

# Moment convergence of regularized least-squares estimator for linear regression model

Yusuke Shimizu<sup>1</sup>

Received: 4 Jan 2016 / Revised: 4 July 2016 / Published online: 9 August 2016  
© The Institute of Statistical Mathematics, Tokyo 2016

**Abstract** In this paper, we study the uniform tail-probability estimates of a regularized least-squares estimator for the linear regression model. We make use of the polynomial type large deviation inequality for the associated statistical random fields, which may not be locally asymptotically quadratic. Our results enable us to verify various arguments requiring convergence of moments of estimator-dependent statistics, such as the mean squared prediction error and the bias correction for AIC-type information criterion.

**Keywords** Moment convergence · Regularized least-squares estimation · Sparse estimation · Large deviation inequality

## 1 Introduction

Assume that we have a sample  $\{(X_i, Y_i)\}_{i=1}^n$ , where  $Y_i \in \mathbb{R}$  and  $X_i = (X_{i,1}, \dots, X_{i,p})^\top \in \mathbb{R}^p$ , obeying the linear regression model:

$$Y_i = \theta_0^\top X_i + \epsilon_i, \quad i = 1, \dots, n, \quad (1)$$

where  $\theta_0$  is a  $p$ -dimensional true value of parameter contained in the interior of a compact parameter space  $\Theta \subset \mathbb{R}^p$  and  $(\epsilon_i)_{i=1}^n$  represent noises. Through this paper, the number of variables  $p$  is fixed. Though not essential, we suppose that the covariate  $X$  is non-random; usually,  $\{(X_i, Y_i)\}_{i=1}^n$  are standardized from the beginning such a way that

---

✉ Yusuke Shimizu  
y-shimizu@math.kyushu-u.ac.jp

<sup>1</sup> Graduate school of Mathematics, Kyushu University, 744 Motoooka Nishi-ku, Fukuoka 819-0395, Japan

$$\sum_{i=1}^n Y_i = 0, \quad \sum_{i=1}^n X_{i,j} = 0 \quad \text{and} \quad \frac{1}{n} \sum_{i=1}^n X_{i,j}^2 = 1, \quad j \in \{1, \dots, p\}.$$

For brevity, we omit the dependence of  $X_i$  and  $Y_i$  on  $n$  from the notation so long as there is no confusion (we use the notation  $X_{in}$  when we emphasize the dependence of  $X_i$  on  $n$ ). In this paper, we deal with the situation

$$\theta_0 = (z_0, \rho_0) = (z_{0,1}, \dots, z_{0,p_0}, \rho_{0,1}, \dots, \rho_{0,p_1}),$$

where  $z_{0,k} = 0$  and  $\rho_{0,l} \neq 0$  for any  $k \in \{1, \dots, p_0\}$  and  $l \in \{1, \dots, p_1\}$ ; divide the compact parameter space  $\Theta = \Theta_0 \times \Theta_1 \subset \mathbb{R}^{p_0} \times \mathbb{R}^{p_1}$  such that  $z_0 = 0 \in \Theta_0$  and  $\rho_0 \in \Theta_1$ . We can rewrite the linear regression model (1) to

$$Y_i = z_0^\top X_i^{(z)} + \rho_0^\top X_i^{(\rho)} + \epsilon_i, \quad i = 1, \dots, n, \tag{2}$$

where  $X_i^{(z)} := (X_{i,1}, \dots, X_{i,p_0})^\top$  and  $X_i^{(\rho)} := (X_{i,p_0+1}, \dots, X_{i,p_0+p_1})^\top$ , representing irrelevant and relevant covariate vectors, respectively. Then, we define the regularized least-squares estimator (regularized-LSE)  $\hat{\theta}_n = (\hat{z}_n, \hat{\rho}_n)$  as the minimizer of the contrast function

$$Z_n(\theta) = Z_n(z, \rho) := \sum_{i=1}^n (Y_i - z^\top X_i^{(z)} - \rho^\top X_i^{(\rho)})^2 + \sum_{j=1}^p p_n(\theta_j) \tag{3}$$

over  $\Theta$ , where  $p_n(\cdot)$  is a non-random and non-negative function such that  $p_n(0) = 0$ . Further conditions on  $p_n$  will be imposed later on. There is a huge literature on the sparse linear regression via regularization, where the estimator  $\hat{z}_n$  of  $z_0 = 0$  satisfies the *sparse consistency*  $P(\hat{z}_n = 0) \rightarrow 1$  as  $n \rightarrow \infty$ , while  $\sqrt{n}(\hat{\rho}_n - \rho_0)$  has a non-trivial asymptotic law. The sparse consistency implies that  $R_n \hat{z}_n = o_p(1)$  for arbitrary  $R_n \rightarrow \infty$ , for example, sparse-bridge (Radchenko 2005), the smoothly clipped absolute deviation (SCAD; Fan and Li 2001) and the Seamless- $L_0$  regularization (Dicker et al 2012). In Sect. 3, we will refer some asymptotic behaviors of these regularized estimators.

Let us mention some basic facts concerning the parametric  $M$ -estimation. Given a statistical model indexed by a finite-dimensional parameter  $\theta \in \Theta \subset \mathbb{R}^p$ , we typically estimate a true parameter value  $\theta_0 \in \Theta$  by a minimum point  $\hat{\theta}_n$  of an appropriate continuous contrast function  $Z_n : \Theta \rightarrow \mathbb{R}$ . To assess the asymptotic performance of  $\hat{\theta}_n$  quantitatively, when  $\sqrt{n}$ -consistency is concerned, we look at the statistical random fields

$$\mathbb{M}_n(w; \theta_0) := Z_n \left( \theta_0 + \frac{w}{\sqrt{n}} \right) - Z_n(\theta_0), \tag{4}$$

where  $w \in \mathbb{R}^p$ . As is well known, the weak convergence of  $\mathbb{M}_n$  to some  $\mathbb{M}_0$  over compact sets, the identifiability condition on  $\mathbb{M}_0$ , and the tightness of the scaled estimator  $\hat{w}_n := \sqrt{n}(\hat{\theta}_n - \theta_0)$  make the “argmin” functional continuous for  $\mathbb{M}_n$ :

$\hat{w}_n \in \operatorname{argmin} \mathbb{M}_n \xrightarrow{\mathcal{L}} \operatorname{argmin} \mathbb{M}_0$ . See, e.g., [van der Vaart \(1998\)](#). Further, when concerned with moments of  $\hat{w}_n$ -dependent statistics such as the mean square error, more than the weak convergence is required. Then, the *polynomial type large deviation inequality (PLDI)* of [Yoshida \(2011\)](#), which estimates the tail of  $\mathcal{L}(\hat{w}_n)$  in such a way that

$$\sup_{r>0} \sup_{n>0} r^L P(|\hat{w}_n| \geq r) < \infty \tag{5}$$

for a given  $L > 0$ , plays an important role. Assume that there exists a random variable  $\hat{w}_0$  such that  $\hat{w}_n \xrightarrow{\mathcal{L}} \hat{w}_0$ . Then, the moment convergence

$$E[|\hat{w}_n|^q] \rightarrow E[|\hat{w}_0|^q], \quad q > 0, \tag{6}$$

holds if there exists a  $q' > q$  such that  $\sup_{n>0} E[|\hat{w}_n|^{q'}] < \infty$ . Suppose that the PLDI (5) holds for some  $L > q'$ . Then, we obtain

$$\sup_{n>0} E[|\hat{w}_n|^{q'}] = \sup_{n>0} \int_0^\infty P(|\hat{w}_n|^{q'} > s) ds < \infty.$$

As the results, we get the moment convergence (6) if we ensure the PLDI (5) for some  $L > q$ . Hence, the main purpose of this paper is to derive the PLDI (5) with  $\hat{w}_n := (\sqrt{n}\hat{z}_n, \sqrt{n}(\hat{\rho}_n - \rho_0))$ .

We should mention the importance of convergence of moments: asymptotic behavior of expected values of statistics depending on estimators. It especially serves as a critical tool when, for example, analyzing the mean squared prediction error and the bias correction for information criteria; see [Chen and Ing \(2011\)](#), [Afendras and Markatou \(2015a, b\)](#), [Findley and Wei \(2002\)](#), [Uchida and Yoshida \(2001, 2006\)](#), [Sakamoto and Yoshida \(2004\)](#), as well as [Yoshida \(2011\)](#). Let us consider a typical scenario. If (6) holds with  $q = 2$  and  $\hat{w}_0 \sim N_p(0, V)$  where  $V$  is a  $p \times p$ -diagonal matrix, the mean squared error of  $\hat{\theta}_n$  can be expressed as:

$$E \left[ |\hat{\theta}_n - \theta_0|^2 \right] = \frac{\operatorname{tr}(V)}{n} + o_p \left( \frac{1}{n} \right),$$

from which the mean squared prediction error can be established with its theoretical justification. Also, the moment convergence provides benefits to AIC type information criteria, which is widely used as a simple and practical estimate of the best model, and is derived from the bias correction procedures. From the point of view of regularized estimation, AIC is used to select tuning parameters contained in regularization terms. Recently, [Umezu et al \(2015\)](#) proposed the bias-corrected AIC for non-concave regularized likelihood estimator of generalized linear model by verifying its moment convergence, and concluded that the proposed AIC performs well through simulation studies. In particular, they studied asymptotic behaviors of the estimator including  $E[|\sqrt{n}\hat{z}_n|^2] \rightarrow 0$  which cannot be deduced from the sparse consistency, and is needed

to derive AIC for sparse-type estimations. Convergence of moments for regularized estimators gives us the validity of AIC for selecting tuning parameters.

This paper is organized as follows. In Sect. 2, we will derive the PLDI (5) for the regularized-LSE of the linear regression model (2). We will give some examples of the regularization term in the contrast function (3) in Sect. 3.

For convenience of reference, we end this section with stating Theorem 1 and Theorem 3(a) of Yoshida (2011), which will play an essential role in our study. We need to introduce some notation. For any fixed  $\theta_0 \in \Theta$ , we define a random function

$$\mathbb{Y}_n(\theta; \theta_0) := -\frac{1}{n}(Z_n(\theta) - Z_n(\theta_0)).$$

Also, let  $\theta \mapsto \mathbb{Y}_0(\theta; \theta_0)$  be a random function. We consider the PLAQ representation of  $\mathbb{M}_n$ :

$$\mathbb{M}_n(w; \theta_0) = \Delta_n(\theta_0)[w] + \frac{1}{2}\Gamma_0(\theta_0)[w, w] + r_n(w; \theta_0) \tag{7}$$

for  $w \in \{w \in \mathbb{R}^p : \theta_0 + w/\sqrt{n} \in \Theta\}$ , where  $\Delta_n(\theta_0) \in \mathbb{R}^p$ ,  $\Gamma_0(\theta_0) \in \mathbb{R}^p \times \mathbb{R}^p$  and  $r_n(w; \theta_0) \in \mathbb{R}$  are random variables.<sup>1</sup> Finally, let  $\alpha \in (0, 1)$ ,  $U_n(r, \theta_0) := \{w \in \mathbb{R}^p : r \leq |w| \leq n^{(1-\alpha)/2}\}$ . We now introduce some conditions:

(A1)  $\exists v_1 > 0, \forall L > 0, \exists c_L > 0$  : constant,  $\forall r > 0$ ,

$$\sup_{n>0} P\left(\sup_{w \in U_n(r, \theta_0)} \frac{|r_n(w; \theta_0)|}{1 + |w|^2} \geq r^{-v_1}\right) \leq \frac{c_L}{r^L}.$$

(A2)  $\Gamma_0(\theta_0)$  is deterministic and positive-definite.

(A3)  $\exists \chi = \chi(\theta_0) > 0$  : non-random,  $\exists v = v(\theta_0) > 0, \forall \theta \in \Theta$ ,

$$\mathbb{Y}_0(\theta; \theta_0) \leq -\chi|\theta - \theta_0|^v.$$

(A4)  $\alpha \in (0, 1), v_1 \in (0, 1), \alpha v < v_2, \beta \in [0, \infty), 1 - 2\beta - v_2 > 0$ .

(A5)  $\forall L > 0, N_1 := L(1 - v_1)^{-1}, N_2 := L(1 - 2\beta - v_2)^{-1}$ ,

$$\begin{aligned} \sup_{n>0} E\left[|\Delta_n(\theta_0)|^{N_1}\right] &< \infty; \\ \sup_{n>0} E\left[\left(\sup_{\theta \in \Theta} n^{1/2-\beta} |\mathbb{Y}_n(\theta; \theta_0) - \mathbb{Y}_0(\theta; \theta_0)|\right)^{N_2}\right] &< \infty. \end{aligned}$$

**Theorem 1** [Yoshida (2011), Theorems 1 and 3(a)] *Assume [A1]–[A5]. Then, the estimate (5) holds.*<sup>2</sup>

<sup>1</sup> The sign in front of the quadratic term  $(1/2)\Gamma_0(\theta_0)[w, w]$  is different from the original PLAQ of Yoshida (2011) since we consider minimization of (4).

<sup>2</sup> The uniform (w.r.t.  $\theta_0$ ) evaluation is not in our scope here.

### 2 Moment convergence

In this section, we discuss the moment convergence of  $\hat{w}_n$  by checking the conditions of Theorem 1. In particular, if we have the weak convergence  $\hat{w}_n \xrightarrow{\mathcal{L}} \hat{w}_0$  for some random vector  $\hat{w}_0$ , then the moment convergence (6) holds. Let  $C_n := n^{-1} \sum_{i=1}^n X_i X_i^\top$ .

**Theorem 2** *Assume that the linear regression model is (2) and the contrast function is (3). Suppose the following conditions:*

$$\epsilon_1, \epsilon_2, \dots \text{ are i.i.d. with } E[\epsilon_i] = 0 \text{ and } \forall k > 0, E[|\epsilon_i|^k] < \infty; \tag{8}$$

$$\exists \delta > 0, \exists C_0 > 0, \sup_{n>0} (n^\delta |C_n - C_0|) < \infty; \tag{9}$$

$$\sup_{n>0} \sup_{i \leq n} |X_{in}| < \infty; \tag{10}$$

$$\exists \beta \in \left(0, \frac{1}{2}\right), \forall K \subset \mathbb{R} : \text{compact}, \sup_{n>0} \sup_{a \in K} \frac{p_n(a)}{n^{1/2+\beta}} < \infty; \tag{11}$$

$$\exists \kappa \in (0, 2), \forall a \neq 0, \exists c_a > 0 : \text{constant}, \forall b \in \mathbb{R},$$

$$\limsup_{n \rightarrow \infty} \left| p_n \left( a + \frac{b}{\sqrt{n}} \right) - p_n(a) \right| \leq c_a |b|^\kappa. \tag{12}$$

Then, the PLDI (5) holds with  $\hat{w}_n = (\sqrt{n}\hat{z}_n, \sqrt{n}(\hat{\rho}_n - \rho_0))$ . Additionally, if we have the weak convergence  $\hat{w}_n \xrightarrow{\mathcal{L}} \hat{w}_0$  for some random vector  $\hat{w}_0$ , then the moment convergence (6) holds.

*Proof* We will check the conditions of Theorem 1 to conclude (5). Set  $w = (u, v) \in \mathbb{R}^{p_0} \times \mathbb{R}^{p_1}$ . We have the statistical random fields

$$\begin{aligned} \mathbb{M}_n(w; \theta_0) &= Z_n \left( \theta_0 + \frac{w}{\sqrt{n}} \right) - Z_n(\theta_0) \\ &= \sum_{i=1}^n \left\{ \left( \epsilon_i - \frac{w^\top}{\sqrt{n}} X_i \right)^2 - \epsilon_i^2 \right\} + \sum_{k=1}^{p_0} p_n \left( \frac{u_k}{\sqrt{n}} \right) \\ &\quad + \sum_{l=1}^{p_1} \left\{ p_n \left( \rho_{0l} + \frac{v_l}{\sqrt{n}} \right) - p_n(\rho_{0l}) \right\} \\ &= - \sum_{i=1}^n \frac{2}{\sqrt{n}} \epsilon_i X_i[w] + \frac{1}{2} (2C_0)[w, w] + (C_n - C_0)[w, w] \\ &\quad + \sum_{k=1}^{p_0} p_n \left( \frac{u_k}{\sqrt{n}} \right) + \sum_{l=1}^{p_1} \left\{ p_n \left( \rho_{0l} + \frac{v_l}{\sqrt{n}} \right) - p_n(\rho_{0l}) \right\}. \end{aligned}$$

Since  $\hat{w}_n$  is a minimum point of  $\mathbb{M}_n(w; \theta_0)$  and  $\mathfrak{p}_n$  is a non-negative function, we have

$$\begin{aligned} P(|\hat{w}_n| \geq r) &\leq P\left[\sup_{|w| \geq r} \{-\mathbb{M}_n(w; \theta_0)\} \geq -\mathbb{M}_n(0; \theta_0) = 0\right] \\ &\leq P\left[\sup_{|w| \geq r} \left\{ \sum_{i=1}^n \frac{2}{\sqrt{n}} \epsilon_i X_i[w] - \frac{1}{2}(2C_0)[w, w] - (C_n - C_0)[w, w] \right. \right. \\ &\quad \left. \left. - \sum_{l=1}^{p_1} \left( \mathfrak{p}_n\left(\rho_{0l} + \frac{v_l}{\sqrt{n}}\right) - \mathfrak{p}_n(\rho_{0l}) \right) \right\} \geq 0\right]. \end{aligned}$$

Hence, we will establish the PLDI

$$\begin{aligned} \sup_{r>0} \sup_{n>0} r^L P\left[\sup_{|w| \geq r} \left\{ \sum_{i=1}^n \frac{2}{\sqrt{n}} \epsilon_i X_i[w] - \frac{1}{2}(2C_0)[w, w] - (C_n - C_0)[w, w] \right. \right. \\ \left. \left. - \sum_{l=1}^{p_1} \left( \mathfrak{p}_n\left(\rho_{0l} + \frac{v_l}{\sqrt{n}}\right) - \mathfrak{p}_n(\rho_{0l}) \right) \right\} \geq 0\right] < \infty \end{aligned} \tag{13}$$

for any  $L > 0$  to ensure the PLDI (5). We have the *PLAQ-like* expression which corresponds to (7) with

$$\Delta_n(\theta_0) = \sum_{i=1}^n \frac{2}{\sqrt{n}} \epsilon_i X_i; \tag{14}$$

$$\Gamma_0(\theta_0) = 2C_0; \tag{15}$$

$$r_n(w; \theta_0) = -(C_n - C_0)[w, w] - \sum_{l=1}^{p_1} \left\{ \mathfrak{p}_n\left(\rho_{0l} + \frac{v_l}{\sqrt{n}}\right) - \mathfrak{p}_n(\rho_{0l}) \right\}. \tag{16}$$

According to (8)–(11), we obtain for any  $\theta \in \Theta$

$$\begin{aligned} \mathbb{Y}_n(\theta; \theta_0) &= -\frac{1}{n}(Z_n(\theta) - Z_n(\theta_0)) \\ &= -\frac{1}{n} \sum_{i=1}^n \left[ \{\epsilon_i - (\theta - \theta_0)^\top X_i\}^2 - \epsilon_i^2 \right] - \frac{1}{n} \sum_{j=1}^p \{ \mathfrak{p}_n(\theta_j) - \mathfrak{p}_n(\theta_{0j}) \} \\ &= \frac{2}{n} \sum_{i=1}^n \epsilon_i X_i[\theta - \theta_0] - C_n[\theta - \theta_0, \theta - \theta_0] - \frac{1}{n} \sum_{j=1}^p \{ \mathfrak{p}_n(\theta_j) - \mathfrak{p}_n(\theta_{0j}) \} \\ &\xrightarrow{P} -C_0[\theta - \theta_0, \theta - \theta_0] =: \mathbb{Y}_0(\theta; \theta_0). \end{aligned}$$

We get  $\mathbb{Y}_0(\theta; \theta_0) \leq -\lambda_{\min}(C_0)|\theta - \theta_0|^2$  where  $\lambda_{\min}(C_0)$  denotes the minimal eigenvalue of the matrix  $C_0$ . Apparently, [A2] holds from (9) and (15), and also [A3] holds with  $\chi = \lambda_{\min}(C_0)$  and  $\nu = 2$ . Hence, it remains to verify [A1], [A4] and [A5].

First, we will verify [A1]. From (16), we have

$$\frac{|r_n(w; \theta_0)|}{1 + |w|^2} \leq \frac{|w|^2}{1 + |w|^2} |C_n - C_0| + \frac{1}{1 + |w|^2} \left| \sum_{l=1}^{p_1} \left\{ \mathfrak{p}_n \left( \rho_{0l} + \frac{v_l}{\sqrt{n}} \right) - \mathfrak{p}_n(\rho_{0l}) \right\} \right|. \tag{17}$$

Let us fix  $\beta, \nu_2, \alpha \in (0, 1)$  and  $\xi$  such that  $0 \vee (1/2 - \delta) \leq \beta < 1/2, 1 - 2\beta > \nu_2 > 2\alpha$  and  $0 < \xi < (2\alpha/(1 - \alpha)) \wedge 1$ . Note that these parameters meet  $\beta - 1/2 + (1 - \alpha)\xi/2 < 0$ . Then for the first term of the right-hand side of (17), we get from (9)

$$\begin{aligned} & \sup_{w \in U_n(r, \theta_0)} \left( \frac{|w|^2}{1 + |w|^2} |C_n - C_0| \right) \\ &= n^{1/2 - \beta - \delta} (n^\delta |C_n - C_0|) \sup_{w \in U_n(r, \theta_0)} \left( \frac{|w|^2}{1 + |w|^2} n^{\beta - 1/2} |w|^\xi |w|^{-\xi} \right) \\ &\lesssim n^{\beta - 1/2} n^{(1 - \alpha)\xi/2} r^{-\xi} \lesssim r^{-\xi}, \end{aligned} \tag{18}$$

where  $A_n \lesssim B_n$  means that  $\sup_n (A_n/B_n) < \infty$ . Next, we will estimate the second term of the right-hand side of (17). We obtain from (12) that there exists a  $\kappa \in (0, 2)$  such that

$$\frac{1}{1 + |w|^2} \left| \sum_{l=1}^{p_1} \left\{ \mathfrak{p}_n \left( \rho_{0l} + \frac{v_l}{\sqrt{n}} \right) - \mathfrak{p}_n(\rho_{0l}) \right\} \right| \lesssim \frac{|v|^\kappa}{1 + |w|^2} \lesssim |w|^{\kappa - 2}, \quad w \in U_n(r, \theta_0);$$

note that  $\sup_{w \in U_n(r, \theta_0)} |v_l|/\sqrt{n} \rightarrow 0$ . Since we can take  $\alpha \in (0, 1)$  such that  $2 - \kappa > \xi$  (note that  $0 < \xi < (2\alpha/(1 - \alpha)) \wedge 1$ ), we get

$$\sup_{w \in U_n(r, \theta_0)} |w|^{\kappa - 2} \lesssim r^{-\xi}. \tag{19}$$

Fix a  $\nu_1 \in (0, \xi)$ . Then from (17)–(19), we have for any  $L > 0$

$$\sup_{n > 0} P \left( \sup_{w \in U_n(r, \theta_0)} \frac{|r_n(w; \theta_0)|}{1 + |w|^2} \geq r^{-\nu_1} \right) \lesssim \frac{1}{r^L}.$$

This means that [A1] holds, and [A4] also holds with taking the parameters as above.

Second, we will verify [A5]. From (14), we define  $\Delta_n(\theta_0) = \sum_{i=1}^n (2/\sqrt{n}) \epsilon_i X_i =: \sum_{i=1}^n \chi_{ni}$ . Then, using Burkholder’s inequality and Jensen’s inequality, we obtain for  $N_1 = L(1 - \nu_1)^{-1} \geq 2$

$$\begin{aligned}
 \sup_{n>0} E \left[ \left| \Delta_n(\theta_0) \right|^{N_1} \right] &\leq \sup_{n>0} E \left[ \max_{j \leq n} \left| \sum_{i=1}^j \chi_{ni} \right|^{N_1} \right] \\
 &\lesssim \sup_{n>0} E \left[ \left( \sum_{i=1}^n \chi_{ni}^2 \right)^{N_1/2} \right] \\
 &\lesssim \sup_{n>0} E \left[ \frac{1}{n} \sum_{i=1}^n |\epsilon_i X_i|^{2 \cdot N_1/2} \right] \\
 &\lesssim E[|\epsilon_1|^{N_1}] \cdot \sup_{n>0} \left( \frac{1}{n} \sum_{i=1}^n |X_i|^{N_1} \right) < \infty. \tag{20}
 \end{aligned}$$

The last boundedness of (20) follows from (8) and (10). Moreover, we get for any  $\theta \in \Theta$

$$\begin{aligned}
 \sum_{i=1}^n \frac{2}{n} \epsilon_i X_i [\theta - \theta_0] - C_n[\theta - \theta_0, \theta - \theta_0] &\xrightarrow{P} -C_0[\theta - \theta_0, \theta - \theta_0]; \\
 \frac{1}{n} \sum_{j=1}^p \{ \mathfrak{p}_n(\theta_j) - \mathfrak{p}_n(\theta_{0j}) \} &\xrightarrow{P} 0.
 \end{aligned}$$

Since  $(a + b)^{N_2} \lesssim a^{N_2} + b^{N_2}$  for any  $a, b \geq 0$  and  $N_2 = L(1 - 2\beta - \nu_2)^{-1} \geq 2$ , we have

$$\begin{aligned}
 \sup_{n>0} E \left[ \sup_{\theta \in \Theta} \left( n^{1/2-\beta} \left| \sum_{i=1}^n \frac{2}{n} \epsilon_i X_i [\theta - \theta_0] - C_n[\theta - \theta_0, \theta - \theta_0] + C_0[\theta - \theta_0, \theta - \theta_0] \right| \right)^{N_2} \right] \\
 \lesssim \sup_{n>0} \left( n^{-\beta N_2} E \left[ \left| \sum_{i=1}^n \frac{1}{\sqrt{n}} \epsilon_i X_i \right|^{N_2} \right] \right) + \left\{ \sup_{n>0} (n^{1/2-\beta-\delta} n^\delta |C_n - C_0|) \right\}^{N_2} < \infty. \tag{21}
 \end{aligned}$$

Note that the parameter space  $\Theta$  is a compact set. Further, we obtain

$$\sup_{n>0} \sup_{\theta \in \Theta} \left[ n^{1/2-\beta} \left| \frac{1}{n} \sum_{j=1}^p \{ \mathfrak{p}_n(\theta_j) - \mathfrak{p}_n(\theta_{0j}) \} \right| \right]^{N_2} < \infty \tag{22}$$

since we assume (11). From (20)–(22), we conclude that [A5] holds. Therefore, the proof of (5) is complete because we established the PLDI (13). The latter claim of the theorem is trivial.  $\square$

*Remark 1* We could deal with random design  $(X_i)$ . Assume for simplicity that  $(X_i)$  and  $(\epsilon_j)$  are independent. Then, in order to conclude (5), we need to change (9) and (10) into (23) and (24), respectively:

$$\exists \delta > 0, \exists C_0 > 0 : \text{constant}, \forall k > 0, \sup_{n>0} E \left[ \left| n^\delta (C_n - C_0) \right|^k \right] < \infty. \tag{23}$$



$$\forall k > 0, \sup_{n>0} \sup_{i \leq n} E[|X_{in}|^k] < \infty. \tag{24}$$

The corresponding proofs are entirely analogous to the case of deterministic  $X$ .  $\square$

### 3 Examples

We will give some examples of the regularization term in (3) satisfying the conditions (11) and (12) in Theorem 2: sparse-bridge (Radchenko 2005), the smoothly clipped absolute deviation (SCAD; Fan and Li 2001) and the Seamless- $L_0$  regularization (Dicker et al 2012). From the previous studies, it is known that these regularized estimators  $\hat{\theta}_n = (\hat{z}_n, \hat{\rho}_n)$  have the sparse consistency  $P(\hat{z}_n = 0) \rightarrow 1$ , which concludes the sparse estimation, and the asymptotic laws of  $\sqrt{n}(\hat{\rho}_n - \rho_0)$  under some appropriate regularity conditions. Also when the number of variables  $p = p_n \rightarrow \infty$  as  $n \rightarrow \infty$ , the asymptotic behavior of the SCAD and the Seamless- $L_0$  estimators are known, but once again note that we consider the case that  $p$  is fixed.

#### 3.1 Sparse-bridge

In this section, we will focus on the sparse-bridge LSE defined the contrast function to be

$$Z_n(\theta) = Z_n(z, \rho) := \sum_{i=1}^n (Y_i - z^\top X_i^{(z)} - \rho^\top X_i^{(\rho)})^2 + \lambda_n \sum_{j=1}^p |\theta_j|^\gamma, \tag{25}$$

where  $\lambda_n \geq 0$  denotes the tuning parameter controlling the degree of regularization together with the bridge index  $\gamma \in (0, 1)$ . This means  $p_n(\cdot) = \lambda_n |\cdot|^\gamma$ . Denote by  $\hat{\theta}_n = (\hat{z}_n, \hat{\rho}_n)$  a minimizer of  $Z_n$  over a compact parameter space  $\Theta = \Theta_0 \times \Theta_1 \subset \mathbb{R}^{p_0} \times \mathbb{R}^{p_1}$ . The asymptotic behavior of  $\hat{\theta}_n$  is studied by Radchenko (2005). He assumed regularity conditions including that the noises  $\epsilon_1, \epsilon_2, \dots$  are i.i.d. with  $E[\epsilon_i] = 0$  and  $E[\epsilon_i^2] =: \sigma^2 > 0$ ,  $C_n \rightarrow C_0$  for some  $C_0 > 0$  and that  $n^{-1} \max_{i \leq n} |X_i|^2 \rightarrow 0$ . Note that these conditions are satisfied with (8)–(10). Then, he proved the following results:

- The sparse consistency of  $\hat{z}_n$ :

$$P(\hat{z}_n = 0) \rightarrow 1 \text{ if } \lambda_n/n^{\gamma/2} \rightarrow \infty \text{ and } \lambda_n/n \rightarrow 0.$$

- The asymptotic laws of  $\hat{\rho}_n$ :

- (i)  $\sqrt{n}(\hat{\rho}_n - \rho_0) \xrightarrow{\mathcal{L}} N_{p_1}(-\lambda_0 B_0^{-1} \Upsilon, \sigma^2 B_0^{-1})$  if  $\lambda_n/n^{\gamma/2} \rightarrow \infty$  and  $\lambda_n/\sqrt{n} \rightarrow \lambda_0 \geq 0$ ;
- (ii)  $n\lambda_n^{-1}(\hat{\rho}_n - \rho_0) \xrightarrow{\mathcal{L}} -B_0^{-1} \Upsilon$  if  $\lambda_n/\sqrt{n} \rightarrow \infty$  and  $\lambda_n/n \rightarrow 0$ ,

where

$$\Upsilon := \frac{\gamma}{2} \{\text{sgn}(\rho_{0,1})|\rho_{0,1}|^{\gamma-1}, \dots, \text{sgn}(\rho_{0,p_1})|\rho_{0,p_1}|^{\gamma-1}\}$$

and  $B_0$  is the  $p_1 \times p_1$  submatrix located in the bottom right corner of the matrix  $C_0$ . We are concerned here with the moment convergence of  $\hat{w}_n$ . With regard to the asymptotic law of the non-zero parameter  $\rho$ , we only consider the case (i), where the asymptotic distribution is non-degenerate. The following Corollary 1 is derived from Theorem 2.

**Corollary 1** *Assume that the linear regression model is (2) and the contrast function is (25), where  $\lambda_n/n^{\gamma/2} \rightarrow \infty$  and  $\lambda_n/\sqrt{n} \rightarrow \lambda_0 \geq 0$  for  $\gamma \in (0, 1)$ . Suppose that we have (8)–(10). Then, the PLDI (5) holds. In particular, the moment convergence (6) holds with  $\hat{w}_0 = (0, \hat{v}_0)$ , where  $\mathcal{L}(\hat{v}_0) = N_{p_1}(-\lambda_0 B_0^{-1} \Upsilon, \sigma^2 B_0^{-1})$ .*

*Proof* Apparently, we only need to check the conditions (11) and (12) in Theorem 2 for  $\mathfrak{p}_n(\cdot) = \lambda_n |\cdot|^\gamma$ . (11) follows easily since for any  $a \in \mathbb{R}$ , we have

$$\frac{\mathfrak{p}_n(a)}{\sqrt{n}} = \frac{\lambda_n}{\sqrt{n}} |a|^\gamma \lesssim 1$$

from  $\lambda_n/\sqrt{n} \rightarrow \lambda_0 \geq 0$ . We will show (12). When  $n$  is large enough, we have for any  $a \neq 0$  and  $b \in \mathbb{R}$

$$\begin{aligned} \left| \mathfrak{p}_n\left(a + \frac{b}{\sqrt{n}}\right) - \mathfrak{p}_n(a) \right| &= \lambda_n \left| \left| a + \frac{b}{\sqrt{n}} \right|^\gamma - |a|^\gamma \right| \\ &\lesssim \frac{\lambda_n}{\sqrt{n}} |b| \lesssim |b|. \end{aligned}$$

This shows that (12) holds for  $\kappa = 1$ , hence we obtain the PLDI (5). The latter claim is trivial since we have  $(\sqrt{n}\hat{z}_n, \sqrt{n}(\hat{\rho}_n - \rho_0)) \xrightarrow{\mathcal{L}} (0, \hat{v}_0)$ , where  $\mathcal{L}(\hat{v}_0) = N_{p_1}(-\lambda_0 B_0^{-1} \Upsilon, \sigma^2 B_0^{-1})$ . □

*Remark 2* Here, we briefly mention the case of the bridge-LSE  $\hat{\theta}_n$  defined as the minimal point of the contrast function

$$Z_n(\theta) := \sum_{i=1}^n (Y_i - \theta^\top X_i)^2 + \lambda_n \sum_{j=1}^p |\theta_j|^\gamma, \tag{26}$$

where  $\lambda_n \geq 0$  and  $\gamma > 0$  satisfy that  $\lambda_n/n^{(1 \wedge \gamma)/2} \rightarrow \lambda_0 \geq 0$ ; then, we do not have the sparse consistency. Note that, different from (25), in (26) we do not divide the true value of parameter  $\theta_0$  into the zero part and the non-zero part: jointly estimate all the components. We assume (8)–(10). Then, Knight and Fu (2000) proved the following asymptotic behavior of  $\hat{\theta}_n$ .

– Consistency:

$$\hat{\theta}_n \xrightarrow{P} \theta_0 \text{ if } \lambda_n/n \rightarrow 0.$$

– Asymptotic laws:

$$\hat{w}_0 = \sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{\mathcal{L}} \operatorname{argmin}(V_0) \text{ if } \lambda_n/n^{(1+\gamma)/2} \rightarrow \lambda_0 \geq 0,$$

where for  $W \sim N_p(0, \sigma^2 C_0)$ ,

$$V_0(w) := \begin{cases} -2W[w] + C_0[w, w] + \gamma\lambda_0 \sum_{j=1}^p w_j \operatorname{sgn}(\theta_{0j})|\theta_{0j}|^{\gamma-1} & (\gamma > 1), \\ -2W[w] + C_0[w, w] & \\ +\lambda_0 \sum_{j=1}^p \{w_j \operatorname{sgn}(\theta_{0j})I(\theta_{0j} \neq 0) + |w_j|I(\theta_{0j} = 0)\} & (\gamma = 1), \\ -2W[w] + C_0[w, w] + \lambda_0 \sum_{j=1}^p |w_j|^\gamma I(\theta_{0j} = 0) & (\gamma < 1). \end{cases}$$

Let  $\hat{w}_0 = \operatorname{argmin}(V_0)$ . We can derive the PLDI for the bridge-LSE by making use of the argument similar to the proof of Theorem 2. In particular, for every continuous  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  of at most polynomial growth,  $E[f(\hat{w}_n)] \rightarrow E[f(\hat{w}_0)]$ . See (Masuda and Shimizu 2014, Section 2) for details. □

### 3.2 SCAD

The SCAD-LSE (Fan and Li 2001) is defined as the minimum point of the contrast function (3), where

$$p_n(\theta_j) = \begin{cases} n\lambda_n|\theta_j| & (|\theta_j| \leq \lambda_n), \\ \frac{-n(\theta_j^2 - 2\tau\lambda_n|\theta_j| + \lambda_n^2)}{2(\tau - 1)} & (\lambda_n < |\theta_j| \leq \tau\lambda_n), \\ \frac{n(\tau + 1)\lambda_n^2}{2} & (|\theta_j| > \tau\lambda_n). \end{cases}$$

$\tau > 2$  is an additional tuning parameter. Let the minimizer be  $\hat{\theta}_n = (\hat{z}_n, \hat{\rho}_n)$ , and (8)–(10) hold. Then, under some conditions including

$$\lambda_n \rightarrow 0, \sqrt{n}\lambda_n \rightarrow \infty, \tag{27}$$

they proved the sparse consistency and the asymptotic law of  $\rho_n$ :

$$\sqrt{n}(\hat{\rho}_n - \rho_0) \xrightarrow{\mathcal{L}} N_{p_1}(0, \mathcal{I}_{p_1}^{-1}(\rho_0)),$$

where  $\mathcal{I}_{p_1}(\rho_0) = \mathcal{I}_{p_1}(0, \rho_0)$  denotes the  $p_1 \times p_1$  Fisher information matrix knowing  $z_0 = 0$ .

Let us take  $\lambda_n \sim n^{\beta-1/2}$ , where  $\beta$  is the same as in the proof of Theorem 2. This meets (27). Now, we will show (11) and (12). First, we establish (11). Obviously, we only need to consider the case  $\lambda_n < |\theta_j| \leq \tau\lambda_n$ . When  $n$  is large enough, we have  $n\theta_j^2/n^{\beta+1/2} \lesssim n\lambda_n^2/n^{\beta+1/2} \sim n^{1+2\beta-1-\beta-1/2} = n^{\beta-1/2} \lesssim 1$ , hence (11) holds. To ensure (12), we use

$$p'_n(\theta_j) = \lambda_n n \left\{ I(\theta_j \leq \lambda_n) + \frac{(\tau\lambda_n - \theta_j)_+}{(\tau - 1)\lambda_n} I(\theta_j > \lambda_n) \right\}, \quad \theta_j > 0,$$

where  $(x)_+ = \max(0, x)$ . When  $n$  is large enough, for any  $a > 0$  and  $b \in \mathbb{R}$

$$\begin{aligned} \left| p_n\left(a + \frac{b}{\sqrt{n}}\right) - p_n(a) \right| &\leq \frac{|b|}{\sqrt{n}} \int_0^1 \left| p'_n\left(a + \frac{b}{\sqrt{n}}t\right) \right| dt \\ &\sim \lambda_n \sqrt{n} |b| \int_0^1 I\left(a + \frac{b}{\sqrt{n}}t \leq \lambda_n\right) dt \\ &\quad + \sqrt{n} |b| \int_0^1 \frac{[\tau\lambda_n - \{a + (b/\sqrt{n})t\}]_+}{\tau - 1} I\left(a + \frac{b}{\sqrt{n}}t > \lambda_n\right) dt \\ &\lesssim |b|. \end{aligned}$$

Similarly, we get the same estimate for  $a < 0$ . As the results, it is possible to take  $\lambda_n$  ensuring (6), where  $\hat{w}_0 = (0, \hat{v}_0)$  and  $\mathcal{L}(\hat{v}_0) = N_{p_1}(0, \mathcal{I}_{p_1}^{-1}(\rho_0))$ .

### 3.3 Seamless- $L_0$

The Seamless- $L_0$  regularization (Dicker et al 2012), which approximates the (technically unpleasant due to its discontinuity at the origin)  $L_0$ -loss, is given by

$$Z_n(\theta) = Z_n(z, \rho) := \sum_{i=1}^n (Y_i - z^\top X_i^{(z)} - \rho^\top X_i^{(\rho)})^2 + \frac{2n\lambda_n}{\log 2} \sum_{j=1}^p \log\left(\frac{|\theta_j|}{|\theta_j| + \tau_n} + 1\right),$$

where  $\tau_n > 0$  is an additional tuning parameter. Let the minimizer be  $\hat{\theta}_n = (\hat{z}_n, \hat{\rho}_n)$  and (8)–(10) hold. Then, under some conditions including

$$\lambda_n = O(1), \quad \lambda_n \sqrt{n} \rightarrow \infty, \quad \tau_n = O(n^{-3/2}), \tag{28}$$

Dicker et al (2012) proved the sparse consistency and the asymptotic law of  $\hat{\rho}_n$ :

$$\sqrt{n}(\hat{\rho}_n - \rho_0) \xrightarrow{\mathcal{L}} N_{p_1}(0, \sigma^2 B_0^{-1}),$$

where  $B_0$  is the same as in Sect. 3.1.

Let us take  $\lambda_n \sim n^{\beta-1/2}$  and  $\tau_n \sim n^{-3/2}$ , where  $\beta$  is the same as in the proof of Theorem 2. This meets (28). Now, we will show (11) and (12) for  $p_n(\cdot) =$

$(2n\lambda_n/\log 2) \log\{|\cdot|/(|\cdot| + \tau_n) + 1\}$ . (11) follows easily since  $p_n/n^{1/2+\beta} \lesssim n^{1+\beta-1/2-1/2-\beta} = 1$ . To ensure (12), we make use of the equation

$$|\log(1+x) - \log(1+x')| = \left| \int_0^1 \frac{ds}{1+x'+(x-x')s} (x-x') \right|$$

where  $x, x' > 0$ . When  $n$  is large enough, for any  $a > 0$  and  $b \in \mathbb{R}$

$$\begin{aligned} \left| p_n\left(a + \frac{b}{\sqrt{n}}\right) - p_n(a) \right| &= \frac{2n\lambda_n}{\log 2} \left| \log\left(\frac{|a + b/\sqrt{n}|}{|a + b/\sqrt{n}| + \tau_n} + 1\right) - \log\left(\frac{|a|}{|a| + \tau_n} + 1\right) \right| \\ &\lesssim n\lambda_n \left| \frac{a + \delta}{a + \delta + \tau_n} - \frac{a}{a + \tau_n} \right| \quad (\delta := b/\sqrt{n}) \\ &= n\lambda_n \frac{|(a + \delta)(a + \tau_n) - a(a + \delta + \tau_n)|}{(a + \delta + \tau_n)(a + \tau_n)} \\ &= n\lambda_n \frac{\tau_n |\delta|}{(a + \delta + \tau_n)(a + \tau_n)} \\ &\sim n^{\beta-3/2} |b| \lesssim |b|. \end{aligned}$$

Similarly, we get the same estimate for  $a < 0$ . As the results, it is possible to take the tuning parameters ensuring (6), where  $\hat{w}_0 = (0, \hat{v}_0)$  and  $\mathcal{L}(\hat{v}_0) = N_{p_1}(0, \sigma^2 B_0^{-1})$ .

**Acknowledgements** The author is grateful to the referees for their helpful comments and suggestions. He also thanks Professor H. Masuda for his valuable comments.

## References

- Afendras, G., Markatou, M. (2015a). *Optimality of training/test size and resampling effectiveness of cross-validation estimators of the generalization error*. [arXiv:1511.02980v1](https://arxiv.org/abs/1511.02980v1)
- Afendras, G., Markatou, M. (2015b). *Uniform integrability of the OLS estimators, and the convergence of their moments*. [arXiv:1511.02962v1](https://arxiv.org/abs/1511.02962v1)
- Chan, N. H., Ing, C.-K. (2011). Uniform moment bounds of Fisher's information with applications to time series. *The Annals of Statistics*, 39(3), 1526–1550.
- Dicker, L., Huang, B., Lin, X. (2012). Variable selection and estimation with the seamless- $L_0$  penalty. *Statistica Sinica*, 23(2), 929–962.
- Fan, J., Li, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96(456), 1348–1360.
- Findley, D. F., Wei, C.-Z. (2002). AIC, overfitting principles, and the boundedness of moments of inverse matrices for vector autoregressions and related models. *Journal of Multivariate Analysis*, 83(2), 415–450.
- Knight, K., Fu, W. (2000). Asymptotics for lasso-type estimators. *The Annals of Statistics*, 28(5), 1356–1378.
- Masuda, H., Shimizu, Y. (2014). *Moment convergence in regularized estimations*. [arXiv:1406.6751v2](https://arxiv.org/abs/1406.6751v2)
- Radchenko, P. (2005). Reweighting the lasso. *2005 Proceedings of the American Statistical Association [CD-ROM]*. <http://www-rcf.usc.edu/~radchenk/Lasso.pdf>
- Sakamoto, Y., Yoshida, N. (2004). Asymptotic expansion formulas for functionals of  $\epsilon$ -Markov processes with a mixing property. *Annals of the Institute of Statistical Mathematics*, 56(3), 545–597.
- Uchida, M., Yoshida, N. (2001). Information criteria in model selection for mixing processes. *Statistical Inference for Stochastic Processes*, 4(1), 73–98.
- Uchida, M., Yoshida, N. (2006). Asymptotic expansion and information criteria. *SUT Journal of Mathematics*, 42(1), 31–58.

- Umezu, Y., Shimizu, Y., Masuda, H., Ninomiya, Y. (2015). *AIC for non-concave penalized likelihood method*. [arXiv:1509.01688](https://arxiv.org/abs/1509.01688). (Submitted for publication).
- van der Vaart, A. W. (1998). Asymptotic statistics. In *Cambridge Series in Statistical and Probabilistic Mathematics* (Vol. 3). Cambridge: Cambridge University Press.
- Yoshida, N. (2011). Polynomial type large deviation inequalities and quasi-likelihood analysis for stochastic differential equations. *Annals of the Institute of Statistical Mathematics*, 63(3), 431–479.