

Asymptotic oracle properties of SCAD-penalized least squares estimators

San-Teng Huang

National Dong Hwa University

2018/09/25

Outline

- 1 Introduction
- 2 SCAD-penalized Least Squares
- 3 Asymptotic properties of the LS-SCAD estimator
- 4 Conclusion
- 5 Reference

Introduction

Model

- Consider a linear model:

$$y_i = \beta_0 + \tilde{\mathbf{x}}_i \boldsymbol{\beta} + \varepsilon_i, \quad i = 1, 2, \dots, n \quad (1)$$

where $\boldsymbol{\beta}$ is $p_n \times 1$ vector,

$\tilde{\mathbf{x}}_i$ is $1 \times p_n$ row vector of \mathbf{X} ,

ε_i *i.i.d* with mean 0, variance σ^2 , for $i = 1, \dots, n$.

- For simplicity, assume $\beta_0 = 0$.

Of interest

- Assume $p_n \rightarrow \infty$ as $n \rightarrow \infty$.

In biomedical studies investigating the relationship between phenotype and genomic.

- Assume only some parameters are non-zero.

the true $\beta_0 = (\beta_{01}, \beta_{02}, \dots, \beta_{0p_n}) = (\beta_{01}^T, \beta_{02}^T)$
 with $\beta_{01}^T = (\beta_{01}, \dots, \beta_{0k_n})$, $\beta_{02}^T = (0_{01}, \dots, 0_{0m_n})$,
 $p_n = k_n + m_n$.

- estimate β_0 where β_0 is sparse.

Variable selection

- Classical variable selection method: best subset selection, forward/backward-stepwise selection.

Disadvantage:

- (i) for high-dimensional data, computation is not feasible.
(p is large \Rightarrow lots of combinations)
- (ii) discrete process: take some subsets of full model and estimate coefficients. the variables are either removed or retained.

Variable selection

- Penalized method: SCAD, LASSO, ridge regression.
 - (i) It's more computationally feasible for high-dimensional data.
 - (ii) continuous process: fit the model with all variables but regularizing the estimated coefficients. achieve variable selection and estimations simultaneously.

SCAD-penalized Least Squares(LS-SCAD)

- To minimize

$$Q_n(\boldsymbol{\beta}; \lambda_n) = \sum_{i=1}^n (y_i - \tilde{\mathbf{x}}_i \boldsymbol{\beta})^2 + n \sum_{j=1}^{p_n} p_{\lambda_n}(|\beta_j|) \quad (2)$$

where $p_{\lambda_n}(|\beta_j|)$ is SCAD penalty function.

- The LS-SCAD estimator is

$$\hat{\boldsymbol{\beta}}_n = \arg \min_{\boldsymbol{\beta}} Q_n(\boldsymbol{\beta}; \lambda_n)$$

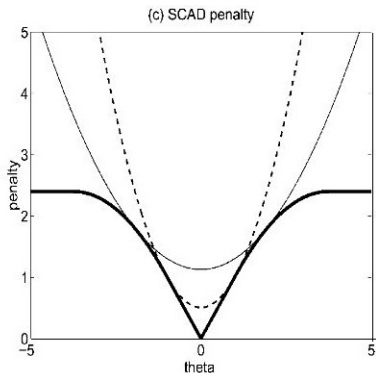
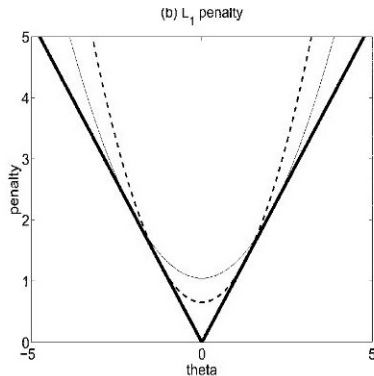
SCAD penalty function

- Smoothly Clipped Absolute Deviation (SCAD):

$$p_{\lambda_n}(|\beta|) = \begin{cases} \lambda_n |\beta| & , |\beta| \leq \lambda_n \\ -\left(\frac{|\beta|^2 - 2a\lambda_n|\beta| + \lambda_n^2}{2(a-1)}\right), & \lambda_n < |\beta| \leq a\lambda_n \\ \frac{(a+1)\lambda_n^2}{2} & , |\beta| > a\lambda_n \end{cases}$$

where $a = 3.7$ (Fan and Li, 2001).

Plot of $p_\lambda(|\theta|)$:



Asymptotic properties of the LS-SCAD estimator

Assumptions(for fixed covariates)

Let k_n be the number of nonzero coefficients, $\rho_{n,1}$ be the smallest eigenvalue of $\frac{1}{n}\mathbf{X}^T\mathbf{X}$. let $\pi_{n,k_n}, \omega_{n,m_n}$ be the largest eigenvalues of $\frac{1}{n}\mathbf{X}_1^T\mathbf{X}_1, \frac{1}{n}\mathbf{X}_2^T\mathbf{X}_2$, respectively.

(A0) (a) ε_i 's are *i.i.d* with mean 0 and variance σ^2 ;

(b) For $j \in \{1, \dots, p_n\}$, $\|\mathbf{x}_j\|^2 = n$. (\mathbf{x}_j column vector of \mathbf{X})

(A1) (a) $\lim_{n \rightarrow \infty} \sqrt{k_n} \lambda_n / \sqrt{\rho_{n,1}} = 0$;

(b) $\lim_{n \rightarrow \infty} \sqrt{p_n} / \sqrt{n \rho_{n,1}} = 0$.

(A2) (a) $\lim_{n \rightarrow \infty} \sqrt{k_n} \lambda_n / (\sqrt{\rho_{n,1}} \min_{1 \leq j \leq k_n} |\beta_j|) = 0$;

(b) $\lim_{n \rightarrow \infty} \sqrt{p_n} / (\sqrt{n \rho_{n,1}} \min_{1 \leq j \leq k_n} |\beta_j|) = 0$.

(A3) $\lim_{n \rightarrow \infty} \sqrt{\max(\pi_{n,k_n}, \omega_{n,m_n}) p_n} / (\sqrt{n \rho_{n,1}} \lambda_n) = 0$.

(A4) $\lim_{n \rightarrow \infty} \max_{1 \leq i \leq n} \tilde{\mathbf{x}}_{i1} (\sum_{i=1}^n \tilde{\mathbf{x}}_{i1}^T \tilde{\mathbf{x}}_{i1})^{-1} \tilde{\mathbf{x}}_{i1}^T = 0$. ($\tilde{\mathbf{x}}_{i1}$ row vector of \mathbf{X}_1)

Assumptions(for random covariates)

Let ρ_1 be the smallest eigenvalue of $E[\mathbf{x}\mathbf{x}^T]$. let π_{k_n}, ω_{m_n} be the largest eigenvalues of $E[\mathbf{X}_1^T \mathbf{X}_1]$, $E[\mathbf{X}_2^T \mathbf{X}_2]$, respectively.

(B0) $(\mathbf{x}_i^T, \varepsilon_i) = (X_{i1}, \dots, X_{ip_n}, \varepsilon_i)$, $i = 1, \dots, n$ are *i.i.d* with

(a) $E[X_{ij}] = 0$, $\text{Var}(X_{ij}) = 1$;

(b) $E[\varepsilon|\mathbf{X}] = 0$, $\text{Var}(\varepsilon|\mathbf{X}) = \sigma^2$.

(B1) (a) $\lim_{n \rightarrow \infty} p_n^2 / (n\rho_1^2) = 0$;

(b) $\lim_{n \rightarrow \infty} k_n \lambda_n^2 / \rho_1 = 0$.

(B2) (a) $\lim_{n \rightarrow \infty} \sqrt{p_n} / (\sqrt{n\rho_1} \min_{1 \leq j \leq k_n} |\beta_j|) = 0$;

(b) $\lim_{n \rightarrow \infty} \lambda_n \sqrt{k_n} / (\sqrt{\rho_1} \min_{1 \leq j \leq k_n} |\beta_j|) = 0$.

(B3) $\lim_{n \rightarrow \infty} \sqrt{\max(\pi_{k_n}, \omega_{m_n}) p_n} / (\sqrt{n\rho_1} \lambda_n) = 0$.

Consistency

Theorem 1 (Consistency in the fixed design setting)

Under (A0) – (A1),

$$\|\hat{\beta}_n - \beta\| \xrightarrow{P} 0 \quad \text{as } n \rightarrow \infty.$$

Theorem 2 (Consistency in the random design setting)

Suppose that there exists an absolute constant M_4 such that for all n , $\max_{1 \leq j \leq p_n} E[X_j^4] \leq M_4 < \infty$. Then under (B0) – (B1),

$$\|\hat{\beta}_n - \beta\| \xrightarrow{P} 0 \quad \text{as } n \rightarrow \infty.$$

Convergency rate

Lemma 1 (Convergency rate in the fixed design setting)

Under (A0) – (A2),

$$\|\hat{\beta}_n - \beta\| = O_P\left(\frac{\sqrt{p_n}}{\sqrt{n\rho_{n,1}}}\right).$$

Lemma 2 (Convergency rate in the random design setting)

Under (B0) – (B2),

$$\|\hat{\beta}_n - \beta\| = O_P\left(\frac{\sqrt{p_n}}{\sqrt{n\rho_1}}\right).$$

Variable selection

Theorem 3 (Variable selection in the fixed design setting)

Under (A0) – (A3),

$\hat{\beta}_{2n} = \mathbf{0}_{m_n}$ *with probability tending to 1.*

Theorem 4 (Variable selection in the random design setting)

Suppose there exists an absolute constant M such that

$\max_{1 \leq j \leq p_n} |X_j| \leq M < \infty$. *Then under (B0) – (B3),*

$\hat{\beta}_{2n} = \mathbf{0}_{m_n}$ *with probability tending to 1.*

Asymptotic normality

Let $\{\mathbf{A}_n, n = 1, 2, \dots\}$ be sequence of $d \times k_n$ matrices with full rank.

Theorem 5 (Asymptotic normality in the fixed design setting)

Under (A0) – (A4),

$$\sqrt{n}\Sigma_n^{-1/2}\mathbf{A}_n(\hat{\beta}_{1n} - \beta_1) \xrightarrow{D} N(\mathbf{0}_d, \mathbf{I}_d),$$

where $\Sigma_n = \sigma^2\mathbf{A}_n(\sum_{i=1}^n \tilde{\mathbf{x}}_{i1}^T \tilde{\mathbf{x}}_{i1} / n)^{-1}\mathbf{A}_n^T$.

Theorem 6 (Asymptotic normality in the random design setting)

Suppose that there exists an absolute constant M such that $\max_{1 \leq j \leq p_n} \|X_j\| \leq M < \infty$ and a σ_4 such that $E[\varepsilon^4 | X_{11}] \leq \sigma_4 < \infty$ for all n . Then under (B0) – (B3),

$$n^{-1/2}\Sigma_n^{-1/2}\mathbf{A}_n E^{-1/2}[\tilde{\mathbf{x}}_{i1}^T \tilde{\mathbf{x}}_{i1}] \sum_{i=1}^n \tilde{\mathbf{x}}_{i1}^T \tilde{\mathbf{x}}_{i1} (\hat{\beta}_{1n} - \beta_1) \xrightarrow{D} N(\mathbf{0}_d, \mathbf{I}_d),$$

where $\Sigma_n = \sigma^2\mathbf{A}_n\mathbf{A}_n^T$.

Lindeberg-Feller multivariate CLT

Theorem

Suppose \mathbf{X}_i is a sequence of independent random vectors with mean $\boldsymbol{\mu}_i$, covariance matrix $\boldsymbol{\Sigma}_i$. Assume that $\frac{1}{n} \sum_{i=1}^n \boldsymbol{\Sigma}_i \rightarrow \boldsymbol{\Sigma}$ as $n \rightarrow \infty$.

If for any $\epsilon > 0$,

$$\frac{1}{n} \sum_{i=1}^n E[||\mathbf{X}_i - \boldsymbol{\mu}_i||^2 I_{\{||\mathbf{X}_i - \boldsymbol{\mu}_i|| > \epsilon \sqrt{n}\}}] \rightarrow 0 \text{ as } n \rightarrow \infty$$

then

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n (\mathbf{X}_i - \boldsymbol{\mu}_i) \xrightarrow{D} N(\mathbf{0}, \boldsymbol{\Sigma}).$$

Order of p_n

- (A2.a) and (A2.b) are identical to (A1.a) and (A1.b),
if $\liminf_{n \rightarrow \infty} \min_{1 \leq j \leq k_n} |\beta_j| > 0$.

$$(A1) \text{ (a) } \lim_{n \rightarrow \infty} \frac{\sqrt{k_n} \lambda_n}{\sqrt{\rho_{n,1}}} = 0;$$

$$\text{(b) } \lim_{n \rightarrow \infty} \frac{\sqrt{p_n}}{\sqrt{n \rho_{n,1}}} = 0.$$

$$(A2) \text{ (a) } \lim_{n \rightarrow \infty} \frac{\sqrt{k_n} \lambda_n}{\sqrt{\rho_{n,1}} \min_{1 \leq j \leq k_n} |\beta_j|} = 0;$$

$$\text{(b) } \lim_{n \rightarrow \infty} \frac{\sqrt{p_n}}{\sqrt{n \rho_{n,1}} \min_{1 \leq j \leq k_n} |\beta_j|} = 0.$$

Order of p_n

▪ (A3) $\lim_{n \rightarrow \infty} \frac{\sqrt{\max(\pi_{n,k_n}, \omega_{n,m_n}) p_n}}{\sqrt{n} \rho_{n,1} \lambda_n} = 0$ is implied by

$$\lim_{n \rightarrow \infty} \frac{p_n}{\sqrt{n} \rho_{n,1} \lambda_n} = 0 \quad (\text{since } \pi_{n,k_n} \leq k_n, \omega_{n,m_n} \leq m_n)$$

▪ Suppose $\liminf_{n \rightarrow \infty} \min_{1 \leq j \leq k_n} |\beta_j| > 0$ and $\liminf_{n \rightarrow \infty} \rho_{n,1} > 0$.

Then,

(i) if $p_n = o(n^{\frac{1}{3}})$ and take suitable λ_n (e.g. $\lambda_n = O(n^{-\frac{1}{6}})$), then (A1) – (A3) can be satisfied.

(However, if $p_n = o(n^{\frac{1}{2}})$, then $\lambda_n = O(1) \Rightarrow$ (A1.a) fails)

(ii) Furthermore, suppose either k_n is fixed, or the largest eigenvalue of $\frac{1}{n} \mathbf{X}^T \mathbf{X}$ is bounded above.

Then $p_n = o(n^{\frac{1}{2}})$ is sufficient.

Conclusion

Conclusion

- Asymptotic properties of LS-SCAD estimator:
consistency, convergence rate, sparsity, asymptotic normality.
- Order of p_n that is sufficient for asymptotic properties.

Reference :

- [1] Fan, Jianqing, and Runze Li. "Variable selection via nonconcave penalized likelihood and its oracle properties." *Journal of the American statistical Association* 96.456 (2001): 1348-1360.
- [2] TIBSHIRANI, R. (1996) Regression Shrinkage and Selection via the Lasso. *J.R.Stat.Soc.B.58*, No.1, 267-288.