A New Scope of Penalized Empirical Likelihood with High-Dimensional Estimating Equations

> Jinyuan Chang University of Melbourne Homepage: www.jinyuanchang.com

This is a joint work with Cheng Yong Tang and Tong Tong Wu

June 14, 2016

# Model setting

- ▶ X<sub>1</sub>,..., X<sub>n</sub>: *d*-dimensional i.i.d. observations.
- $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)^{\mathrm{T}}$ : *p*-dimensional parameter.
- ▶  $\mathbf{g}(\mathbf{X}; \boldsymbol{\theta}) = (g_1(\mathbf{X}; \boldsymbol{\theta}), \dots, g_r(\mathbf{X}; \boldsymbol{\theta}))^{\mathrm{T}}$ : an *r*-dimensional estimating function.
- The true parameter  $\theta_0$  is determined by the moment condition

 $\mathbb{E}\{\mathbf{g}(\mathbf{X}_i;\boldsymbol{\theta}_0)\}=\mathbf{0}.$ 

• The goal of this paper is to construct the estimation of  $\theta_0$  when p and/or  $r \gg n$ .

Both p and r are fixed:

- ► Hansen (1982): Generalized Methods of Moments (GMM).
- Qin and Lawless (1994): Using the idea of Empirical Likelihood (EL) proposed by Owen (1988, 1990).
- Newey and Smith (2004): Investigate higher order properties of GMM and generalized EL.

Some other advantages of EL: Wilks' theorems, Bartlett-correctable, ...

Define an EL with estimating equations as

$$L(\boldsymbol{\theta}) = \sup \left\{ \prod_{i=1}^{n} \pi_i : \pi > 0, \ \sum_{i=1}^{n} \pi_i = 1, \ \sum_{i=1}^{n} \pi_i \mathbf{g}(\mathbf{X}_i; \boldsymbol{\theta}) = \mathbf{0} \right\},\$$

and estimate  $\theta_0$  as  $\widehat{\theta} = \arg \max_{\theta \in \Theta} L(\theta)$ , which is equivalent to

$$\widehat{\boldsymbol{\theta}} = \arg\min_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \max_{\boldsymbol{\lambda} \in \widehat{\Lambda}_n(\boldsymbol{\theta})} \sum_{i=1}^n \log\{1 + \boldsymbol{\lambda}^{\mathrm{T}} \mathbf{g}(\mathbf{X}_i; \boldsymbol{\theta})\},\$$

where  $\widehat{\Lambda}_n(\boldsymbol{\theta}) = \{ \boldsymbol{\lambda} \in \mathbb{R}^r : \boldsymbol{\lambda}^{\mathrm{T}} \mathbf{g}(\mathbf{X}_i; \boldsymbol{\theta}) \in \mathcal{V}, i = 1, \dots, n \}$  for  $\boldsymbol{\theta} \in \boldsymbol{\Theta}$  and  $\mathcal{V}$  is an open interval containing zero.

•  $\hat{\theta}$  is called the maximum EL estimator.

When p or/and r diverging,

- Leng and Tang (2012), and Chang, Chen and Chen (2015) investigated the asymptotic behaviors of the EL estimator θ in high-dimensional settings.
- Hjort, McKeague and Van Keilegom (2009), Chen, Peng and Qin (2009), and Tang and Leng (2010) considered EL for some other high-dimensional problems.
- However, EL only works when both p and r are growing at some rate slower than n. Since |λ|<sub>2</sub> is required to be o<sub>p</sub>(1) in existing analyses of EL.

Such kind of restrictions on the diverging rates of p and r do not match the real practical ultra high-dimensional problems.

# New scope of our work

- ▶ The relationship between *r* and *p* in the new high-dimensional paradigm with sparsity.
  - $r \ge p$  is no longer required, facilitated with penalized EL.
- Penalizing the Lagrange multiplier  $\lambda$ :
  - to obtain sparse  $\lambda$  with many zero components;
  - to effectively reduce the effective dimensionality of the number of estimating functions.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

 Solving the problem of EL with ultra high-dimensional parameters and estimating functions. Why EL works?

## Proposition 1

Assume that there exist uniform constants  $C_1 > 0$ ,  $C_2 > 1$  and  $\gamma > 2$  such that

$$\max_{1 \le j \le r} \mathbb{E} \bigg\{ \sup_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} |g_j(\mathbf{X}_i; \boldsymbol{\theta})|^{\gamma} \bigg\} \le C_1,$$

and

$$\mathbb{P}\left[C_{2}^{-1} \leq \inf_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \lambda_{\min}\left\{\frac{1}{n} \sum_{i=1}^{n} \mathbf{g}(\mathbf{X}_{i}; \boldsymbol{\theta}) \mathbf{g}(\mathbf{X}_{i}; \boldsymbol{\theta})^{\mathrm{T}}\right\}\right]$$
$$\leq \sup_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \lambda_{\max}\left\{\frac{1}{n} \sum_{i=1}^{n} \mathbf{g}(\mathbf{X}_{i}; \boldsymbol{\theta}) \mathbf{g}(\mathbf{X}_{i}; \boldsymbol{\theta})^{\mathrm{T}}\right\} \leq C_{2}\right] \rightarrow 1.$$

If  $r = o(n^{1/2-1/\gamma})$ , then the maximum EL estimator  $\widehat{\theta}$  satisfies  $|\bar{\mathbf{g}}(\widehat{\theta})|_2 = O_p(r^{1/2}n^{-1/2})$  where  $\bar{\mathbf{g}}(\widehat{\theta}) = n^{-1}\sum_{i=1}^n \mathbf{g}(\mathbf{X}_i; \widehat{\theta})$ .

- Proposition 1 shows that regardless the dimensionality p of the parameter  $\theta$ , with r unbiased estimating functions,  $\hat{\theta}$  always ensures that  $|\bar{\mathbf{g}}(\hat{\theta})|_2$  to be of  $O_p(r^{1/2}n^{-1/2})$ .
- If sup<sub>θ∈Θ</sub> |g(θ) E{g(X<sub>i</sub>; θ)}|<sub>∞</sub> →<sub>p</sub> 0 holds, we have |E{g(X<sub>i</sub>; θ̂)}|<sub>∞</sub> = o<sub>p</sub>(1) if r = o(n<sup>1/2-1/γ</sup>). With the unique identification condition, θ<sub>0</sub> is the unique solution of E{g(X<sub>i</sub>; θ)} = 0 (which is usually guaranteed by the condition r ≥ p), we can obtain the consistency of θ̂.
- ▶ However,  $\hat{\theta}$  is clearly not uniquely defined when r < p rendering inapplicability of the method in practice.

(日) (同) (三) (三) (三) (○) (○)

► To resolve the ambiguity in the estimator  $\hat{\theta}$ , we take the approach that searches for a sparse optimizer with appropriate penalization on  $\theta$  under extra assumption that  $\theta_0$  is sparse. Write  $\theta_0 = (\theta_1^0, \ldots, \theta_p^0)^T$  and let  $S = \{1 \le k \le p : \theta_k^0 \ne 0\}$  with s = |S|. Sparsity means  $s \ll n$ . Without lose of generality, we write  $\theta_0 = (\theta_{0,(1)}^T, \theta_{0,(2)}^T)^T$  where  $\theta_{0,(1)} \in \mathbb{R}^s$  being nonzero components and  $\theta_{0,(2)} = \mathbf{0} \in \mathbb{R}^{p-s}$ .

$$\widetilde{\boldsymbol{\theta}}_n = \arg\min_{\boldsymbol{\theta}\in\boldsymbol{\Theta}} \max_{\boldsymbol{\lambda}\in\widehat{\Lambda}_n(\boldsymbol{\theta})} \bigg[ \sum_{i=1}^n \log\{1 + \boldsymbol{\lambda}^{\mathrm{T}} \mathbf{g}(\mathbf{X}_i; \boldsymbol{\theta})\} + n \sum_{k=1}^p p_{1,\pi}(|\boldsymbol{\theta}_k|) \bigg],$$

where  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)^{\mathrm{T}}$ ,  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_r)^{\mathrm{T}}$ , and  $p_{1,\pi}(\cdot)$  is a penalty function with tuning parameter  $\pi$ .

(日) (同) (三) (三) (三) (○) (○)

#### **Proposition 2**

Let  $a_n = \sum_{k=1}^p p_{1,\pi}(|\theta_k^0|)$  and  $b_n = \max\{rn^{-1}, a_n\}$ . Let  $\Theta_* = \{\theta \in \Theta : |\theta - \theta_0|_{\infty} \le \epsilon\}$  for some sufficiently small  $\epsilon > 0$ . Assume that

$$\max_{1 \le j \le r} \mathbb{E} \bigg\{ \sup_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_*} |g_j(\mathbf{X}_i; \boldsymbol{\theta})|^{\gamma} \bigg\} \le C_1,$$

$$\mathbb{P}\left[C_2^{-1} \leq \inf_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_*} \lambda_{\min} \left\{ \frac{1}{n} \sum_{i=1}^n \mathbf{g}(\mathbf{X}_i; \boldsymbol{\theta}) \mathbf{g}(\mathbf{X}_i; \boldsymbol{\theta})^{\mathrm{T}} \right\} \\ \leq \sup_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_*} \lambda_{\max} \left\{ \frac{1}{n} \sum_{i=1}^n \mathbf{g}(\mathbf{X}_i; \boldsymbol{\theta}) \mathbf{g}(\mathbf{X}_i; \boldsymbol{\theta})^{\mathrm{T}} \right\} \leq C_2 \right] \to 1.$$

and

$$\inf_{\boldsymbol{\theta} \in \{\boldsymbol{\theta} \in \boldsymbol{\Theta}_* : |\boldsymbol{\theta} - \boldsymbol{\theta}_0|_{\infty} > \varepsilon\}} |\mathbb{E}\{\mathbf{g}(\mathbf{X}_i; \boldsymbol{\theta})\}|_{\infty} \ge \Delta(\varepsilon)$$

for any  $0 < \varepsilon < \epsilon$ , where  $\Delta(\cdot)$  is a positive function satisfying  $\liminf_{\varepsilon \to 0^+} \varepsilon^{-\beta} \Delta(\varepsilon) \ge C_3$  for some uniform constants  $C_3 > 0$  and  $\beta > 0$ . If  $r = o(n^{1/2-1/\gamma})$  and  $b_n = o(n^{-2/\gamma})$ , then there exists a local solution  $\tilde{\theta}_n \in \Theta_*$ such that  $|\tilde{\theta}_n - \theta_0|_{\infty} = O_p \{b_n^{1/(2\beta)}\}$ .

- Proposition 2 implies that if facilitated with sparsity of the model parameter, the consistency of a *p*-dimensional estimator does not necessarily require  $r \ge p$ . In particular, under the local identification condition near the truth, a consistent and sparse estimator of  $\theta_0$  is feasible with no explicit requirement on the relationship between r and p, as long as the r unbiased estimating functions provide information for the zero and nonzero components in  $\theta_0$ .
- This result ensures the important flexibility in the number of estimating functions for obtaining a sparse estimator using penalized EL. That is, to obtain a sparse estimator in penalized EL whose nonzero components are consistent to the truth, one can actually opt to use smaller number of estimating equations. Further, we can show that  $\tilde{\theta}_n$  can estimate zero components of  $\theta_0$  as zero with a high probability as in the following proposition.

## **Proposition 3**

Under the same conditions as in Proposition 2, if we further assume  $g_j(\mathbf{X}; \boldsymbol{\theta})$  to be continuously differentiable with respect to  $\boldsymbol{\theta} \in \boldsymbol{\Theta}_*$  for any  $\mathbf{X}$  and  $j = 1, \ldots, r$  satisfying the condition

$$\sup_{\boldsymbol{\theta}\in\boldsymbol{\Theta}_{*}} \max_{1\leq j\leq r} \max_{k\notin\mathcal{S}} \left\{ \frac{1}{n} \sum_{i=1}^{n} \left| \frac{\partial g_{j}(\mathbf{X}_{i};\boldsymbol{\theta})}{\partial \theta_{k}} \right| \right\} = O_{p}(\varphi_{n})$$

for some  $\varphi_n > 0$ , which may diverge with n. Suppose there exist a constant 0 < c < 1 and  $\chi_n \to 0$  such that

$$\max_{k \in \mathcal{S}} \sup_{c \mid \theta_k^0 \mid < t < c^{-1} \mid \theta_k^0 \mid} p'_{1,\pi}(t) = O(\chi_n).$$

In addition to the restrictions imposed on r and  $b_n$  for the consistency of  $\hat{\theta}_n$ above, if r and  $b_n$  also satisfy the restriction conditions such that  $b_n = o(\min_{k \in S} |\theta_k^0|^{2\beta})$ ,  $s\chi_n b_n^{1/(2\beta)} = O(rn^{-1})$  and  $rn^{-1/2}\varphi_n = o(\pi)$ , then  $\mathbb{P}\{\tilde{\theta}_{n,(2)} = \mathbf{0}\} \to 1$  where  $\tilde{\theta}_{n,(2)}$  is the corresponding estimation for the zero components of  $\theta_0$  in  $\tilde{\theta}_n$ . • Though Propositions 2 and 3 are general, there are still restrictions because conditions are violated when r is large. With Proposition 1 ensuring the behavior of the estimating equations at the maximum EL estimate, and Propositions 2 and 3 breaking the requirement of  $r \ge p$  for a consistent and sparse estimator  $\tilde{\theta}_n$ .

Define

$$\begin{aligned} \widehat{\boldsymbol{\theta}}_n &= \arg\min_{\boldsymbol{\theta}\in\boldsymbol{\Theta}} \max_{\boldsymbol{\lambda}\in\widehat{\Lambda}_n(\boldsymbol{\theta})} \bigg[ \sum_{i=1}^n \log\{1 + \boldsymbol{\lambda}^{\mathrm{T}} \mathbf{g}(\mathbf{X}_i; \boldsymbol{\theta})\} - n \sum_{j=1}^r p_{2,\nu}(|\lambda_j|) \\ &+ n \sum_{k=1}^p p_{1,\pi}(|\boldsymbol{\theta}_k|) \bigg], \end{aligned}$$

where  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)^{\mathrm{T}}$ ,  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_r)^{\mathrm{T}}$ , and  $p_{1,\pi}(\cdot)$  and  $p_{2,\nu}(\cdot)$  are two penalty functions with tuning parameters  $\pi$  and  $\nu$ , respectively.

# Main results - Consistency

## Condition 1

There exist some  $K_1 > 0$  and  $\gamma > 4$  such that

$$\max_{1 \leq j \leq r} \mathbb{E} \bigg\{ \sup_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} |g_j(\mathbf{X}_i; \boldsymbol{\theta})|^{\gamma} \bigg\} \leq K_1.$$

## Condition 2

There exists a positive function  $\Delta(\cdot)$  such that

$$\inf_{\boldsymbol{\theta} \in \{\boldsymbol{\theta} \in \boldsymbol{\Theta} : |\boldsymbol{\theta} - \boldsymbol{\theta}_0|_{\infty} > \varepsilon\}} |\mathbb{E}\{\mathbf{g}(\mathbf{X}_i; \boldsymbol{\theta})\}|_{\infty} \ge \Delta(\varepsilon)$$

for any  $\varepsilon > 0$ , and there exist  $K_2 > 0$  and  $\beta > 0$  such that  $\liminf_{\varepsilon \to 0^+} \varepsilon^{-\beta} \Delta(\varepsilon) \ge K_2.$  • For  $\theta \in \Theta$  and  $\lambda \in \widehat{\Lambda}_n(\theta)$ , we define

$$f(\boldsymbol{\lambda};\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^{n} \log\{1 + \boldsymbol{\lambda}^{\mathrm{T}} \mathbf{g}(\mathbf{X}_{i};\boldsymbol{\theta})\} - \sum_{j=1}^{r} p_{2,\nu}(|\lambda_{j}|),$$
$$S_{n}(\boldsymbol{\theta}) = \max_{\boldsymbol{\lambda} \in \widehat{\Lambda}_{n}(\boldsymbol{\theta})} f(\boldsymbol{\lambda};\boldsymbol{\theta}) + \sum_{k=1}^{p} p_{1,\pi}(|\theta_{k}|).$$

• Then  $\widehat{\boldsymbol{\theta}}_n = \arg \min_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} S_n(\boldsymbol{\theta}).$ 

▶ Let  $\widehat{\lambda}(\theta) = \arg \max_{\lambda \in \widehat{\Lambda}_n(\theta)} f(\lambda; \theta)$  be the Lagrange multiplier defined at  $\theta \in \Theta$ . For any subset  $\mathcal{A} \subset \{1, ..., r\}$ , we denote by  $\mathbf{g}_{\mathcal{A}}(\mathbf{X}_i; \theta)$  the subvector of  $\mathbf{g}(\mathbf{X}_i; \theta)$  with components indexed by  $\mathcal{A}$ . We write  $\overline{\mathbf{g}}_{\mathcal{A}}(\theta) = n^{-1} \sum_{i=1}^{n} \mathbf{g}_{\mathcal{A}}(\mathbf{X}_i; \theta), \widehat{V}_{\mathcal{A}}(\theta) = n^{-1} \sum_{i=1}^{n} \mathbf{g}_{\mathcal{A}}(\mathbf{X}_i; \theta) \mathbf{g}_{\mathcal{A}}(\mathbf{X}_i; \theta)^{\mathrm{T}}$  and  $V_{\mathcal{A}}(\theta_0) = \mathbb{E}\{\mathbf{g}_{\mathcal{A}}(\mathbf{X}_i; \theta_0)\mathbf{g}_{\mathcal{A}}(\mathbf{X}_i; \theta_0)^{\mathrm{T}}\}.$ 

#### Proposition 4

Suppose that Condition 1 hold, and the penalty function  $p_{2,\nu}(\cdot)$  is a convex function. Assume that  $\max_{1 \le j \le r} n^{-1} \sum_{i=1}^{n} |g_j(\mathbf{X}_i; \boldsymbol{\theta}_0)|^2 = O_p(\varpi_n)$  for some  $\varpi_n > 0$  that may diverge with n, and  $\min_{1 \le j \le r} \mathbb{E}\{|g_j(\mathbf{X}_i; \boldsymbol{\theta}_0)|^2\}$  is uniformly bounded away from zero. If  $(n^{-1}\varpi_n \log r)^{1/2}/\nu \to 0$  and  $\log r = o(\min\{n^{1/3}, n^{1-2/\gamma}\varpi_n^{-1}\})$ , then  $\mathbb{P}\{S_n(\boldsymbol{\theta}_0) = \sum_{k=1}^{p} p_{1,\pi}(|\boldsymbol{\theta}_k^0|)\} \to 1$  as  $n \to \infty$  and  $\lambda = \mathbf{0}$  is a local maximizer of  $f(\boldsymbol{\lambda}; \boldsymbol{\theta}_0)$  with probability approaching one.

- Based on the convexity of  $p_{2,\nu}(\cdot)$ ,  $f(\boldsymbol{\lambda};\boldsymbol{\theta})$  is concave w.r.t.  $\boldsymbol{\lambda}$ . To compute  $S_n(\boldsymbol{\theta})$ , we only need to find a local maximizer for  $f(\boldsymbol{\lambda};\boldsymbol{\theta})$ .
- Moreover, for any  $\theta$  such that  $|\theta \theta_0|_{\infty} > \varepsilon_n$ , Condition 2 indicates  $|\mathbb{E}\{\mathbf{g}(\mathbf{X}_i; \theta)\}|_{\infty} \ge \Delta(\varepsilon_n)$ . We show in our proof that when  $\theta$  takes value departing from the truth  $\theta_0$ , i.e.,  $\Delta(\varepsilon_n)$  decays to zero at some slow enough rate,  $S_n(\theta)$  takes a value larger than  $\varrho S_n(\theta_0)$  for some diverging  $\varrho$  with probability tending to 1; see also Chang, Tang and Wu (2013, 2016) for such a phenomenon in marginal EL. This feature, together with the Proposition 4, leads to the consistency of the penalized EL estimator  $\hat{\theta}_n$ .

#### Theorem 1

Let  $a_n = \sum_{k=1}^p p_{1,\pi}(|\theta_k^0|)$ . Assume that conditions in Proposition 4 and Condition 2 hold. Define  $b_n = \max\{n^{-1}, a_n, \nu^2\}$ . If  $b_n = o(n^{-2/\gamma})$ , then  $|\widehat{\theta}_n - \theta_0|_{\infty} = O_p\{b_n^{1/(2\beta)}\}.$ 

Main results – Asymptotic normality

#### Condition 3

Assume that

$$\sup_{\boldsymbol{\theta} \in \{\boldsymbol{\theta} \in \boldsymbol{\Theta}: |\boldsymbol{\theta} - \boldsymbol{\theta}_0|_{\infty} < \epsilon_n b_n^{1/(2\beta)}\}} |\bar{\mathbf{g}}(\boldsymbol{\theta}) - \mathbb{E}\{\mathbf{g}(\mathbf{X}_i; \boldsymbol{\theta})\}|_{\infty} = O_p(\zeta_n)$$

for some  $\zeta_n \to 0$  and  $\epsilon_n \to \infty$ , where  $b_n$  is defined in Theorem 1 and  $\beta$  is specified in Condition 2. Define  $\Theta_n = \{ \theta \in \Theta : |\theta - \theta_0|_{\infty} < \epsilon_n b_n^{1/(2\beta)} \}$ . We also assume that

$$\sup_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_n} \max_{1 \le j \le r} \left\{ \frac{1}{n} \sum_{i=1}^n |g_j(\mathbf{X}_i; \boldsymbol{\theta})|^2 \right\} = O_p(\varpi_n)$$

for some  $\varpi_n > 0$  that may diverge with n.

Define  $\mathcal{M}_{\theta} = \{1 \leq j \leq r : |\mathbb{E}\{g_j(\mathbf{X}_i; \boldsymbol{\theta})\}| \geq \nu \rho'_2(0^+)/2\}$  for  $\boldsymbol{\theta} \in \boldsymbol{\Theta}$ . Proposition 5 below shows that for any  $\boldsymbol{\theta}$  in a small neighborhood of  $\boldsymbol{\theta}_0$ , the support of the Lagrange multiplier  $\widehat{\boldsymbol{\lambda}}(\boldsymbol{\theta})$  is a subset of  $\mathcal{M}_{\theta}$  with probability approaching one.

#### Proposition 5

Let  $\{\theta_n\}$  be a sequence of points in  $\Theta_n$  with  $\Theta_n$  being defined in Condition 3, and the penalty function  $p_{2,\nu}(\cdot)$  is a convex function. Assume Conditions 1 and 3 hold,  $\mathbb{P}[\lambda_{\min}\{\widehat{V}_{\mathcal{M}_{\theta_n}}(\theta_n)\} \geq \xi_n] \rightarrow 1$  for some  $\xi_n > 0$  ( $\xi_n$  may decay to zero as  $n \rightarrow \infty$ ), and  $|\overline{\mathbf{g}}_{\mathcal{M}_{\theta_n}}(\theta_n)|_{\infty} - \nu \rho'_2(0^+) = O_p(u_n)$  for some  $u_n \rightarrow 0$ . If  $\max\{\zeta_n, m_n u_n \varpi_n \xi_n^{-1}\} = o(\nu)$  and  $m_n^{\gamma+1} u_n^{\gamma} n = o(\xi_n^{\gamma})$  where  $m_n = |\mathcal{M}_{\theta_n}|$ , then with probability approaching one there exists a sparse local maximizer  $\widehat{\lambda}(\theta_n)$  for  $f(\lambda; \theta_n)$  satisfying  $|\widehat{\lambda}(\theta_n)|_1 = O_p(m_n u_n \xi_n^{-1})$  and  $\mathbb{P}\{\operatorname{supp}(\widehat{\lambda}(\theta_n)) \subset \mathcal{M}_{\theta_n}\} \rightarrow 1$  as  $n \rightarrow \infty$ .

- Proposition 5 implies that when θ is around θ<sub>0</sub>, the penalty on λ effectively conducts a moments selection by choosing the estimating functions in a way that E{g<sub>j</sub>(X<sub>i</sub>; θ)} has large absolute deviation from 0 when θ ≠ θ<sub>0</sub>.
- ▶ By taking  $\theta_n$  as  $\hat{\theta}_n$ , Proposition 5 shows that with probability approaching one, there exists a sparse Lagrange multiplier  $\hat{\lambda}(\hat{\theta}_n)$  which is a local maximizer for  $f(\lambda; \hat{\theta}_n)$ .
- ▶ We also note here that  $\mathcal{M}_{\hat{\theta}_n}$  is a random set depending on the observations  $\{\mathbf{X}_1, \ldots, \mathbf{X}_n\}$ . To clearly characterize the property of  $\hat{\theta}_n$ , we define  $\mathcal{R}_n \subset \{1, \ldots, r\}$  to be the support of  $\hat{\lambda}(\hat{\theta}_n)$ , i.e., the set of effective estimating functions that ultimately contribute to estimating the nonzero components of  $\hat{\theta}_n$ , and  $q_n = |\mathcal{R}_n|$ .

- By definition, R<sub>n</sub> is a random set depending on multiple sources including the penalty function, the tuning parameter ν, and also the realization of the sample X<sub>1</sub>,..., X<sub>n</sub>. Clearly, the asymptotic variance of θ̂<sub>n</sub> depends on the set R<sub>n</sub> collecting the effective estimating functions for the ultimate estimations.
- For the sake of simplicity in presenting the theoretical development, we assume that the set  $\mathcal{R}_n$  remains invariant when  $\theta$  takes a value in a neighborhood of  $\hat{\theta}_n$ . In other words, the set of estimating functions does not vary in the final steps of optimizing  $S_n(\theta)$ . Actually, such a requirement is not restrictive and is easy to fulfill. In practice, if the set of estimating functions varies substantially in optimizing  $S_n(\theta)$ , one can always drop the penalty on some chosen estimating functions and retain a stable  $\mathcal{R}_n$ , especially when the iterations are closing to convergence.

## Condition 4

There exist uniform constants  $0 < K_3 < K_4$  such that

$$\lim_{n \to \infty} \mathbb{P} \bigg[ K_3 < \inf_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_n} \lambda_{\min} \{ \widehat{V}_{\mathcal{R}_n}(\boldsymbol{\theta}) \} \le \sup_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_n} \lambda_{\max} \{ \widehat{V}_{\mathcal{R}_n}(\boldsymbol{\theta}) \} < K_4 \bigg] = 1$$

where  $\Theta_n$  is defined in Condition 3.

#### Condition 5

(i) There exists a positive constant c < 1 such that

 $\max_{k \in \mathcal{S}} \sup_{c|\theta_k^0| < t < c^{-1}|\theta_k^0|} p_{1,\pi}'(t) = O(\min\{s^{-1}b_n^{-1/(2\beta)}n^{-1}, sq_n^{1/2}\nu^2, q_n^{1/\gamma+1}n^{-1/2+1/\gamma}\nu\}).$ 

(ii) The penalty function  $p_{2,\nu}(\cdot)$  is convex, and there exits a uniform constant  $K_5 > 0$  such that  $|\rho'_2(t;\nu) - \rho'_2(0^+)| \le K_5 t$  for any  $t \in (0,\epsilon)$  and  $\tau > 0$  where  $\epsilon > 0$  is a sufficiently small constant.

#### Condition 6

For each j = 1, ..., p,  $g_j(\mathbf{X}; \boldsymbol{\theta})$  is continuously differentiable with respect to  $\boldsymbol{\theta}$ in  $\Theta_n$  for any  $\mathbf{X}$ , where  $\Theta_n$  is defined in Condition 3, and there exist a function  $B_{n,jk}$  with  $\mathbb{E}\{B_{n,jk}^2(\mathbf{X}_i)\} \leq K_6$  for some uniform constant  $K_6 > 0$ and  $|\partial g_j(\mathbf{X}; \boldsymbol{\theta})/\partial \theta_k| \leq B_{n,jk}(\mathbf{X})$  for any  $\boldsymbol{\theta} \in \Theta_n$ . In addition, it holds that

$$\sup_{\boldsymbol{\theta}\in\boldsymbol{\Theta}_n} \max_{j\in\mathcal{R}_n} \max_{k\notin\mathcal{S}} \left\{ \frac{1}{n} \sum_{i=1}^n \left| \frac{\partial g_j(\mathbf{X}_i; \boldsymbol{\theta})}{\partial \theta_k} \right| \right\} = O_p(\varphi_n)$$

for some  $\varphi_n > 0$  that may diverge with n.

## Condition 7

For each j = 1, ..., p,  $g_j(\mathbf{X}; \boldsymbol{\theta})$  is twice continuously differentiable with respect to  $\boldsymbol{\theta}$  in  $\boldsymbol{\Theta}_n$  for any  $\mathbf{X}$ , where  $\boldsymbol{\Theta}_n$  is defined in Condition 3, and there exist some functions  $B_{n,jkl}$  with  $\mathbb{E}\{B_{n,jkl}^2(\mathbf{X}_i)\} \leq K_7$  for a uniform constant  $K_7 > 0$  and  $|\partial^2 g_j(\mathbf{X}; \boldsymbol{\theta})/\partial \theta_k \partial \theta_l| \leq B_{n,jkl}(\mathbf{X})$  for any  $\boldsymbol{\theta} \in \boldsymbol{\Theta}_n$ . In addition, for  $\mathcal{R}_n$  given in Condition 4, it holds that

$$\lim_{n \to \infty} \mathbb{P} \bigg[ \inf_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_n} \lambda_{\min}( [\nabla_{\boldsymbol{\theta}_{(1)}} \bar{\mathbf{g}}_{\mathcal{R}_n}(\boldsymbol{\theta})]^{\mathrm{T}} [\nabla_{\boldsymbol{\theta}_{(1)}} \bar{\mathbf{g}}_{\mathcal{R}_n}(\boldsymbol{\theta})]) > K_8 \bigg] = 1$$

where  $\theta_{(1)}$  is an s-dimensional subvector of  $\theta$  with components indexed by  $S = \{1 \le k \le p : \theta_k^0 \ne 0\}$  and  $K_8 > 0$  is a uniform constant.

# We define

$$\begin{split} \mathbf{J} &= [\mathbb{E}\{\nabla_{\boldsymbol{\theta}_{(1)}}\mathbf{g}_{\mathcal{R}_n}(\mathbf{X}_i;\boldsymbol{\theta}_0)\}]^{\mathrm{T}}V_{\mathcal{R}_n}^{-1}(\boldsymbol{\theta}_0)[\mathbb{E}\{\nabla_{\boldsymbol{\theta}_{(1)}}\mathbf{g}_{\mathcal{R}_n}(\mathbf{X}_i;\boldsymbol{\theta}_0)\}],\\ & \widehat{\boldsymbol{\psi}}_{\mathcal{R}_n} = \widehat{\mathbf{J}}^{-1}[\nabla_{\boldsymbol{\theta}_{(1)}}\overline{\mathbf{g}}_{\mathcal{R}_n}(\widehat{\boldsymbol{\theta}}_n)]^{\mathrm{T}}\widehat{V}_{\mathcal{R}_n}^{-1}(\widehat{\boldsymbol{\theta}}_n) \bigg[\frac{1}{n}\sum_{i=1}^n \frac{\mathbf{g}_{\mathcal{R}_n}(\mathbf{X}_i;\widehat{\boldsymbol{\theta}}_n)}{1+\widehat{\boldsymbol{\lambda}}(\widehat{\boldsymbol{\theta}}_n)^{\mathrm{T}}\mathbf{g}(\mathbf{X}_i;\widehat{\boldsymbol{\theta}}_n)}\bigg]\\ & \text{with } \widehat{\mathbf{J}} \ = \ [\nabla_{\boldsymbol{\theta}_{(1)}}\overline{\mathbf{g}}_{\mathcal{R}_n}(\widehat{\boldsymbol{\theta}}_n)]^{\mathrm{T}}\widehat{V}_{\mathcal{R}_n}^{-1}(\widehat{\boldsymbol{\theta}}_n)[\nabla_{\boldsymbol{\theta}_{(1)}}\overline{\mathbf{g}}_{\mathcal{R}_n}(\widehat{\boldsymbol{\theta}}_n)]. \ \text{We have the following asymptotic properties for } \widehat{\boldsymbol{\theta}}_n. \end{split}$$

(日) (個) (三) (三) (三) (○) (○)

#### Theorem 2

Let  $\hat{\theta}_{n,(1)}$  and  $\hat{\theta}_{n,(2)}$  be the corresponding estimations of  $\theta_{0,(1)}$  and  $\theta_{0,(2)}$ , respectively. Then the following properties hold.

- ► Under Conditions 1–6, if  $\log r = o(\min\{n^{1/3}, n^{1-2/\gamma}\varpi_n^{-1}\})$ ,  $q_n = o(n^{(\gamma-2)/(2\gamma+2)})$  $b_n = o(\min_{k \in S} |\theta_k^0|^{2\beta})$ , and the tuning parameters  $\pi$  and  $\nu$  satisfy the conditions that  $q_n n^{-1/2} \varphi_n = o(\pi)$  and  $\max\{\zeta_n, q_n n^{-1/2} \varpi_n, (n^{-1} \varpi_n \log r)^{1/2}\} = o(\nu)$ , then  $\mathbb{P}\{\widehat{\theta}_{n,(2)} = \mathbf{0}\} \to 1$  as  $n \to \infty$ .
- Assume the eigenvalues of  $[\mathbb{E}\{\nabla_{\theta_{(1)}}\mathbf{g}_{\mathcal{R}_n}(\mathbf{X}_i;\theta_0)\}]^{\mathrm{T}}[\mathbb{E}\{\nabla_{\theta_{(1)}}\mathbf{g}_{\mathcal{R}_n}(\mathbf{X}_i;\theta_0)\}]$ and  $V_{\mathcal{R}_n}(\theta_0)$  are uniformly bounded away from zero and infinity. Under Conditions 1–7, for any  $\alpha \in \mathbb{R}^s$  with unit  $L_2$ -norm, if  $\log r = o(\min\{n^{1/3}, n^{1-2/\gamma}\varpi q_n = o(n^{(\gamma-2)/(2\gamma+2)}), b_n = o(\min_{k \in S} |\theta_k^0|^{2\beta}), and \pi and \nu$  satisfy the conditions that  $q_n n^{-1/2} \varphi_n = o(\pi), \max\{\zeta_n, q_n n^{-1/2} \varpi_n, (n^{-1} \varpi_n \log r)^{1/2}\} = o(\nu)$  and  $\nu = o(\min\{s^{-1/2}q_n^{-1/2}n^{-1/4}, q_n^{-1/\gamma-3/2}n^{-1/\gamma}\})$ , then

$$n^{1/2} \boldsymbol{\alpha}^{\mathrm{T}} \mathbf{J}^{1/2} \{ \widehat{\boldsymbol{\theta}}_{n,(1)} - \boldsymbol{\theta}_{0,(1)} - \widehat{\boldsymbol{\psi}}_{\mathcal{R}_n} \} \xrightarrow{d} N(0,1)$$

as  $n \to \infty$ .

- The first part of Theorem 2 indicates that the zero components of  $\theta_0$  can be exactly estimated as zero by  $\hat{\theta}_n$  with probability approaching one. The bias term  $\hat{\psi}_{\mathcal{R}_n}$  is due to the use of the penalty function  $p_{2,\nu}(\cdot)$ . It can be shown that  $\hat{\psi}_{\mathcal{R}_n} = \hat{\mathbf{J}}^{-1} [\nabla_{\theta_{(1)}} \bar{\mathbf{g}}_{\mathcal{R}_n}(\hat{\theta}_n)]^{\mathrm{T}} \hat{V}_{\mathcal{R}_n}^{-1}(\hat{\theta}_n) \hat{\eta}_{\mathcal{R}_n}$  where  $\hat{\boldsymbol{\eta}} = (\hat{\eta}_1, \dots, \hat{\eta}_r)^{\mathrm{T}}$  with  $\hat{\eta}_j = \nu \rho'_2(|\hat{\lambda}_j|; \nu) \mathrm{sgn}(\hat{\lambda}_j)$  for  $\hat{\lambda}_j \neq 0$  and  $\hat{\eta}_j \in [-\nu \rho'_2(0^+), \nu \rho'_2(0^+)]$  for  $\hat{\lambda}_j = 0$ .
- The second part suggests that the the rate of the convergence is n<sup>1/2</sup> for each nonzero component of θ<sub>0</sub>.

# Algorithm

For ease and stability in implementations, we calculate the penalized EL estimator θ̂<sub>n</sub> by minimizing the following slightly modified objective function:

$$\widehat{\boldsymbol{\theta}}_{n} = \arg\min_{\boldsymbol{\theta}\in\boldsymbol{\Theta}} \max_{\boldsymbol{\lambda}\in\widehat{\Lambda}_{n}(\boldsymbol{\theta})} \bigg[ \underbrace{\sum_{i=1}^{n} \log_{\star}\{1 + \boldsymbol{\lambda}^{\mathrm{T}} \mathbf{g}(\mathbf{X}_{i}; \boldsymbol{\theta})\} - n \sum_{j=1}^{r} p_{2,\nu}(|\lambda_{j}|)}_{f(\boldsymbol{\lambda}; \boldsymbol{\theta})} + n \sum_{k=1}^{p} p_{1,\pi}(|\theta_{k}|) \bigg],$$

where  $\log_{\star}(z)$  is a twice differentiable pseudo-logarithm function with bounded support adopted from Owen (2001):

$$\log_{\star}(z) = \begin{cases} \log(z) & \text{if } z \ge \epsilon;\\ \log(\epsilon) - 1.5 + 2z/\epsilon - z^2/(2\epsilon^2) & \text{if } z \le \epsilon; \end{cases}$$

where  $\epsilon$  is chosen as 1/n in our implementations.

In the optimization, we apply the quadratic approximation (Fan and Li, 2001) to the penalty functions p<sub>1,π</sub>(·) and p<sub>2,ν</sub>(·). More specifically, for a penalty function p<sub>τ</sub>(·), the quadratic approximation states

$$p_{\tau}(|t|) \approx p_{\tau}(|t_0|) + \frac{1}{2} \frac{p'_{\tau}(|t_0|)}{|t_0|} (t^2 - t_0^2)$$
(1)

for t being in a small neighborhood of  $t_0$ . The first and second derivatives are approximated by

$$p'_{\tau}(|t|) \approx \frac{p'_{\tau}(|t_0|)}{|t_0|} \cdot t \text{ and } p''_{\tau}(|t|) \approx \frac{p'_{\tau}(|t_0|)}{|t_0|}$$

• The computation of EL is a challenging aspect, especially with high-dimensional p and r. To compute the penalized EL estimator  $\hat{\theta}_n$ , we propose to apply a modified two-layer coordinate decent algorithm. The inner layer of the algorithm solves for  $\lambda$  with given  $\theta$  by maximizing  $f(\lambda; \theta)$ . This layer only involves maximizing a concave function, and hence is stable. The outer layer of the algorithm searches for the optimizer  $\hat{\theta}_n$ . Both layers can be solved using coordinate descent by cycling through and updating each of the coordinates; see Tang and Wu (2014).

▶ In the inner layer,  $\lambda$  is solved at a given  $\theta$ , which can be done by optimizing  $f(\lambda; \theta)$  with respect to  $\lambda$  using coordinate descent. Suppose that  $\lambda$  starts at an initial value  $\hat{\lambda}^{(0)}$ . With the other coordinates fixed, the (m + 1)th Newton's update for  $\lambda_j$  (j = 1, ..., r), the *j*th component of  $\lambda$ , is given by

$$\widehat{\lambda}_{j}^{(m+1)} = \widehat{\lambda}_{j}^{(m)} - \frac{\sum_{i=1}^{n} \log_{\star}'(t_{i}^{(m)}) g_{j}(\mathbf{X}_{i}; \boldsymbol{\theta}) - n p_{2,\nu}'(|\widehat{\lambda}_{j}^{(m)}|)}{\sum_{i=1}^{n} \log_{\star}''(t_{i}^{(m)}) \{g_{j}(\mathbf{X}_{i}; \boldsymbol{\theta})\}^{2} - n p_{2,\nu}'(|\widehat{\lambda}_{j}^{(m)}|)},$$

where  $t_i^{(m)} = 1 + \mathbf{g}(\mathbf{X}_i; \boldsymbol{\theta})^{\mathrm{T}} \widehat{\boldsymbol{\lambda}}^{(m)}$  with  $\widehat{\boldsymbol{\lambda}}^{(m)} = (\widehat{\lambda}_1^{(m)}, \dots, \widehat{\lambda}_r^{(m)})^{\mathrm{T}}$ . The procedure cycles through all the r components of  $\boldsymbol{\lambda}$  and is repeated until convergence.

During this process, the objective function needs to be checked to ensure it gets optimized in each step. If not, the step size continues to be halved until the objective function gets driven in the right direction.

▶ The outer layer of the algorithm is to optimize the objective function with respect to the parameter  $\theta$ , the main interest of the penalized EL, using the coordinate descent algorithm. At a given  $\lambda$ , the algorithm updates  $\theta_k$  (k = 1, ..., p), by minimizing  $S_n(\theta)$  with respect to  $\theta_k$  with other  $\theta_l$   $(l \neq k)$  fixed. Suppose that  $\theta$  starts at an initial value  $\hat{\theta}^{(0)}$ . The (m+1)th update for  $\theta_k$  is given by

$$\begin{split} \widehat{\theta}_{k}^{(m+1)} &= \widehat{\theta}_{k}^{(m)} - \frac{\sum_{i=1}^{n} \log_{\star}'(s_{i}^{(m)}) w_{ik}^{(m)} + np'_{1,\tau}(|\widehat{\theta}_{k}^{(m)}|)}{\sum_{i=1}^{n} [\log_{\star}''(s_{i}^{(m)}) \{w_{ik}^{(m)}\}^{2} + \log_{\star}'(s_{i}^{(m)}) z_{ik}^{(m)}] + np''_{1,\tau}(|\widehat{\theta}_{k}^{(m)}|)}, \\ \text{where } s_{i}^{(m)} &= 1 + \lambda^{\mathrm{T}} \mathbf{g}(\mathbf{X}_{i}; \widehat{\boldsymbol{\theta}}^{(m)}), w_{ik}^{(m)} = \lambda^{\mathrm{T}} \partial \mathbf{g}(\mathbf{X}_{i}; \widehat{\boldsymbol{\theta}}^{(m)}) / \partial \theta_{k} \text{ and } z_{ik}^{(m)} = \\ \lambda^{\mathrm{T}} \partial^{2} \mathbf{g}(\mathbf{X}_{i}; \widehat{\boldsymbol{\theta}}^{(m)}) / \partial \theta_{k}^{2} \text{ with } \widehat{\boldsymbol{\theta}}^{(m)} = (\widehat{\theta}_{1}^{(m)}, \dots, \widehat{\theta}_{p}^{(m)})^{\mathrm{T}}. \end{split}$$

Since quadratic approximations are applied in the algorithms, we follow Fan and Li (2001) and set a component λ<sup>(m)</sup><sub>j</sub> or θ<sup>(m)</sup><sub>k</sub> as zero when it is less than a threshold level say 10<sup>-3</sup> in an iteration.

We summarize the computation procedure for θ and λ in the following pseudo-code. Suppose ξ is a pre-defined small number, say, ξ = 10<sup>-4</sup>.

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

# Numerical results

- ▶ Linear regression model  $Y_i = \mathbf{Z}_i^{\mathrm{T}} \boldsymbol{\theta}_0 + \varepsilon_i$ , where  $\boldsymbol{\theta}_0 = (3, 1.5, 0, 0, 2, 0, \dots, 0)^{\mathrm{T}}$ ,  $\mathbf{Z}_i \in \mathbb{R}^p$  are generated from  $N(\mathbf{0}, \mathbf{\Sigma})$  with  $\sigma_{kk} = 1$  for any  $k = 1, \dots, p$ and  $\sigma_{kl} = 0.5$  for any  $k \neq l$ , where  $\mathbf{\Sigma} = (\sigma_{kl})_{p \times p}$ , and  $\varepsilon_i$  is a standard normal distributed random variable. Write  $\mathbf{X}_i = (Y_i, \mathbf{Z}_i^{\mathrm{T}})^{\mathrm{T}}$ . The estimating function is  $\mathbf{g}(\mathbf{X}; \boldsymbol{\theta}) = \mathbf{Z}(Y - \mathbf{Z}^{\mathrm{T}} \boldsymbol{\theta})$  with p = r.
- The proposed penalized EL with two penalties (namely, PEL2) is compared to the single penalty approach (PEL) discussed in Tang and Leng (2010). Three information criteria for choosing the tuning parameters π and ν in the penalty function – BIC (Schwarz, 1978), BICC (Wang, Li and Leng, 2009), and EBIC (Chen and Chen, 2008) – are used.

• The results from MLE for the three true variables (i.e., MLE-Oracle) is also reported. The column of  $\theta_{nonzero}$  reports the average number of selected predictors. The column of  $\theta_{true}$  reports the average number of true predictors that are selected. The difference is the average number of false predictors that get selected. The next column reports the model error (ME), which is defined by  $ME = (\hat{\theta} - \theta)^T (\hat{\theta} - \theta)$  for a given estimator  $\hat{\theta}$ .

-					
(n, p, r)	Method	$ heta_{ m nonzeros}$	$ heta_{ ext{true}}$	ME	No. EE's
n = 50	MLE-Oracle	3 (0)	NA	0.104 (0.009)	NA
p = 100	PEL-BIC	0 (0)	0 (0)	15.25 (0)	100 (0)
r = 100	PEL-BIC	0 (0)	0 (0)	15.25 (0)	100 (0)
	PEL-EBIC	0 (0)	0 (0)	15.25 (0)	100 (0)
	PEL2-BIC	4.94 (0.53)	2.92 (0.04)	0.988 (0.160)	5.73 (0.21)
	PEL2-BIC	4.73 (0.48)	2.92 (0.04)	1.025 (0.165)	5.69 (0.22)
	PEL2-EBIC	4.42 (0.43)	2.90 (0.04)	1.239 (0.211)	5.63 (0.23)
n = 100	MLE-Oracle	3 (0)	NA	0.047 (0.005)	NA
p = 200	PEL-BIC	0 (0)	0 (0)	15.25 (0)	200 (0)
r = 200	PEL-BIC	0 (0)	0 (0)	15.25 (0)	200 (0)
	PEL-EBIC	0 (0)	0 (0)	15.25 (0)	200 (0)
	PEL2-BIC	9.22 (1.27)	3 (0)	1.070 (0.225)	5.38 (0.17)
	PEL2-BIC	9.28 (1.28)	3 (0)	1.079 (0.227)	5.39 (0.17)
	PEL2-EBIC	8.38 (1.03)	3 (0)	1.056 (0.228)	5.34 (0.17)
n = 100	MLE-Oracle	3 (0)	NA	0.039 (0.003)	NA
p = 500	PEL-BIC	0 (0)	0 (0)	15.25 (0)	500 (0)
r = 500	PEL-BIC	0 (0)	0 (0)	15.25 (0)	500 (0)
	PEL-EBIC	0 (0)	0 (0)	15.25 (0)	500 (0)
	PEL2-BIC	6.28 (1.31)	3 (0)	0.946 (0.153)	5.48 (0.16)
	PEL2-BIC	5.96 (1.31)	3 (0)	0.930 (0.155)	5.38 (0.17)
		6.04 (1.32)	3 (0)		5 11 (0 16) °

## Discussion

- We study a new penalized EL approach with two penalties, with one encouraging sparsity of the estimator and the other encouraging sparsity of the Lagrange multiplier in the optimizations associated with the EL. Such an approach utilizes sparsity in the target parameters and effectively achieves a moment selection procedure for estimating the sparse parameter. Both theory and numerical examples confirm the merits of the new penalized EL.
- One interesting extension of the approach is to explore inferences with estimating equations after the variable selection procedure. Such a direction is a suitable stage for EL method with estimating equations who takes advantage of adaptivity to various moment conditions with less stringent distributional assumptions. (Done)
- The other interesting and challenging problem is to explore the optimality of the sparse estimator using estimating equations with high data dimensionality. (Working on)

# Thank you!

◆□ ▶ < 圖 ▶ < 圖 ▶ < 圖 ▶ < 圖 • 의 Q @</p>