

Bootstapping Lasso Estimators

Alberto Quaini

June 2017

Presentation

Presentation of:

“Bootstrapping Lasso Estimator” – A. Chatterjee, S. N. Lahiri [2011], JASA.

for PEF UNISG course:

“Resampling methods and forecasting” – L. Camponovo

Additional literature:

+ Chatterjee, Lahiri [2010]

Outline

1. Introduction
2. The Modified Bootstrap method
 - ▶ Background and motivation
 - ▶ A Modified Bootstrap method
3. Bootstrapping the Lasso estimator
 - ▶ Consistency and the distributional approximation
 - ▶ Bootstrap bias and variance estimation
4. Bootstrapping the Adaptive Lasso estimator
 - ▶ A residual Bootstrap method for the Adaptive Lasso estimator
 - ▶ Main results

Outline

5. Data-based choice of the regularization parameter

- ▶ The optimal regularization parameter
- ▶ Data-based selection of the optimal regularization parameter
- ▶ Jackknife-After-Bootstrap based choice of the regularization parameter

6. Numerical results

- ▶ Choice of optimal penalization and thresholding parameters
- ▶ Coverage accuracy of confidence regions
- ▶ Variance estimation

Introduction

Linear regression model with iid errors:

$$y_i = x_i^T \beta + \epsilon_i, \quad i = 1, \dots, n \quad (1)$$

Lasso estimator:

$$\hat{\beta}_n = \underset{u \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^n (y_i - x_i^T u)^2 + \lambda_n \sum_{j=1}^p |u_j| \quad (2)$$

- ▶ estimation and variable selection method (Tibshirani [1996])
- ▶ computationally feasible (Friedman et Al. [2007])
- ▶ model consistency (Wainwright [2006], Zhao and Yu [2006] and Zou [2006])
- ▶ estimation consistency (Knight and Fu [2000])

Introduction

Problems:

1. Consistency

- ▶ Knight and Fu [2000] show that the limiting distribution of the Lasso estimator is complex in the situation of sparse underlying parameter vector
- ▶ in practice alternative approximations are needed to carry on inference for the Lasso
- ▶ Chatterjee and Lahiri [2010] show that the Bootstrapped Lasso estimator based on the residual Bootstrap method is inconsistent whenever at least one component of the parameter vector is zero

2. Confidence intervals and testing

- ▶ proposals of Tibshirani [1996] and Osborne et Al. [2000] have the drawback of considering the Lasso an approximately linear transformation
- ▶ proposals of Tibshirani [1996], Fan and Li [2001] and Fan and Peng [2004] only provide CI for underlying non-zero parameters

Introduction

Results and proposals in Chatterjee and Lahiri [2011]:

1. Consistency

- ▶ construct a suitable modification to the residual-based Bootstrap
- ▶ show consistency under mild regularity conditions even when some of the underlying parameters are zero

2. Confidence interval and testing

- ▶ the modified Bootstrap method provides consistent estimate of the variance of the Lasso estimator for both zero and non-zero parameter components

Introduction

3. choice of the regularization parameter λ_n

- ▶ accuracy of the Lasso critically depends on the regularization parameter
- ▶ the modified Bootstrap is consistent for the MSE of the Lasso, hence it can be used for selecting λ_n

4. Adaptive Lasso estimator (Zou [2006])

- ▶ adaptive weights are used for penalizing different coefficients in the L_1 penalty
- ▶ it enjoys the oracle property, i.e. performs as well as if the true underlying model were given in advance
- ▶ Chatterjee and Lahiri [2011] show that the simple residual Bootstrap can consistently estimate the distribution of the adaptive Lasso estimator

The Modified Bootstrap method

background and motivation

The **residual Bootstrap** method (standard in linear regression setting with nonrandom x_i , see Efron [1979], Freedman [1981]) proceeds as follows in the context of the Lasso (Knight and Fu [2000]):

1. Consider the set of centered residuals
 $E = \{e_i = \tilde{e}_i - \bar{e}, \text{ for } i = 1, \dots, n\}$, where $\bar{e} = n^{-1} \sum_i \tilde{e}_i$ and \tilde{e}_i 's are the residuals of the Lasso fit on the original sample.
2. Construct B bootstrap samples of size n selecting with replacement from E : $E_b^* = \{e_{i,b}^* : i = 1, \dots, n\}$ and compute $y_{i,b}^* = x_i^T \hat{\beta}_n + e_{i,b}^*$, for $i = 1, \dots, n$, and $b = 1, \dots, B$, where $\hat{\beta}_n$ is the Lasso estimator for the original sample.

The Modified Bootstrap method

3. Compute the bootstrap version of $T_n = n^{1/2}(\hat{\beta}_n - \beta)$, i.e. $T_n^* = n^{1/2}(\hat{\beta}_{n,b}^* - \hat{\beta}_n)$, where $\hat{\beta}_{n,b}^*$ is the Lasso estimator for bootstrap sample b .
4. The residual Bootstrap estimator of the distribution G_n of T_n is $\hat{G}_n(B) = P_*(T_n^* \in B)$, where $B \in \mathcal{B}(R^p)$ and P_* is the probability of T_n^* given errors ϵ_i 's.

Chatterjee and Lahiri [2010] show that:

- ▶ the estimators of the zero parameters fail to capture the target sign value, which is zero
- ▶ because of that, \hat{G}_n , instead of converging to the deterministic limit of G_n converges weakly to a random probability measure
- ▶ i.e. it fails to provide a valid approximation to G_n

The Modified Bootstrap method

A Modified Bootstrap method

Objective: capture the signs of the parameters, especially the zero components, with probability tending to 1, as the sample size n goes to infinity.

Idea: force components of the Lasso estimator $\hat{\beta}_n$ to be exactly zero whenever they are close to zero using the fact that the Lasso estimator is root- n consistent.

To this end:

1. Form a sequence $\{a_n\}$ of real numbers such that $a_n + (n^{-1/2} \log(n))a_n^{-1} \rightarrow 0$ asymptotically.
2. Threshold the components of the Lasso estimator β_n at a_n , and define the modified Lasso estimator

$$\tilde{\beta}_{n,j} = \beta_{n,j} \mathbb{1}(\beta_{n,j} \geq a_n), \text{ for } j = 1, \dots, p. \quad (3)$$

The Modified Bootstrap method

Note that with probability tending to 1 (as $n \rightarrow \infty$):

- ▶ $|\hat{\beta}_{n,j}| = |\beta_j| + O(n^{-1/2}) > |\beta_j|/2 \geq a_n$, for a nonzero component β_j
- ▶ $|\hat{\beta}_{n,j}| = |\beta_j| + O(n^{-1/2}) = O(n^{-1/2}) \in [-a_n, a_n]$, for a zero component β_j

Then proceed as before:

3. Consider the set of centered residuals
 $R = \{r_i = \tilde{r}_i - \bar{r}, \text{ for } i = 1, \dots, n\}$, where $\bar{r} = n^{-1} \sum_i \tilde{r}_i$ and \tilde{r}_i 's are the residuals of the modified Lasso fit on the original sample.
4. Construct B bootstrap samples of size n selecting with replacement from R : $R_b^{**} = \{r_{i,b}^{**} : i = 1, \dots, n\}$ and compute $y_{i,b}^{**} = x_i^T \tilde{\beta}_n + r_{i,b}^{**}$, for $i = 1, \dots, n$, and $b = 1, \dots, B$, where $\tilde{\beta}_n$ is the modified Lasso estimator for the original sample.

The Modified Bootstrap method

5. Compute the bootstrap version of $T_n = n^{1/2}(\hat{\beta}_n - \beta)$, i.e. $T_n^{**} = n^{1/2}(\hat{\beta}_{n,b}^{**} - \tilde{\beta}_n)$, where $\hat{\beta}_{n,b}^{**}$ is the Lasso estimator (not the modified one) for bootstrap sample b .
6. The residual Bootstrap estimator of the distribution G_n of T_n is $\tilde{G}_n(B) = P_{**}(T_n^{**} \in B)$, where $B \in \mathcal{B}(R^p)$ and P_{**} is the probability of T_n^{**} given errors ϵ_i 's.

Remarks:

- ▶ Centering the residuals ensures the Bootstrap analogue of the condition $E[e_i] = 0$
- ▶ A rescaling factor $(1 - p/n) - 1/2$ is sometimes used in the construction of the residuals (see Efron [1982]) to improve finite sample accuracy
- ▶ It is possible to replace $\hat{\beta}_n$ by any other \sqrt{n} -consistent estimator of β , e.g. least squares

Bootstrapping the Lasso estimator

Consistency and the distributional approximation

Theorem 1: Consistency of Modified Bootstrap

Assume:

- ▶ (C1) $n^{-1} \sum_i x_i x_i^T \rightarrow C$, p.d. matrix. Furthermore $n^{-1} \sum_i \|x_i\|^3 \rightarrow O(1)$.
- ▶ (C2) $\lambda_n n^{-1/2} \rightarrow \lambda_0 \geq 0$.
- ▶ (C3) errors ϵ_i 's are iid with $E[\epsilon_i] = 0$ and $\text{VAR}[\epsilon_i] = \sigma^2 < \infty$.

Then:

$$\mathcal{P}(\tilde{G}_n, G_n) \rightarrow 0, \text{ as } n \rightarrow \infty, \text{ with probability 1,}$$

where $\mathcal{P}(\cdot, \cdot)$ denotes the Prohorov probability metric.

Bootstrapping the Lasso estimator

Remarks:

- ▶ Chatterjee and Lahiri [2010] shows that under the same set of regularity assumptions, if β has at least one zero component and if \hat{G}_n is the residual bootstrap estimate of G_n , then

$$\mathcal{P}(\hat{G}_n, G_n) \not\rightarrow 0, \text{ in probability, as } n \rightarrow \infty$$

- ▶ Theorem 1 states strong consistency of the modified Bootstrap distribution estimator
- ▶ From Theorem 1 it follows that the modified bootstrap method can be used to approximate the distribution of the Lasso estimator T_n for any $\beta \in R^p$. Hence, it can be used to construct valid large sample confidence set estimators of β

Bootstrapping the Lasso estimator

Definitions:

- ▶ let $t(\alpha)$ denote the $\alpha \in (0, 1)$ quantile of $\|T_\infty\|$, where T_∞ denotes the limiting random vector such that $T_n \rightarrow T_\infty$ and has distribution G_∞ .
- ▶ let $\hat{t}_n(\alpha)$ denote the $\alpha \in (0, 1)$ quantile of the bootstrap distribution of $\|T_n^{**}\|$. Then the set

$$I_{n,\alpha} \equiv \{t \in R^p : \|t - \hat{\beta}_n\| \leq n^{-1/2} \hat{t}_n(\alpha)\}.$$

Bootstrapping the Lasso estimator

Corollary 1: Modified Bootstrap Confidence Interval

Assume (C1), (C2) and (C3) hold. Then:

- i if $\alpha \in (0, 1)$ is such that $P(\|T_\infty\| \leq t(\alpha) + \eta) > \alpha, \forall \eta > 0$, then for all $\beta \in R^p$:

$$P(\beta \in I_{n,\alpha}) \rightarrow \alpha, \quad \text{as } n \rightarrow \infty \quad (4)$$

- ii if there is at least one nonzero component of β , then (4) holds for all $\alpha \in (0, 1)$.

- ▶ Corollary 1 justifies the use of the modified Bootstrap method to construct valid large sample confidence regions for β
- ▶ Corollary 1 can also be used to test the null hypothesis $H_0 : \beta_j = 0$ for all $j \in J$ for a given $J \subset \{1, \dots, p\}$

Bootstrapping the Lasso estimator

Remarks:

- ▶ Leeb and Pötscher [2006, 2008] and Pötscher and Schneider [2009] show that it is impossible to consistently estimate the distribution function of the Lasso estimator in a uniform sense
- ▶ Problems arise especially when some underlying nonzero parameters get close to zero as n gets large
- ▶ Theorem 1 provides a method to obtain a consistent estimator in case the underlying parameters are fixed
- ▶ Andrews and Guggenberger [2009] show that uniform consistency is not necessary for producing uniformly valid confidence intervals
- ▶ Corollary 1 asserts that the modified Bootstrap method can control the asymptotic size of confidence intervals, however it is not clear if the latter are uniformly valid in the parameter values

Bootstrapping the Lasso estimator

Bootstrap bias and variance estimation

Theorem 2: Bias and Variance Consistency

Assume (C1), (C2) and (C3) hold. Then with probability 1:

$$E_*[T_n^{**}] \rightarrow E[T_\infty] \quad \text{and} \quad (5)$$

$$(VAR_*[T_n^{**}])_{p \times p} \rightarrow (VAR_*[T_\infty])_{p \times p} \quad (6)$$

- ▶ for $\lambda_0 \neq 0$ in assumption (C2), T_n may be asymptotically biased and standard MSE estimation methods are unreliable
- ▶ The modified Bootstrap method produces strongly consistent estimators of the asymptotic bias and variance of T_n
- ▶ hence, Theorem 2 allows to estimate the MSE of a Lasso estimate and quantify the associated uncertainty for all values of β

Boostrapping the Adaptive Lasso estimator

The **adaptive Lasso** estimator (Zou [2006]):

$$\check{\beta}_n = \operatorname{argmin}_{u \in R^p} \sum_{i=1}^n (y_i - x_i^T u)^2 + \lambda_n \sum_{j=1}^p \frac{|u_j|}{|\bar{\beta}_{j,n}|^\gamma}, \quad (7)$$

where $\bar{\beta}_n$ denote an initial consistent estimator of β (e.g. least squares), $\lambda_n \geq 0$ is the penalty and $\gamma > 0$.

Oracle property (Zou [2006]):

$$P(B_n = A) \rightarrow 1, \quad \text{as } n \rightarrow 1 \quad (8)$$

$$\sqrt{n}(\check{\beta}_n^{nz} - \beta^{nz}) \rightarrow N(0, \sigma^2 C_{nz}) \quad (9)$$

where $A = \{j : \beta_j = 0\}$, $B_n = \{j : \check{\beta}_{j,n} = 0\}$ and $\sigma^2 C_{nz}$ is the var-cov matrix between nonzero estimated and underlying parameters

Boostrapping the Adaptive Lasso estimator

A residual Bootstrap method for Adaptive Lasso

The algorithm is similar to the residual Bootstrap described earlier with few adjustments:

- ▶ E becomes the set of centered residuals of the adaptive Lasso fit on the original sample
- ▶ for each bootstrap sample, the adaptive Lasso estimator becomes:

$$\check{\beta}_n^+ = \operatorname{argmin}_{u \in \mathbb{R}^p} \sum_{i=1}^n (y_i^+ - x_i^T u)^2 + \lambda_n \sum_{j=1}^p \frac{|u_j|}{|\bar{\beta}_{j,n}^+|^{\gamma}}, \quad (10)$$

where $y_i^+ = x_i^T \check{\beta}_n + \epsilon_i^+$, $i = 1, \dots, n$, ϵ_i^+ 's are bootstrapped from E and $\bar{\beta}_n^+$ is defined by replacing y_i 's with y_i^+ 's in the definition of $\bar{\beta}_n$

Boostrapping the Adaptive Lasso estimator

Remarks:

- ▶ the penalty of the adaptive Lasso incorporates a built-in soft-thresholding for the zero parameters
- ▶ Hence, the Bootstrap procedure just described does not need an initial truncation as it is the case for the Lasso

Denote:

- ▶ $\check{T}_n \equiv \sqrt{n}(\check{\beta}_n - \beta)$ with distribution H_n , where $H_n(x) = P(\check{T}_n \leq x), x \in R$
- ▶ $\check{T}_n^+ \equiv \sqrt{n}(\check{\beta}_n^+ - \check{\beta}_n)$ with distribution H_n^+ conditional on the ϵ_i 's

Boostrapping the Adaptive Lasso estimator

Theorem 3:

Assume (C1), and (C3) hold and suppose

$$\frac{\lambda_n}{\sqrt{n}} \rightarrow 0 \quad \text{and} \quad \lambda_n n^{(\gamma-1)/2} \rightarrow \infty. \quad (11)$$

Then,

$$\mathcal{P}(\hat{H}_n, H_n) \xrightarrow{P} 0, \quad \text{as } n \rightarrow \infty, \quad (12)$$

Remarks:

- ▶ (12) can be used to construct valid confidence intervals for β
- ▶ a corollary of the form of corollary 1 can be formulated for the adaptive Lasso residual bootstrap method

Boostrapping the Adaptive Lasso estimator

- ▶ estimation of the MSE of the adaptive Lasso estimator can be difficult
- ▶ adaptive Lasso residual Bootstrap method provides a consistent estimator of the MSE matrix of the scaled adaptive Lasso estimator $\check{\beta}_n$ given by $MSE[\check{T}_n] \equiv nE[(\check{\beta}_n - \beta)(\check{\beta}_n - \beta)^T]$

Corollary 3

Assume (C1) and (C2) hold, Then:

$$MSE_*(\check{T}_n^+) - MSE(\check{T}_n) \xrightarrow{P} 0, \quad \text{as } n \rightarrow \infty. \quad (13)$$

Data-based choice of the regularization parameter

The optimal regularization parameter

Remarks:

- ▶ it can be shown that the distribution of T_n depends on λ_n only through λ_0
- ▶ note that $MSE(\hat{\beta}_n)$ can be expressed as $n^{-1}E\|T_n\|^2$ and that $nMSE(\hat{\beta}_n)$ converges to the MSE of the limiting random variable T_∞
- ▶ The effect of the penalization by λ_n on the overall accuracy of $\hat{\beta}_n$ is reflected by its MSE
- ▶ for what comes next, consider the natural reparametrization $\lambda_n = \lambda_0 n^{1/2}$, $\lambda_0 \in [0, \infty)$

Data-based choice of the regularization parameter

Then define the **optimal penalization** parameter as

$$\lambda_0^{opt} \equiv \operatorname{argmin} \phi(\lambda_0) \quad (14)$$

where $\phi(\lambda_0) = E \|T_\infty\|^2$.

- ▶ Thus, choosing $\lambda_0 = \lambda_0^{opt}$ yields a Lasso estimator that minimizes the *MSE* in large samples.

Data-based choice of the regularization parameter

Data-based selection of the optimal regularization parameter

λ_0 can be estimated using the modified bootstrap:

- ▶ for any $\lambda_0 \in [0, \infty)$ and thresholding value $a \in (0, \infty)$ (as defined in (3)), rewrite $T_n^{**} \equiv T_n^{**}(\lambda_0, a)$
- ▶ the modified Bootstrap estimator of $\phi(\lambda_0)$ is

$$\hat{\phi}_n(\lambda_0, a) \equiv E_* \|T_n^{**}(\lambda_0, a)\|^2 \quad (15)$$

- ▶ by Theorem 1, $\hat{\phi}_n(\lambda_0, a)$ is a strongly consistent estimator of $\phi(\lambda_0)$

Therefore, an accurate estimator in large sample of λ_0^{opt} is the Bootstrap estimator:

$$\lambda_{0,n}^* = \operatorname{argmin}_{\lambda_0, a} \hat{\phi}_n(\lambda_0, a) \quad (16)$$

Data-based choice of the regularization parameter

Jackknife-After-Bootstrap based choice of the regularization parameter

The authors propose to estimate the thresholding parameter a via Jackknife-After-Bootstrap (JAB, see Efron [1992]). The basic idea of the method is:

- ▶ compute replicates of the Bootstrap estimator $\hat{\phi}_n(\lambda_0, a)$ by means of delete-1 JAB over the Bootstrap samples $\{T_n^{**}(b : \lambda_0, a) : b = 1, \dots, B\}$
- ▶ construct an estimate of the error of the Bootstrap estimator $\hat{\phi}_n(\lambda_0, a)$ based on the JAB replicates
- ▶ select a thresholding value a that minimize this estimate of error over all possible values of λ_0 and a

Numerical results

Choice of optimal penalization and thresholding parameters

For the numerical results regarding the choice of the regularization and thresholding parameters the authors consider:

- ▶ fixed-design matrix and errors generated by independent standard normal distribution with sample size $n = 250$
- ▶ $p = 10$ overall regression coefficients which 6 are nonzero organized in three Cases:
 - i large nonzero coefficients
 - ii some small nonzero coefficients
 - iii some nonzero coefficients shrinking as n gets large

Numerical results

Figures (reference) show the behaviour of the naive residual Bootstrap (NB) and the modified Bootstrap (MB) estimates of the optimal λ_0 in Case (i) and (ii), respectively.

Note:

- ▶ the vertical solid line represents the true optimal λ_0
- ▶ MB estimates are much better than the NB for Case (i)
- ▶ in Case (ii) the NB seem to performs better than the MB
- ▶ however if we stratify the performance of the MB by different choices of a we sometimes get much better results
- ▶ the MB is very sensible to the choice of a , as it influences the exclusion of true nonzero parameters

Numerical results

Figure 1

Numerical results

Figure 2

Numerical results

Regarding the choice of the threshold parameter a , the authors consider a grid of six values over which the optimum values of a are computed.

Numerical results are found in Table 1 in the paper.

The main results are:

- ▶ the JAB estimates of a are very good for Case (i) but deteriorates for the other two cases
- ▶ the choice of the grid is very important
- ▶ the JAB method is good in situations where β has distinct large nonzero coefficients

Numerical results

Coverage accuracy of confidence regions

The authors compare the finite sample performance of the confidence regions for β obtained by using the naive and the modified Bootstrap procedures.

Numerical results are found in Tables 2 – 4 in the paper.

The main findings are:

- ▶ both methods perform well in terms of achieving the desired nominal coverage rate for all three Cases
- ▶ the inconsistency of the NB does not have an effect on the coverage accuracy of the confidence regions
- ▶ however the inconsistency of the NB makes it an unreliable choice

Numerical results

Variance estimation

The authors then compare the NB, the MB and the m-out-of-n Bootstrap (MNB) (Hall, Lee and Park [2009]) and the sandwich formula (see Zou [2006] and Feng and Peng [2004]) for the estimating the variance of the Lasso estimator.

Numerical results are reported in Table 5 and 6 in the paper.

The main findings are:

- ▶ NB and MB have very comparable results
- ▶ MB has some advantages over the NB when the underlying value of the parameter is zero
- ▶ MNB has very poor performances, which gradually improve as m (the Bootstrap sample size) increases
- ▶ the sandwich estimator has very good performances, but they can be tremendously impacted by extreme observations in the data