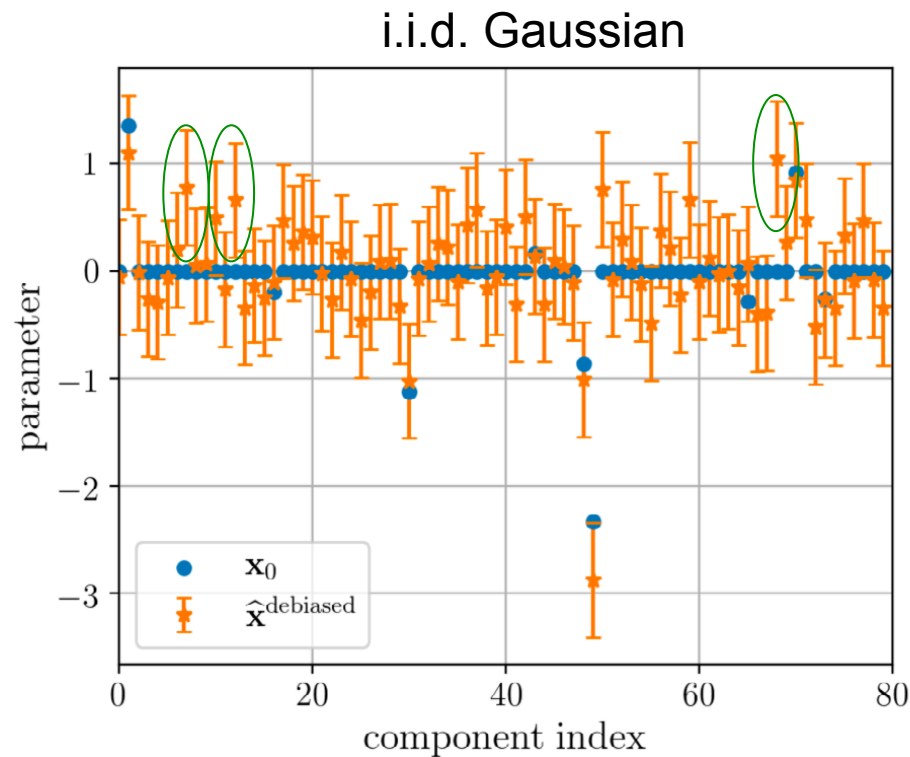


random DCT



Each point represents average over 1000 sets of A and y

Confidence interval (CI)



Error bars indicate 95% confidence intervals evaluated from a single data sample A, \mathbf{y} .

$$\left[\hat{x}_i^{\text{debiased}} - 1.96 \frac{\sqrt{\hat{\chi}}}{\Lambda}, \hat{x}_i^{\text{debiased}} + 1.96 \frac{\sqrt{\hat{\chi}}}{\Lambda} \right]$$

The CIs actually cover true parameters $x_{0,i}$ with a probability of about 95%.

Statistical testing

$$\left\{ \begin{array}{ll} \text{Null hypothesis} & H_{0,i} : x_{0,i} = 0 \\ \text{Alternative hypothesis} & H_{1,i} : x_{0,i} \neq 0 \end{array} \right.$$

#Relevant for “variable selection”

Testing procedure

$$\hat{T}_i(y, A; \lambda) = \begin{cases} 1, & \hat{x}_i^{\text{debiased}} \in \text{RR}(\alpha_{\text{sig}}) \text{ (reject)} \\ 0, & \text{otherwise (accept)} \end{cases}$$

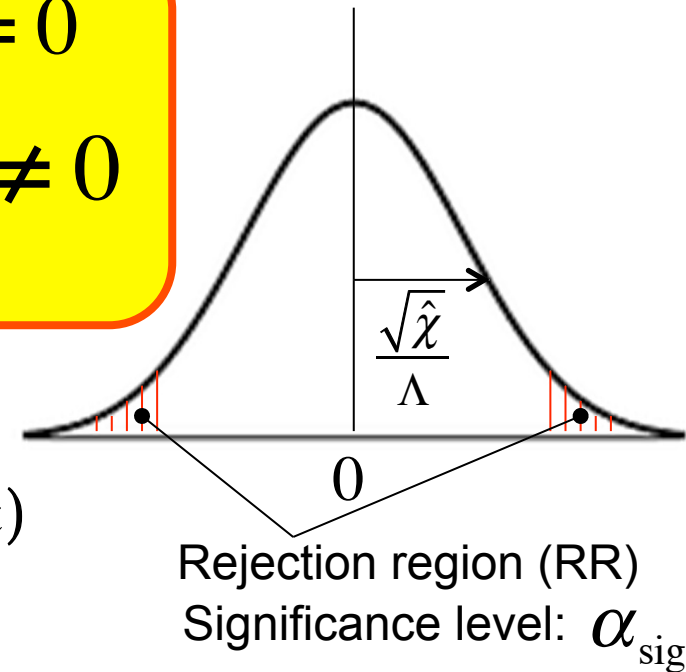
Performance measure

$$\text{FPR} = \frac{\#\{\hat{T}_i = 1 \text{ and } x_{0,i} = 0\}}{\#\{x_{0,i} = 0\}} \equiv \alpha_{\text{sig}}$$

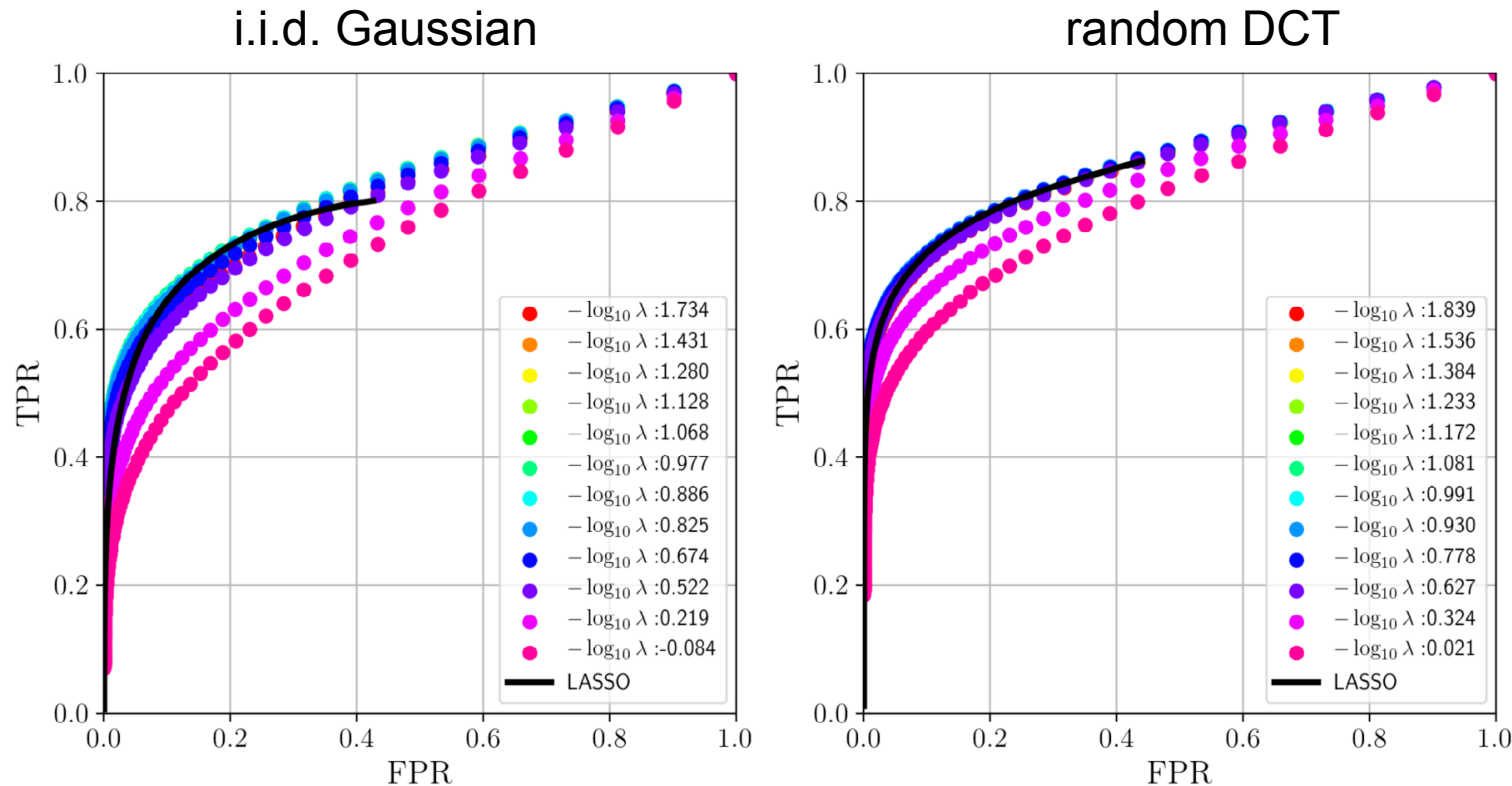
(control parameter)

$$\text{TPR} = \frac{\#\{\hat{T}_i = 1 \text{ and } x_{0,i} = 1\}}{\#\{x_{0,i} = 1\}}$$

(desired to be maximized)



Statistical testing



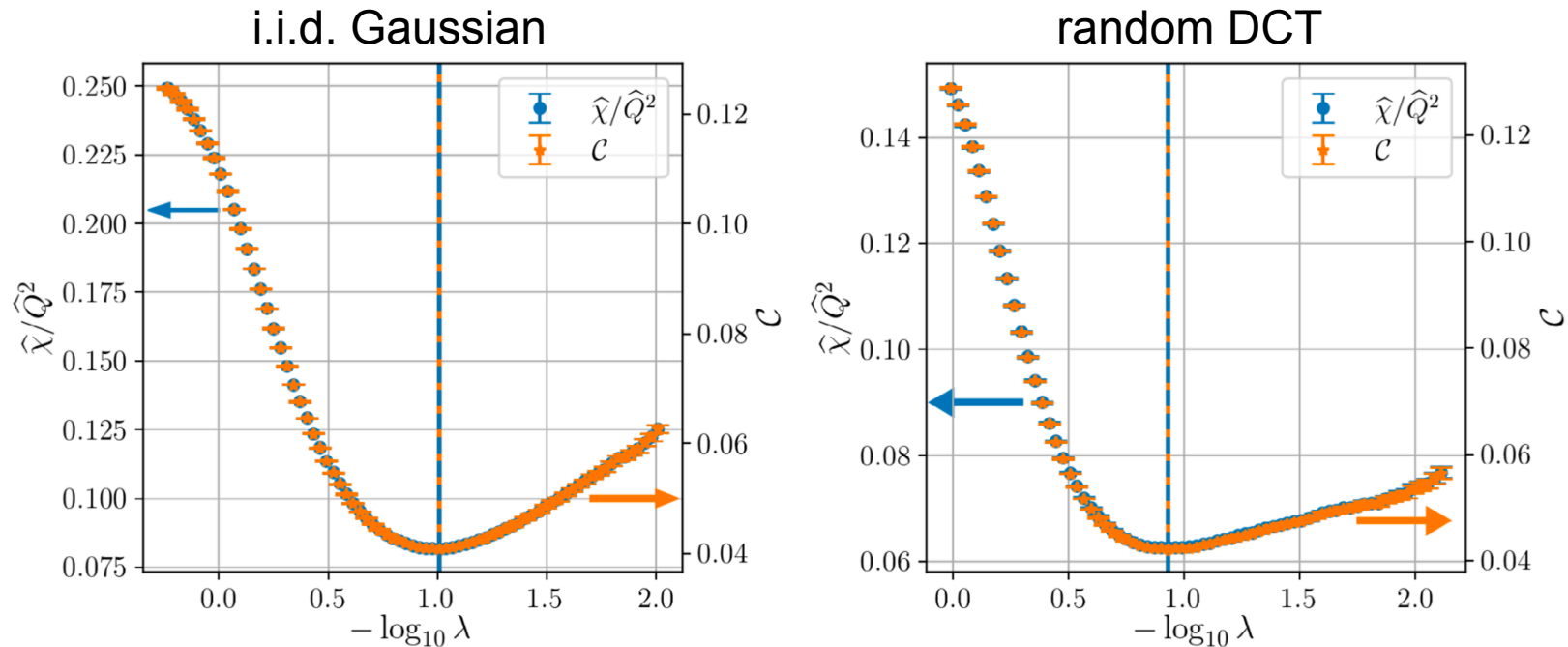
ROC curve (FPR-TPR curve) is maximized when variance $\frac{\hat{\chi}}{\Lambda^2}$ is minimized.

The resulting detection performance is (slightly) better than that achieved by naively employing LASSO varying l_1 strength λ .

CI size versus CV error

- Minimization of cross-validation (CV) error is one of the standard approaches for determining λ in LASSO.
- On the other hand, maximizing the estimation accuracy by minimizing the size of CI would offer an alternative criterion for selecting λ .
- Which criterion is better?

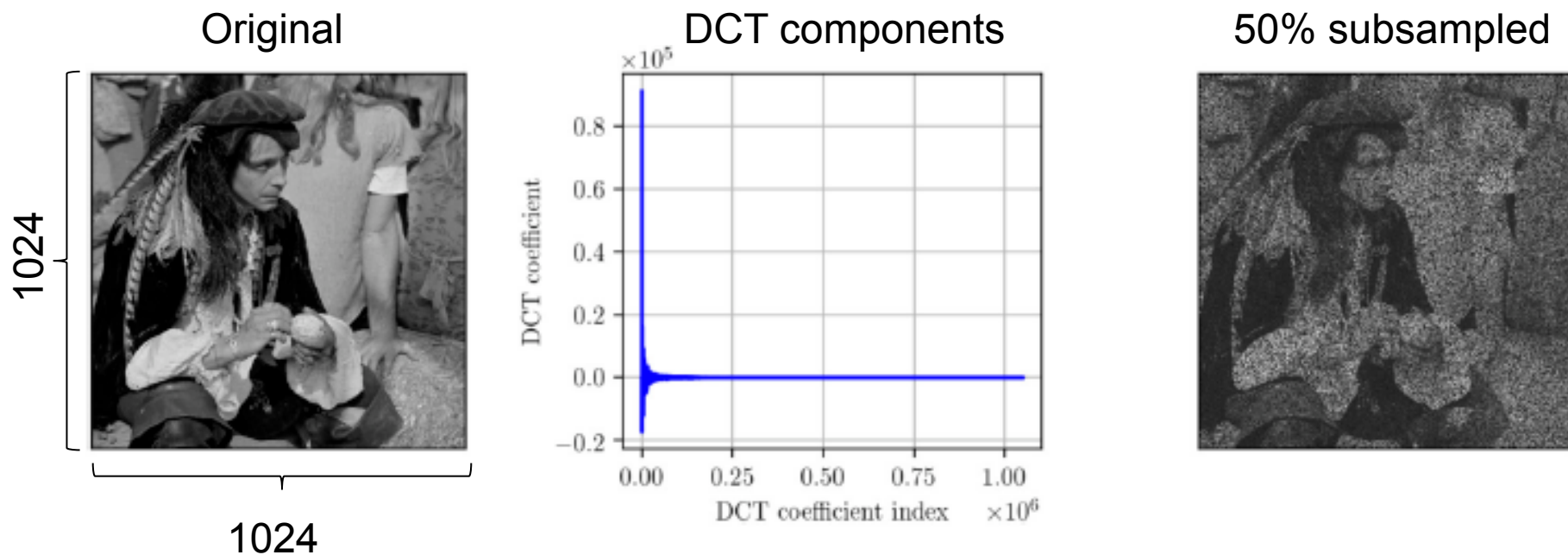
CI size versus CV error



- For i.i.d. Gaussian and random DCT ensembles, CI size and CV error are linearly related, so that the two criteria always lead to the same result!
- The relation between CI size and CV error does not necessarily hold for other ensembles.
- However, as far as we examined, these two criteria are always minimized at the same value of λ when the matrix ensemble is rotationally invariant.

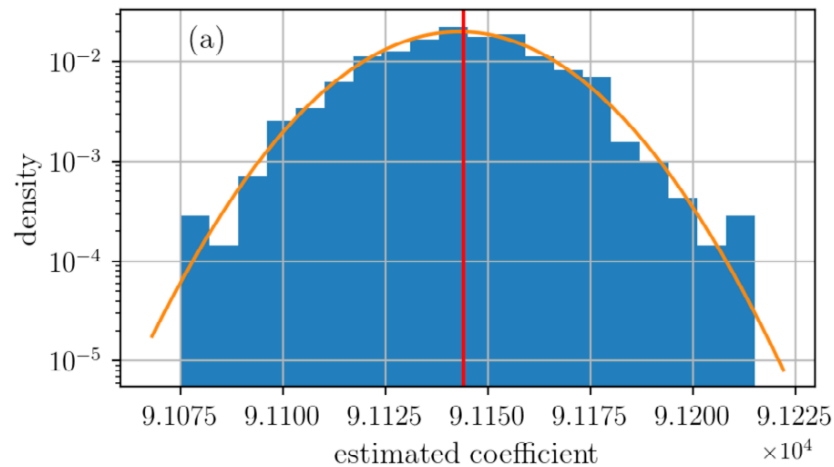
Demonstration on a real world data set

- Task: Estimation of Fourier (DCT) components from sub-sampled real world signals
 - Many demands in spectral analysis
 - Here, we handle a megapixel image for ease of visual understanding

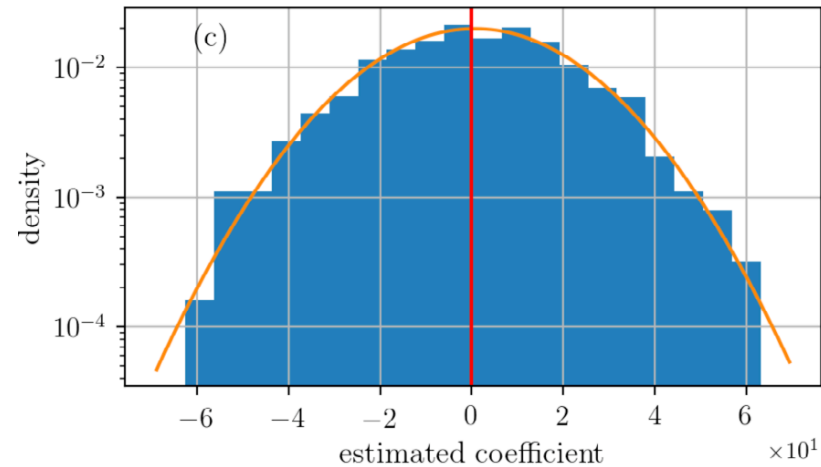


Demonstration on a real world data set

Largest component



Smallest component



Histograms: evaluated from 1000 samples of random DCTs and noises
Curves: evaluated from a single sample of random DCT and noise

Summary

- We developed a method for de-biasing and uncertainty estimation in LASSO in the case of rotationally invariant design/observation matrices.
 - Earlier studies: fixed matrix, vanishing sparsity ratio
 - Javanmard and Montanari (2014), van de Geer et al (2014), Zhang and Zhang (2014)
 - Ours: random matrix, finite sparsity ratio
- Although we here focused on LASSO, the extension to other regularized estimation is straightforward.

Summary

- Advantages of the method
 - Computationally feasible
 - Ability for constructing confidence interval
 - Better performance for statistical testing (signal detection) than naïve LASSO
- Possible unsatisfaction
 - Limited applicability to a special class of random matrices although random Fourier ensembles may potentially have wide application domains
- Semi-analytic resampling (bootstrap) may be useful for resolving this issue.
 - Obuchi and YK, arXiv:1802.10254 => ``pair bootstrap''
 - Takahashi and YK, in preparation => ``residual bootstrap''

References

- T. Takahashi and YK,
“A statistical mechanics approach to de-biasing and
uncertainty estimation in LASSO for random
measurements”,
J. Stat. Mech. (2018) 073405 (open access)
 - Demo:
https://github.com/takashi-takahashi/debiasing_lasso_demo