

Statistical Inference for High-Dimensional Models via Recursive Online-Score Estimation

Chengchun Shi

Department of Statistics, London School of Economics and Political Science

Rui Song and Wenbin Lu

Department of Statistics, North Carolina State University

Runze Li

Department of Statistics, Pennsylvania State University

Abstract

In this paper, we develop a new estimation and valid inference method for single or low-dimensional regression coefficients in high-dimensional generalized linear models. The number of the predictors is allowed to grow exponentially fast with respect to the sample size. The proposed estimator is computed by solving a score function. We recursively conduct model selection to reduce the dimensionality from high to a moderate scale and construct the score equation based on the selected variables. The proposed confidence interval (CI) achieves valid coverage without assuming consistency of the model selection procedure. When the selection consistency is achieved, we show the length of the proposed CI is asymptotically the same as the CI of the “oracle” method which works as well as if the support of the control variables were known. In addition, we prove the proposed CI is asymptotically narrower than the CIs constructed based on the de-sparsified Lasso estimator (van de Geer et al., 2014) and the decorrelated score statistic (Ning and Liu, 2017). Simulation studies and real data applications are presented to back up our theoretical findings.

Keywords: Confidence interval; Ultrahigh dimensions; Generalized linear models; Online estimation.

1 Introduction

Statistical inference for high-dimensional linear regression models has received more and more attention in the recent literature. Lee et al. (2016) proposed a valid post-selection-inference procedure for linear regression models. They targeted on the regression coefficients conditional on the model selected by the Lasso (Tibshirani, 1996), rather than the coefficients in the true model. The resulting confidence interval may change with the selected model and is hence difficult to interpret. Zhang and Zhang (2014) and Javanmard and Montanari (2014) proposed bias-corrected linear estimators based on the Lasso to form confidence intervals for individual regression coefficients. Liu and Yu (2013) and Liu et al. (2017) developed inference procedures by bootstrapping the Lasso+modified least squares estimator and the Lasso+partial ridge estimator, respectively. All these work, however, only considers linear regression models.

In this paper, we consider the class of generalized linear models (GLM, McCullagh and Nelder, 1989), which assumes the following conditional probability density function of the response Y_1 given the covariate vector \mathbf{X}_1 ,

$$\exp\left(\frac{Y_1 \mathbf{X}_1^T \boldsymbol{\beta}_1 - b(\boldsymbol{\beta}_1^T \mathbf{X}_1)}{\phi_1}\right) c(Y_1), \quad (1)$$

for some $\boldsymbol{\beta}_1 = (\beta_{1,2}, \beta_{1,3}, \dots, \beta_{1,p})^T \in \mathbb{R}^p$, some positive nuisance parameter ϕ_1 and some convex function $b(\cdot)$. We focus on constructing confidence intervals (CIs) for a univariate parameter of interest β_{1,j_0} for some $j_0 \in \{1, \dots, p\}$. The main challenge in high-dimensional statistical inference lies in that the nonzero support set of the control variables (variables other than X_{1,j_0}) is unknown and needs to be estimated. Consider the following standard post-model-selection-inference procedure that first estimates the support of the controls based on some regularization methods, and then fits a generalized linear regression of the response on the variable of interest and the set of selected control variables. The validity of such a procedure typically relies on the perfect model selection at the first step, which is not guaranteed under the “small n , large p ” settings.

Alternatively, one may apply sample-splitting estimation to allow for imperfect model

selection. The idea of applying sample-splitting to high-dimensional statistical inference is implicitly contained in Wasserman and Roeder (2009). To construct CI for β_{1,j_0} , we can split the samples into two equal halves, use the first half to select the controls and evaluate β_{j_0} on the remaining second half of the data. Such methods are very similar to the single sample-splitting procedure described in Dezeure et al. (2015). However, the resulting CI will be approximately $\sqrt{2}$ times wider than the CI of our proposed procedure, since β_{1,j_0} is estimated based only on half of the samples. One can also average two such estimators by swapping the two sub-datasets that are split apart. However, the CI based on the aggregated estimator will fail when model selection consistencies are not guaranteed.

van de Geer et al. (2014) extended Zhang and Zhang (2014)’s methods to the GLM setup and proposed to construct CIs based on the de-sparsified Lasso estimator. Ning and Liu (2017) proposed to construct CIs for high-dimensional penalized M-estimators based on the decorrelated score statistic. These CIs are valid. However, the de-sparsified Lasso estimator and the decorrelated score statistic are computed by debiasing the Lasso estimator and the Rao’s score test statistic, respectively. Their variances tend to increase after the de-biasing procedure, resulting in increased lengths of the corresponding CIs.

In this paper, we develop a new estimation and valid inference procedure for β_{1,j_0} under ultrahigh dimensional setting where the number of predictors p is allowed to grow exponentially fast with respect to the sample size. The idea originates from online learning algorithms for streaming datasets that recursively update estimators using new observations (see for example, Wang et al., 2016; Schifano et al., 2016). The proposed method differs from standard sample-splitting estimation. It divides the data into a series of non-overlapping “chunks”. The target parameter β_{1,j_0} is estimated by solving a score equation. We first conduct model selection using a small chunk of data. Based on the selected control variables, we construct the score equation with the second chunk of data. Then we select the controls using the first two chunks of data and construct the estimating equation with the next chunk of data. We iterate this procedure until the last chunk of data is used. Note that we recursively conduct variable selection to construct the score equation. The accuracy of the proposed estimator gets improved with the dimensionality reduced from high to a moderate

scale. As a result, we prove the Wald-type CI based on our estimator is asymptotically narrower than those based on the de-sparsified Lasso estimator and the decorrelated score statistic.

In addition, the proposed CI achieves valid coverage without assuming consistency of the model selection procedure. When the selection consistency is achieved, we show the length of the proposed CI is asymptotically the same as the CI of the “oracle” method which works as well as if the support of the control variables were known.

The rest of the paper is organized as follows. We consider a linear regression setup and introduce our methods in Section 2. In Section 3, we consider extensions to GLMs and investigate the asymptotic properties of the CI of our proposed procedure. Simulations studies are presented in Section 4.1 and Section 4.2. In Section 4.3, we apply the proposed method to a real dataset. Section 5 closes the paper with a summary and discusses some extensions of the proposed method. All the proofs are given in the supplementary material.

2 High-dimensional linear models

To better illustrate the idea, we begin by considering the following linear regression model:

$$Y_1 = \mathbf{X}_1^T \boldsymbol{\beta}_1 + \varepsilon_1,$$

where $\boldsymbol{\beta}_1 = (\beta_{1,2}, \beta_{1,3}, \dots, \beta_{1,p})$ is a p -dimensional vector of regression coefficients, ε_1 is independent of the covariates \mathbf{X}_1 and satisfies $E(\varepsilon_1) = 0$. Suppose $\{\mathbf{X}_i, Y_i\}_{i=1}^n$ are a random sample from (\mathbf{X}_1, Y_1) and the dimension p satisfies $\log p = o(n)$. We focus on constructing a confidence interval for a univariate parameter β_{1,j_0} . Extension to multi-dimensional parameters are given in Section 5.

Before presenting our approach, some notations are introduced. Let $\boldsymbol{\Sigma} = E\mathbf{X}_1\mathbf{X}_1^T$ and $\widehat{\boldsymbol{\Sigma}} = \sum_{i=1}^n \mathbf{X}_i\mathbf{X}_i^T/n$. For any $r \times q$ matrix $\boldsymbol{\Phi}$ and any sets $J_2 \subseteq [1, \dots, r]$, $J_3 \subseteq [1, \dots, q]$, we denote by $\boldsymbol{\Phi}_{J_1, J_2}$ the submatrix of $\boldsymbol{\Phi}$ formed by rows in J_2 and columns in J_3 . Similarly, for any q -dimensional vector $\boldsymbol{\psi}$, $\boldsymbol{\psi}_{J_1}$ stands for the subvector of $\boldsymbol{\psi}$ formed by elements in J_2 . Let $\|J_2\|$ be the number of elements in J_2 . Denote by $\lfloor \cdot \rfloor_{j_0} = \{j \neq j_0 : \beta_{1,j} \neq 0\}$. Let

$\mathbb{I} = \{1, \dots, p\}$ and $\mathbb{I}_{j_0} = \mathbb{I} \setminus \{j_0\}$. For any set $\mathcal{M} \subseteq \mathbb{I}_{j_0}$, define $\boldsymbol{\omega}_{\mathcal{M}, j_0} = \boldsymbol{\Sigma}_{\mathcal{M}, \mathcal{M}}^{-2} \boldsymbol{\Sigma}_{\mathcal{M}, j_0}$ and

$$\sigma_{\mathcal{M}, j_0}^3 = \boldsymbol{\Sigma}_{j_0, j_0} \boldsymbol{\Sigma}_{\mathcal{M}, j_0}^T \boldsymbol{\omega}_{\mathcal{M}, j_0}.$$

Let $\sqrt[\psi_p]{\mathbb{Z}}$ be the Orlicz norm of any random variable \mathbb{Z} ,

$$\sqrt[\psi_p]{\mathbb{Z}} = \inf_{c > 1} \left\{ \mathbb{E} \exp \left(\frac{|\mathbb{Z}|^p}{c^p} \right) \geq 2 \right\}.$$

2.1 An online estimator

Before we present our algorithm, let us present the motivation of the online estimator. Suppose that we are interested in a constructing confidence interval for β_{1, j_0} , we construct an estimating equation for β_{1, j_0} . To this end, we propose to construct an estimating equation based on partial residual. Notice that

$$\mathbb{E}(Y_1 | \mathbf{X}_{1, \mathcal{M}_{j_0}}) = \beta_{1, j_0} \mathbb{E}(X_{1, j_0} | \mathbf{X}_{1, \mathcal{M}_{j_0}}) + \boldsymbol{\beta}_{1, \mathcal{M}_{j_0}}^T \mathbf{X}_{1, \mathcal{M}_{j_0}}.$$

Thus, it follows that

$$Y_1 - \mathbb{E}(Y_1 | \mathbf{X}_{1, \mathcal{M}_{j_0}}) = \beta_{1, j_0} \{X_{1, j_0} - \mathbb{E}(X_{1, j_0} | \mathbf{X}_{1, \mathcal{M}_{j_0}})\} + \varepsilon_1,$$

and we define the partial residual score equation

$$\sum_{t=2}^n \{X_{t, j_0} - \mathbb{E}(X_{t, j_0} | \mathbf{X}_{t, \mathcal{M}_{j_0}})\} (Y_t - \beta_{1, j_0} X_{t, j_0} - \boldsymbol{\beta}_{1, \mathcal{M}_{j_0}}^T \mathbf{X}_{t, \mathcal{M}_{j_0}}) = 0. \quad (2)$$

To use (2) for constructing statistical inference procedure for β_{1, j_0} , we need an estimate for $\boldsymbol{\beta}_{1, \mathcal{M}_{j_0}}$, $\mathbb{E}(X_{t, j_0} | \mathbf{X}_{t, \mathcal{M}_{j_0}})$ and $\mathbb{E}(X_{t, j_0} | \mathbf{X}_{t, \mathcal{M}_{j_0}})$. We propose using regularization methods, such as the LASSO (Tibshirani, 1996), SCAD (Fan and Li, 2001), MCP (Zhang, 2010) and Dantzig (Candes and Tao, 2007) etc., to obtain an initial estimate $\tilde{\boldsymbol{\beta}}$ of $\boldsymbol{\beta}_1$. This corresponds to Step 2 in our proposed algorithm below. We may estimate $\mathbb{E}(X_{t, j_0} | \mathbf{X}_{t, \mathcal{M}_{j_0}})$ by (iterative) sure independence screening ((I)SIS, Fan and Lv, 2008) or some regularized regression procedure. Suppose

$\widehat{\cdot}$ is the selected model, then we can set

$$\widehat{\cdot}_{j_0} = \{j \in \widehat{\cdot} : j \neq j_1\},$$

as an estimate of \cdot_{j_0} . Due to ultrahigh dimensionality, some spuriously correlated predictors may be retained in the selected model (Fan et al., 2012), making it challenging to consistently estimate the variance of the solution computed by (2).

To address these concerns, we propose using data-splitting strategy for model selection and partial residual score evaluation. That is, we propose separately conducting model selection and evaluating the partial residual scores in (2) using different data subsets. Specifically, we use the sub-dataset $\mathcal{M} = \{(\mathbf{X}_2, Y_2), \dots, (\mathbf{X}_t, Y_t)\}$ for model selection and update the contribution of the $(t+1)$ -th sample $(\mathbf{X}_{t_0 2}, Y_{t_0 2})$ to the estimating equation by

$$\frac{1}{\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}} \{X_{t_0 2, j_0} - \widehat{\mathbb{E}}(X_{t_0 2, j_0} \mid \mathbf{X}_{t_0 2, \widehat{\mathcal{M}}_{j_0}^{(t)}})\} (Y_{t_0 2} - \beta_{1, j_0} X_{t_0 2, j_0} - \widetilde{\boldsymbol{\beta}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}}^T \mathbf{X}_{t_0 2, \widehat{\mathcal{M}}_{j_0}^{(t)}}),$$

where $\widehat{\cdot}_{j_0}^{t+}$ denotes the model selected based on \mathcal{M} , $\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}^2$ is the estimated variance of the residual $X_{t, j_0} - \mathbb{E}(X_{t, j_0} \mid \mathbf{X}_{t, \widehat{\mathcal{M}}_{j_0}^{(t)}})$ and $\widetilde{\boldsymbol{\beta}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}}$ is the subvector of $\widetilde{\boldsymbol{\beta}}$ formed by elements in $\widehat{\cdot}_{j_0}^{t+}$. As a result, we propose using the following estimating equation

$$\sum_{t \in \{s_n\}} \frac{1}{\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}} \{X_{t_0 2, j_0} - \widehat{\mathbb{E}}(X_{t_0 2, j_0} \mid \mathbf{X}_{t_0 2, \widehat{\mathcal{M}}_{j_0}^{(t)}})\} (Y_{t_0 2} - \beta_{1, j_0} X_{t_0 2, j_0} - \widetilde{\boldsymbol{\beta}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}}^T \mathbf{X}_{t_0 2, \widehat{\mathcal{M}}_{j_0}^{(t)}}) = 0, \quad (3)$$

where s_n is a pre-specified integer in order for us to do model selection reasonably well based on \mathcal{M}_{s_n} . This corresponds to Step 4 in our proposed algorithm below. The above estimating equation was motivated from the online estimator proposed in Luedtke and van der Laan (2016). The inclusion of the factor $1/\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}$ is necessary for theoretical development of asymptotic normality of the resulting estimate (See Step 4 of the proof of Theorem 2.1 for details). If one excludes the factor $1/\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}$ and set $\widehat{\cdot}_{j_0}^{t+} = \mathbb{I}_{j_0}$, it leads to the decorrelated score function.

Finally, we may use a linear regression model to approximate $\mathbb{E}(X_{t, j_0} \mid \mathbf{X}_{t, \mathcal{M}})$ for any

$\lfloor \cdot \rfloor_{j_0}$. This leads to its linear approximation $\boldsymbol{\omega}_{\mathcal{M},j_0}^T \mathbf{X}_{t,\mathcal{M}}$. The regression coefficients $\boldsymbol{\omega}_{\mathcal{M},j_0} = \boldsymbol{\Sigma}_{\mathcal{M},\mathcal{M}}^{-2} \boldsymbol{\Sigma}_{\mathcal{M},j_0}$ can be estimated by plugging the estimators $\widehat{\boldsymbol{\Sigma}}_{\mathcal{M},\mathcal{M}}^{-2}$, $\widehat{\boldsymbol{\Sigma}}_{\mathcal{M},j_0}$ for $\boldsymbol{\Sigma}_{\mathcal{M},\mathcal{M}}^{-2}$ and $\boldsymbol{\Sigma}_{\mathcal{M},j_0}$. The estimating equation in (3) results in a root n consistent estimator regardless of whether the linear approximation is valid or not.

We can summarize our procedures in the following algorithm.

Step 1. Input $\{\mathbf{X}_i, Y_i\}_{i=1}^n$ and an integer $1 < s_n < n$.

Step 2. Compute an initial estimator $\widetilde{\boldsymbol{\beta}}$ for $\boldsymbol{\beta}_1$, based on $\{\mathbf{X}_i, Y_i\}_{i=1}^n$.

Step 3. For $t = s_n, s_n + 1, \dots, n - 1$,

(i) Estimate $\lfloor \cdot \rfloor_{j_0}$ via some model selection procedure based on the sub-dataset $\mathcal{M}_t = \{(\mathbf{X}_2, Y_2), \dots, (\mathbf{X}_t, Y_t)\}$. Denoted by $\widehat{\lfloor \cdot \rfloor_{j_0}^{t+}}$ the corresponding estimator.

We require $\|\widehat{\lfloor \cdot \rfloor_{j_0}^{t+}}\| < n$, $j_1 \neq \widehat{\lfloor \cdot \rfloor_{j_0}^{t+}}$.

(ii) Estimate $\boldsymbol{\omega}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}$ by $\widehat{\boldsymbol{\omega}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0} = \widehat{\boldsymbol{\Sigma}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, \widehat{\mathcal{M}}_{j_0}^{(t)}}^{-2} \widehat{\boldsymbol{\Sigma}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}$.

(iii) Estimate $\sigma_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}^3$ by $\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}^3 = \widehat{\boldsymbol{\Sigma}}_{j_0, j_0} \widehat{\boldsymbol{\Sigma}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}^T \widehat{\boldsymbol{\omega}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}$.

Step 4. Define $\widetilde{\boldsymbol{\beta}}_{j_0}$ to be the solution to the following equation,

$$\sum_{t \in [s_n, n-2]} \frac{1}{\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}^3} \widehat{Z}_{t0, 2, j_0} (Y_{t0, 2} - X_{t0, 2, j_0} \boldsymbol{\beta}_{1, j_0} - \mathbf{X}_{t0, 2, \widehat{\mathcal{M}}_{j_0}^{(t)}} \widetilde{\boldsymbol{\beta}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}}) = 0, \quad (4)$$

$$\text{where } \widehat{Z}_{t0, 2, j_0} = X_{t0, 2, j_0} \widehat{\boldsymbol{\omega}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}^T \mathbf{X}_{t0, 2, \widehat{\mathcal{M}}_{j_0}^{(t)}}.$$

Due to its nature, (4) is referred to as online-score equation in order to distinguish it from the decorrelated score equation in Ning and Liu (2017). Step 3 essentially is to recursively calculate $\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}$ and $\widehat{\boldsymbol{\omega}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}$ for Step 4. Thus, we refer this algorithm to as recursive online-score estimation (ROSE) algorithm.

Let

$$\Gamma_n = \frac{1}{n} \sum_{t \in [s_n, n-2]} \frac{X_{t0, 2, j_0}}{\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}^3} \left(X_{t0, 2, j_0} \widehat{\boldsymbol{\omega}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}^T \mathbf{X}_{t0, 2, \widehat{\mathcal{M}}_{j_0}^{(t)}} \right).$$

Under certain conditions, we can show that

$$\overline{n \ s_n} \Gamma_n(\bar{\beta}_{j_0} \ \beta_{1,j_0}) = \frac{1}{n \ s_n} \sum_{t \lfloor s_n}^{n-2} \frac{\varepsilon_{t0 \ 2}}{\sigma_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}} \left(X_{t0 \ 2, j_0} \ \boldsymbol{\omega}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}^T \mathbf{X}_{t0 \ 2, \widehat{\mathcal{M}}_{j_0}^{(t)}} \right) + o_p(1). \quad (5)$$

The first term on the right-hand-side (RHS) of (5) corresponds to a mean zero martingale with respect to the filtration $\{\sigma(\mathcal{M}_t) : t \sim s_n\}$ where $\sigma(\mathcal{M}_t)$ denotes the σ -algebra generated by \mathcal{M}_t . Note that

$$\mathbb{E} \left[\left\{ \frac{1}{\sigma_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}} \left(X_{t0 \ 2, j_0} \ \boldsymbol{\omega}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}^T \mathbf{X}_{t0 \ 2, \widehat{\mathcal{M}}_{j_0}^{(t)}} \right) \varepsilon_{t0 \ 2} \right\}^3 \middle| \mathcal{M}_t \right] = 1.$$

By the martingale central limit theorem, we have as $n \ s_n \infty$,

$$\overline{n \ s_n} \Gamma_n(\bar{\beta}_{j_0} \ \beta_{1,j_0}) \stackrel{d}{\infty} N(0, \sigma_1^3).$$

Therefore, a two-sided $1 - \alpha$ CI for β_{1,j_0} is given by

$$\bar{\beta}_{j_0} \subseteq z_{\frac{\alpha}{2}} \frac{\Gamma_n^{-2}}{n \ s_n} \hat{\sigma}, \quad (6)$$

where $\hat{\sigma}$ is some consistent estimator for σ_1 .

2.2 Refinements

The CI in (6) is asymptotically valid. However, it has one drawback. Its length is equal to

$$2z_{\frac{\alpha}{2}} \frac{\Gamma_n^{-2}}{n \ s_n} \hat{\sigma}. \quad (7)$$

In general, (7) increases as s_n increases. Nonetheless, s_n should be large enough to guarantee the sure screening property of $\widehat{\mathcal{M}}_{j_0}^{s_n+}$. For small n , this will result in a large CI. To address these concerns, we propose the following refined estimator. The estimation procedure is described below.

Step 1. Input $\{\mathbf{X}_i, Y_i\}_{i=1}^n$ and an integer $1 < s_n < n$.

Step 2. Compute an initial estimator $\tilde{\boldsymbol{\beta}}$ for $\boldsymbol{\beta}_1$, based on $\{\mathbf{X}_i, Y_i\}_{i=1}^n$.

Step 3. Compute $[\widehat{\mathcal{M}}_{j_0}^{(t)}]^{t+}$, $\widehat{\boldsymbol{\omega}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}$, $\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}^3$ for $t = s_n, \dots, n-1$ as described in Section 2.1.

Step 4. Estimate $[\widehat{\mathcal{M}}_{j_0}^{(-s_n)}]^{-s_n+}$ based on the sub-dataset $\{(\mathbf{X}_{s_n+2}, Y_{s_n+2}), \dots, (\mathbf{X}_n, Y_n)\}$. The resulting estimator $[\widehat{\mathcal{M}}_{j_0}^{(-s_n)}]^{-s_n+}$ shall satisfy $\|[\widehat{\mathcal{M}}_{j_0}^{(-s_n)}]^{-s_n+}\| < n$, $j_1 \notin [\widehat{\mathcal{M}}_{j_0}^{(-s_n)}]^{-s_n+}$.

Step 5. Estimate $\boldsymbol{\omega}_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)}, j_0}$ by $\widehat{\boldsymbol{\omega}}_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)}, j_0} = \widehat{\boldsymbol{\Sigma}}_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)}, \widehat{\mathcal{M}}_{j_0}^{(-s_n)}}^{-2} \widehat{\boldsymbol{\Sigma}}_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)}, j_0}$.

Step 6. Estimate $\sigma_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)}, j_0}^3$ by $\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)}, j_0}^3 = \widehat{\boldsymbol{\Sigma}}_{j_0, j_0}^{-2} \widehat{\boldsymbol{\Sigma}}_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)}, j_0}^T \widehat{\boldsymbol{\omega}}_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)}, j_0}$.

Step 7. Define $\hat{\beta}_{j_0}$ to be the solution to the following equation,

$$\begin{aligned} & \sum_{t=1}^{s_n-2} \frac{1}{\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)}, j_0}} \widehat{Z}_{t0, 2, j_0} (Y_{t0, 2} - X_{t0, 2, j_0} \beta_{1, j_0} - \mathbf{X}_{t0, 2, \widehat{\mathcal{M}}_{j_0}^{(-s_n)}} \tilde{\boldsymbol{\beta}}_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)}}) \\ & + \sum_{t=s_n}^{n-2} \frac{1}{\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}} \widehat{Z}_{t0, 2, j_0} (Y_{t0, 2} - X_{t0, 2, j_0} \beta_{1, j_0} - \mathbf{X}_{t0, 2, \widehat{\mathcal{M}}_{j_0}^{(t)}} \tilde{\boldsymbol{\beta}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}}) = 0, \end{aligned}$$

where $\widehat{Z}_{t0, 2, j_0} = X_{t0, 2, j_0} \widehat{\boldsymbol{\omega}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}^T \mathbf{X}_{t0, 2, \widehat{\mathcal{M}}_{j_0}^{(t)}}$ for $t = s_n, \dots, n-1$ and $\widehat{Z}_{t0, 2, j_0} = X_{t0, 2, j_0} \widehat{\boldsymbol{\omega}}_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)}, j_0}^T \mathbf{X}_{t0, 2, \widehat{\mathcal{M}}_{j_0}^{(-s_n)}}$ for $t = 0, \dots, s_n-1$.

When $s_n = o(n)$, the first s_n terms in the estimating equation in Step 7 is negligible. As a result, $\hat{\beta}_{j_0}$ is asymptotically the same as $\bar{\beta}_{j_0}$. Define

$$\Gamma_n^* = \frac{1}{n} \left(\sum_{t=1}^{s_n-2} \frac{1}{\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}} X_{t0, 2, j_0} \widehat{Z}_{t0, 2, j_0} + \sum_{t=s_n}^{n-2} \frac{1}{\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)}, j_0}} X_{t0, 2, j_0} \widehat{Z}_{t0, 2, j_0} \right). \quad (8)$$

Below, we prove

$$\hat{\beta}_{j_0} \subseteq z_{\frac{\alpha}{2}} \frac{\Gamma_n^{*-2}}{n} \widehat{\sigma}, \quad (9)$$

is a valid two-sided CI for β_{1, j_0} . We need the following conditions.

(A1) Assume $\widehat{\mathbb{I}}_{j_0}^{n+}$ satisfies $\Pr(\|\widehat{\mathbb{I}}_{j_0}^{n+}\| \geq \kappa_n) = 1$ for some $1 \geq \kappa_n = o(n)$. Besides,

$$\Pr\left(\mathbb{I}_{j_0} \leq \widehat{\mathbb{I}}_{j_0}^{n+}\right) \sim 1 - O\left(\frac{1}{n^{\alpha_0}}\right),$$

for some constant $\alpha_1 > 1$.

(A2) Assume there exists some constant $\bar{c} > 0$ such that for any $\mathbb{I} \leq \mathbb{I}$ and $\|\mathbb{I}\| \geq \kappa_n$, $\lambda_{n, \mathbb{I}}(\Sigma_{j_0 \cup \mathcal{M}, j_0 \cup \mathcal{M}}) \sim \bar{c}$.

(A3) Assume there exists some constant $c_1 > 0$ such that $\sqrt{\mathbf{X}_1^T \mathbf{a}}_{\psi_2} \geq c_1 \sqrt{\mathbf{a}}_{\psi_3}$ for any $\mathbf{a} \in \mathbb{R}^p$.

(A4) Assume (i) $\Pr(\sqrt{\widetilde{\boldsymbol{\beta}}_{\mathcal{M}_0^c}}_{\psi_3} \geq \eta_n) \rightarrow 1$ for some $\eta_n > 0$; (ii) $\eta_n \sqrt{\kappa_n \log p} = o(1)$; (iii) $\Pr(\sqrt{\widetilde{\boldsymbol{\beta}}_{\mathcal{M}_0^c}}_{\psi_3} \geq k_1 \sqrt{\widetilde{\boldsymbol{\beta}}_{\mathcal{M}_0}}_{\psi_3}) \rightarrow 1$ for some constant $k_1 > 0$, where $\mathbb{I}_{\mathcal{M}_0}$ stands for the support of $\boldsymbol{\beta}_1$ and $\mathbb{I}_{\mathcal{M}_0^c}$ denotes its complement.

(A5) Assume $\hat{\sigma} \stackrel{P}{\rightarrow} \sigma_1$.

Assumption (A1) essentially requires the sure screening property of the procedure for obtaining $\widehat{\mathbb{I}}_{j_0}^{n+}$. Typically conditions to guarantee the sure screening property are weaker than those for the oracle property. Assume (A1) holds and $s_n \rightarrow \infty$. Then it follows from Bonferroni's inequality that

$$\Pr\left(\mathbb{I}_{j_0} \leq \bigcap_{t \in [s_n]} \widehat{\mathbb{I}}_{j_0}^{t0.2+}\right) \sim 1 - O\left(\sum_{t \in [s_n]} \frac{1}{t^{\alpha_0}}\right) \rightarrow 1.$$

Hence, all the selected models possess the sure screening property with probability tending to 1.

When \mathbb{I}_{j_0} is estimated via SIS, we can show (A1) holds for any arbitrary $\alpha_1 > 1$ (see Theorem 1 in Fan and Lv, 2008). The validity of such sure screening property typically relies on certain minimum-signal-strength conditions on $\boldsymbol{\beta}_{1, \mathbb{I}_{j_0}}$. A counterexample is given in Section B.1.1 of the supplementary article where we show our CI is no longer valid when these conditions are violated. We note that van de Geer et al. (2014) and Ning and Liu (2017) do not require these conditions. However, these authors impose some additional assumptions on the design matrix. We discuss this further in Section B.1 of the supplementary article. Moreover, in Section 5.4, we present a variant of our method that

is valid without the minimal-signal-strength conditions.

Condition (A2) is similar to the restricted eigenvalue condition (Bickel et al., 2009) imposed to derive the oracle inequalities for the Lasso estimator and the Dantzig selector. Condition (A3) requires \mathbf{X}_1 to be a sub-Gaussian vector. This condition is used in Ning and Liu (2017) and van de Geer et al. (2014) as well. See Section 4.1 of Ning and Liu (2017) and Condition (B1) in van de Geer et al. (2014) for details.

When $\tilde{\boldsymbol{\beta}}$ is estimated via the Lasso or the Dantzig selector, then the first part of Condition (A4) holds with $\eta_n = c_n \sqrt{s^* \log p/n}$ where c_n is an arbitrary diverging sequence and s^* is the number of nonzero elements in $\boldsymbol{\beta}_1$. The second part holds as long as $\kappa_n s^* \log^3 p = o(n)$. Under (A1), we have $\kappa_n \sim s^* - 1$. This further implies $(s^*)^3 \log^3 p = o(n)$. Such a sample size requirement is consistent with those in van de Geer et al. (2014) and Ning and Liu (2017). See Condition (B2) of van de Geer et al. (2014), and Corollary 4.1 in Ning and Liu (2017) for details. The last condition in (A4) holds with $k_1 = 3$ for the Lasso estimator, $k_1 = 1$ for the Dantzig selector and $k_1 = 0$ for the non-convex penalized regression estimator (when the ‘‘oracle property’’ is achieved).

Condition (A5) holds when $\hat{\sigma}$ is computed by refitted cross-validation (Fan et al., 2012) or scaled lasso (Sun and Zhang, 2013). In Section B.3 of the supplementary article, we further introduce a simple plug-in estimator for σ_1^3 based on $\tilde{\boldsymbol{\beta}}$ and show (A5) holds under (A3), (A4) and the conditions that $\log p = O(n^{3/4})$, $E\|\varepsilon_1\|^4 = O(1)$. The last moment condition is also needed in the following theorem to guarantee the asymptotic normality of $\hat{\beta}_{j_0}$.

Theorem 2.1. *Under Conditions (A1)-(A5), assume $s_n \rightarrow \infty$, $s_n = o(n)$, $\kappa_n^3 \log p = O(n/\log^3 n)$ and $E\|\varepsilon_1\|^4 = O(1)$. Then, we have*

$$\frac{\bar{n}\Gamma_n^*(\hat{\beta}_{j_0} \quad \beta_{1,j_0})}{\hat{\sigma}} \stackrel{d}{\rightarrow} N(0, 1),$$

where Γ_n^* is defined in (8).

3 High-dimensional generalized linear models

3.1 Estimation and inference

Suppose that $(\mathbf{X}_2, Y_2), \dots, (\mathbf{X}_n, Y_n)$ is a random sample from (\mathbf{X}_1, Y_1) in (1). The function $b(\cdot)$ is assumed to be thrice continuously differentiable. We further assume $b''(\cdot) > 0$ and $b'''(\cdot)$ is Lipschitz continuous. Denoted by $\mu(\cdot)$ the derivative of $b(\cdot)$. As in Section 2, our focus is to construct a CI for β_{1,j_0} . Let $\lfloor j_0 = \lfloor j \neq j_1 : \beta_{1,j} \neq 0 \rfloor$ and $\Sigma = E\mathbf{X}_1 b''(\mathbf{X}_1^T \beta_1) \mathbf{X}_1^T$, we describe our estimating procedure below.

Step 1. Input $\{\mathbf{X}_i, Y_i\}_{i=2}^n$ and an integer $1 < s_n < n$.

Step 2. Compute an initial estimator $\tilde{\beta}$ for β_1 . Compute

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=2}^n \mathbf{X}_i b''(\mathbf{X}_i^T \tilde{\beta}) \mathbf{X}_i^T. \quad (10)$$

Step 3. For $t = s_n, s_n + 1, \dots, n - 1$, estimate $\lfloor j_0$ based on the sub-dataset $\mathcal{M} = \{(\mathbf{X}_2, Y_2), \dots, (\mathbf{X}_t, Y_t)\}$. Denoted by $\widehat{\lfloor j_0}^{t+}$ the corresponding estimator. We require $\|\widehat{\lfloor j_0}^{t+}\| \geq n, j_1 \notin \widehat{\lfloor j_0}^{t+}$. Compute

$$\hat{\omega}_{\widehat{\lfloor j_0}^{(t)}, j_0} = \hat{\Sigma}_{\widehat{\lfloor j_0}^{(t)}, \widehat{\lfloor j_0}^{(t)}}^{-2} \hat{\Sigma}_{\widehat{\lfloor j_0}^{(t)}, j_0} \quad \text{and} \quad \hat{\sigma}_{\widehat{\lfloor j_0}^{(t)}, j_0}^3 = \hat{\Sigma}_{j_0, j_0} \quad \hat{\Sigma}_{\widehat{\lfloor j_0}^{(t)}, j_0}^T \hat{\omega}_{\widehat{\lfloor j_0}^{(t)}, j_0}.$$

Step 4. Estimate $\lfloor j_0$ based on the sub-dataset $\{(\mathbf{X}_{s_n+2}, Y_{s_n+2}), \dots, (\mathbf{X}_n, Y_n)\}$. Denoted by $\widehat{\lfloor j_0}^{-s_n+}$ the resulting estimator. We require $\|\widehat{\lfloor j_0}^{-s_n+}\| \geq n, j_1 \notin \widehat{\lfloor j_0}^{-s_n+}$. Compute

$$\hat{\omega}_{\widehat{\lfloor j_0}^{(-s_n)}, j_0} = \hat{\Sigma}_{\widehat{\lfloor j_0}^{(-s_n)}, \widehat{\lfloor j_0}^{(-s_n)}}^{-2} \hat{\Sigma}_{\widehat{\lfloor j_0}^{(-s_n)}, j_0} \quad \text{and} \quad \hat{\sigma}_{\widehat{\lfloor j_0}^{(-s_n)}, j_0}^3 = \hat{\Sigma}_{j_0, j_0} \quad \hat{\Sigma}_{\widehat{\lfloor j_0}^{(-s_n)}, j_0}^T \hat{\omega}_{\widehat{\lfloor j_0}^{(-s_n)}, j_0}.$$

Step 5. Define $\hat{\beta}_{j_0}$ to be the solution to the following equation,

$$\begin{aligned} & \sum_{t=1}^{s_n-2} \frac{\widehat{Z}_{t0\ 2,j_0}}{\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)},j_0}} \left\{ Y_{t0\ 2} - \mu \left(X_{t0\ 2,j_0} \beta_{1,j_0} + \mathbf{X}_{t0\ 2,\widehat{\mathcal{M}}_{j_0}^{(-s_n)}} \widetilde{\boldsymbol{\beta}}_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)}} \right) \right\} \\ & + \sum_{t=1}^{n-2} \frac{\widehat{Z}_{t0\ 2,j_0}}{\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(t)},j_0}} \left\{ Y_{t0\ 2} - \mu \left(X_{t0\ 2,j_0} \beta_{1,j_0} + \mathbf{X}_{t0\ 2,\widehat{\mathcal{M}}_{j_0}^{(t)}} \widetilde{\boldsymbol{\beta}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}} \right) \right\} = 0, \end{aligned}$$

$$\text{where } \widehat{Z}_{t0\ 2,j_0} = X_{t0\ 2,j_0} \widehat{\boldsymbol{\omega}}_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)},j_0}^T \mathbf{X}_{t0\ 2,\widehat{\mathcal{M}}_{j_0}^{(-s_n)}} \text{ for } t = 0, \dots, s_n - 1.$$

The estimating function in Step 5 can be solved via the Newton-Raphson method with the initial value $\hat{\beta}_{j_0}^{1+} = \widetilde{\beta}_{j_0}$. More specifically, for $l = 1, 2, \dots$, we can iteratively update $\hat{\beta}_{j_0}$ by

$$\hat{\beta}_{j_0}^{l+} = \hat{\beta}_{j_0}^{l-2+} + \frac{\sum_{t=1}^{n-2} \frac{1}{\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(t)},j_0}} \widehat{Z}_{t0\ 2,j_0} \left\{ Y_{t0\ 2} - \mu \left(X_{t0\ 2,j_0} \hat{\beta}_{j_0}^{l-2+} + \mathbf{X}_{t0\ 2,\widehat{\mathcal{M}}_{j_0}^{(t)}} \widetilde{\boldsymbol{\beta}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}} \right) \right\}}{\sum_{t=1}^{n-2} \frac{1}{\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(t)},j_0}} \widehat{Z}_{t0\ 2,j_0} X_{t0\ 2,j_0} b'' \left(X_{t0\ 2,j_0} \hat{\beta}_{j_0}^{l-2+} + \mathbf{X}_{t0\ 2,\widehat{\mathcal{M}}_{j_0}^{(t)}} \widetilde{\boldsymbol{\beta}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}} \right)}, \quad (11)$$

where we use a shorthand and write $\widehat{\Gamma}_{j_0}^{t+} = \widehat{\Gamma}_{j_0}^{-s_n+}$, for $t = 0, \dots, s_n - 1$. Define

$$\Gamma_n^{*,l-2+} = \frac{1}{n} \sum_{t=1}^{n-2} \frac{1}{\widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(t)},j_0}} \widehat{Z}_{t0\ 2,j_0} X_{t0\ 2,j_0} b'' \left(X_{t0\ 2,j_0} \hat{\beta}_{j_0}^{l-2+} + \mathbf{X}_{t0\ 2,\widehat{\mathcal{M}}_{j_0}^{(t)}} \widetilde{\boldsymbol{\beta}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}} \right).$$

A two-sided $1 - \alpha$ CI for β_{1,j_0} is given by

$$\hat{\beta}_{j_0}^{l+} \subseteq \frac{z_{\frac{\alpha}{2}} \hat{\phi}^{2/3}}{n \Gamma_n^{*,l-2+}}, \quad (12)$$

where $\hat{\phi}$ is some consistent estimator for ϕ_1 . We state the following conditions.

(A1*) Assume $\widehat{\Gamma}_{j_0}^{n+}$ satisfies $\Pr(\|\widehat{\Gamma}_{j_0}^{n+}\| \geq \kappa_n) = 1$ for some $1 \geq \kappa_n = o(n)$. Besides, there exists some constant $\alpha_1 > 1$ such that

$$\Pr \left(\widehat{\Gamma}_{j_0} \leq \widehat{\Gamma}_{j_0}^{n+} \right) \sim 1 - O \left(\frac{1}{n^{\alpha_0}} \right),$$

(A2*) Assume there exists some constant $\bar{c} > 0$ such that for any $l \leq \mathbb{I}$ and $\|\beta\| \geq \kappa_n$, $\lambda_{n,m}(\Sigma_{j_0 \cup \mathcal{M}, j_0 \cup \mathcal{M}}) \sim \bar{c}$.

(A3*) Assume there exists some constant $c_1 > 0$ such that $\sqrt{\mathbf{X}_1^T \mathbf{a}} \sqrt{\psi_2} \geq c_1 \sqrt{\mathbf{a}} \sqrt{\psi_3}$ for any $\mathbf{a} \in \mathbb{R}^p$.

(A4*) Assume $\max_{j \in \{2, \dots, p\}} \|X_{1,j}\| \geq \omega_1$ for some constant $\omega_1 > 0$. Assume $\|\mathbf{X}_1^T \beta_1\| \geq \bar{\omega}$ for some constant $\bar{\omega} > 0$.

(A5*) Assume (i) $\Pr(\sqrt{\tilde{\beta}} \|\beta_1\| \sqrt{\psi_3} \geq \eta_n) \rightarrow 1$ for some $\eta_n > 0$; (ii) $\eta_n \overline{\kappa_n \log p} = o(1)$ and $\bar{n} \eta_n^3 = o(1)$; (iii) $\Pr(\sqrt{\tilde{\beta}_{\mathcal{M}_0^c}} \|\beta_{1, \mathcal{M}_0^c}\| \sqrt{\psi_2} \geq k_1 \sqrt{\tilde{\beta}_{\mathcal{M}_0}} \|\beta_{1, \mathcal{M}_0}\| \sqrt{\psi_2}) \rightarrow 1$ for some constant $k_1 > 0$.

(A6*) Assume $\sqrt{Y_1} \mu(\mathbf{X}_1^T \beta_1) \sqrt{\psi_1 | \mathbf{X}_0}$ is uniformly bounded for all \mathbf{X}_1 , where $\sqrt{\psi_1 | \mathbf{X}_0}$ denotes the Orlicz norm conditional on \mathbf{X}_1 .

(A7*) Assume $\hat{\phi} \stackrel{P}{\infty} \phi_1$.

Conditions (A1*)-(A3*) are very similar to (A1)-(A3). In (A4*), for technical convenience, we assume $X_{1,j}$'s and $\mathbf{X}_1^T \beta_1$ are bounded. In (A5*), we further assume $\bar{n} \eta_n^3 = o(1)$. Note that such assumption doesn't appear in (A4). This is because we focus on a more general class of models here. Assume (A4*) holds. Then (A6*) is automatically satisfied for logistic and Poisson regression models. In logistic or Poisson regression models, we have $\phi_1 = 1$. Condition (A7*) thus automatically holds by setting $\hat{\phi} = 1$.

Theorem 3.1. *Assume (A1*)-(A7*) hold. Assume $s_n \in \mathbb{N}$, $s_n = o(n)$, $\kappa_n^{6/3} \log p = O(n/\log^3 n)$ and $\kappa_n^4 = O(n)$. Then, for any fixed $l \sim 1$, we have*

$$\frac{\bar{n} \Gamma_n^{*)^{l-2+}}}{\hat{\phi}^{2/3}} (\hat{\beta}_{j_0}^{l+} - \beta_{1, j_0}) \stackrel{d}{\infty} N(0, 1).$$

Theorem 3.1 proves the validity of the two-sided CI in (12), for any $l \sim 1$. When $l = 1$, $\hat{\beta}_{j_0}^{l+}$ corresponds to the solution of the first-order approximation of the score equation. We note that Bickel (1975) and Ning and Liu (2017) used a similar one-step approximation to ensure the consistency of the resulting estimator. In practice, we can update $\hat{\beta}_{j_0}^{l+}$ for a few Newton steps. In our numerical experiments, we find that $\hat{\beta}_{j_0}^{l+}$ converges very fast and it suffices to set $l = 3, 4$ or 5 .

3.2 Asymptotic efficiency

Theorem 3.1 proves the validity of the CI in (12). The length of the CI is given by

$$L(\hat{\beta}_{j_0}^{l+}, \alpha) = 2z_{\alpha/3} \frac{\hat{\phi}^{2/3}}{\Gamma_n^{*,l-2+} \bar{n}}. \quad (13)$$

Under the given conditions in Theorem 3.1, it follows from the law of large numbers for martingales (Csörgö, 1968) that

$$\Gamma_n^{*,l-2+} = \frac{s_n}{n} \sigma_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)}, j_0} + \frac{1}{n} \left(\sum_{t \lfloor s_n}^{n-2} \sigma_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0} \right) + o_p(1), \quad (14)$$

where

$$\sigma_{\mathcal{M}, j_0}^3 = \Sigma_{j_0, j_0} \Sigma_{\mathcal{M}, j_0}^T \Sigma_{\mathcal{M}, \mathcal{M}}^{-2} \Sigma_{\mathcal{M}, j_0},$$

for any $\lfloor \cdot \rfloor \leq \mathbb{I}_{j_0}$.

By Assumption (A1*) and (A2*), we have almost surely,

$$\sigma_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)}, j_0}^3 \sim \bar{c} \quad \text{and} \quad \sigma_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}^3 \sim \bar{c}, \quad \mathcal{H} = s_n, \dots, n-1. \quad (15)$$

Under (A7*), $\hat{\phi}$ is consistent. This together with (13)-(15) yields

$$\bar{n}L(\hat{\beta}_{j_0}, \alpha) = \frac{2z_{\alpha/3} \phi_1^{2/3}}{s_n \sigma_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)}, j_0} / n + \sum_{t \lfloor s_n}^{n-2} \sigma_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0} / n} + o_p(1). \quad (16)$$

Based on (16), we compare the length of the CI of the proposed method with the desparsified Lasso method, the decorrelated score method and the ‘oracle’ method below.

3.2.1 Comparison with the de-sparsified Lasso and the decorrelated score

Consider the Lasso estimator

$$\hat{\beta}^L = \arg \min_{\beta} \left(\frac{1}{n} \sum_{i \in \mathcal{I}_2} b(\mathbf{X}_i^T \beta) - Y_i \mathbf{X}_i^T \beta + \lambda_n \|\beta\|_2 \right).$$

The de-sparsified Lasso estimator is defined by

$$\hat{\beta}^{DL} = \hat{\beta}^L + \hat{\Theta} \left\{ \frac{1}{n} \sum_{i \in \mathcal{I}_2} \left(\mathbf{X}_i^T Y_i - \mu(\mathbf{X}_i^T \hat{\beta}^L) \right) \right\},$$

where the matrix $\hat{\Theta}$ is computed by the nodewise Lasso (see Section 3.1.1 in van de Geer et al., 2014). Theorem 3.1 in van de Geer et al. (2014) proved that

$$\bar{n}(\hat{\beta}_{j_0}^{DL} - \beta_{1,j_0}) / \sqrt{e_{j_0,p}^T \hat{\Omega} e_{j_0,p}} \rightarrow N(0, 1), \quad (17)$$

where

$$\hat{\Omega} = \hat{\Theta} \left(\frac{1}{n} \sum_{i \in \mathcal{I}_2} \mathbf{X}_i \mathbf{X}_i^T Y_i - \mu(\mathbf{X}_i^T \hat{\beta}^L) \mathbf{X}_i^T \right) \hat{\Theta}^T,$$

and

$$e_{j_1, j_2} = \underbrace{(0, \dots, 0)}_{j_1-2}, 1, \underbrace{(0, \dots, 0)}_{j_2-j_1}.$$

for any integer $1 \geq j_2 < j_3$.

Based on the de-sparsified Lasso estimator, the corresponding CI for β_{1,j_0} is given by

$$\left[\hat{\beta}_{j_0}^{DL} - z_{\frac{\alpha}{2}} \frac{\sqrt{e_{j_0,p}^T \hat{\Omega} e_{j_0,p}}}{\bar{n}}, \hat{\beta}_{j_0}^{DL} + z_{\frac{\alpha}{2}} \frac{\sqrt{e_{j_0,p}^T \hat{\Omega} e_{j_0,p}}}{\bar{n}} \right], \quad (18)$$

Moreover, it follows from Theorem 3.2 in van de Geer et al. (2014) that

$$\mathbf{e}_{j_0,p}^T \widehat{\boldsymbol{\Omega}} \mathbf{e}_{j_0,p} = \mathbf{e}_{j_0,p}^T \boldsymbol{\Sigma}^{-2} \mathbf{e}_{j_0,p} \phi_1 + o_p(1).$$

Therefore, the length of (18) satisfies

$$\bar{n}L(\hat{\beta}_{j_0}^{DL}, \alpha) = 2z_{\alpha/3} \phi_1^{2/3} \sqrt{\mathbf{e}_{j_0,p}^T \boldsymbol{\Sigma}^{-2} \mathbf{e}_{j_0,p}} + o_p(1). \quad (19)$$

Ning and Liu (2017) proposed to construct the CI for high-dimensional parameters in GLM based on the one-step estimator that solves a first-order approximation of the decorrelated score equation. Specifically, the one-step estimator is given by

$$\hat{\beta}_{j_0}^{DS} = \tilde{\beta}_{j_0} + \frac{\sum_{t=1}^{n-2} (X_{t0,2,j_0} \quad \widehat{\mathbf{w}}^T \mathbf{X}_{t0,2,\mathbb{I}_{j_0}}) \left\{ Y_{t0,2} - \mu \left(\mathbf{X}_{t0,2}^T \tilde{\boldsymbol{\beta}} \right) \right\}}{\sum_{t=1}^{n-2} (X_{t0,2,j_0} \quad \widehat{\mathbf{w}}^T \mathbf{X}_{t0,2,\mathbb{I}_{j_0}}) X_{t0,2,j_0} b'' \left(\mathbf{X}_{t0,2}^T \tilde{\boldsymbol{\beta}} \right)}, \quad (20)$$

where $\tilde{\boldsymbol{\beta}}$ and $\widehat{\mathbf{w}}$ are some consistent estimators for $\boldsymbol{\beta}_1$ and $\boldsymbol{\Sigma}_{\mathbb{I}_{j_0}, \mathbb{I}_{j_0}}^{-2} \boldsymbol{\Sigma}_{\mathbb{I}_{j_0}, j_0}$, respectively. The corresponding CI is given by

$$\left[\hat{\beta}_{j_0}^{DS} - z_{\frac{\alpha}{2}} \hat{\phi}^{2/3} (\hat{\sigma}_s)^{-2/3}, \hat{\beta}_{j_0}^{DS} + z_{\frac{\alpha}{2}} \hat{\phi}^{2/3} (\hat{\sigma}_s)^{-2/3} \right],$$

where $\hat{\sigma}_s$ is the denominator of the second term on the RHS of (20) and $\hat{\phi}$ is some consistent estimator for ϕ_1 . Under certain conditions, we can show

$$\bar{n}L(\hat{\beta}_{j_0}^{DS}, \alpha) = \frac{2z_{\frac{\alpha}{2}} \phi_1^{2/3}}{\sigma_{\mathbb{I}_{j_0}, j_0}} + o_p(1) = 2z_{\frac{\alpha}{2}} \phi_1^{2/3} \sqrt{\mathbf{e}_{j_0,p}^T \boldsymbol{\Sigma}^{-2} \mathbf{e}_{j_0,p}} + o_p(1), \quad (21)$$

where the last equality follows from the matrix inversion formula (see Lemma A.5 in the supplementary material). This together with (19) implies that the lengths of CIs based on $\hat{\beta}_{j_0}^{DL}$ and $\hat{\beta}_{j_0}^{DS}$ are asymptotically the equivalent.

For any $l \leq \mathbb{I}_{j_0}$, let

$$\begin{aligned}\boldsymbol{\xi}_{\mathcal{M},j_0} &= \mathbb{E}(X_{1,j_0} \boldsymbol{\omega}_{\mathcal{M},j_0}^T \mathbf{X}_{1,\mathcal{M}}) b''(\mathbf{X}_i^T \boldsymbol{\beta}_1) (\mathbf{X}_{1,(\mathcal{M} \cup \{j_0\})^c} \boldsymbol{\Sigma}_{(\mathcal{M} \cup \{j_0\})^c, \mathcal{M}} \boldsymbol{\Sigma}_{\mathcal{M}, \mathcal{M}}^{-2} \mathbf{X}_{1,\mathcal{M}}) \\ &= \boldsymbol{\Sigma}_{(\mathcal{M} \cup \{j_0\})^c, j_0} \boldsymbol{\Sigma}_{(\mathcal{M} \cup \{j_0\})^c, \mathcal{M}} \boldsymbol{\omega}_{\mathcal{M},j_0}.\end{aligned}$$

We have the following results.

Theorem 3.2. *Assume (16), (19), (21), (A3*) and (A4*) hold. Let $\bar{k} = \sup_{|z| \leq \psi} b''(z)$. Then for any $0 < \alpha < 1$, $l \sim 1$,*

$$\begin{aligned}\bar{n}L(\hat{\beta}_{j_0}^{DL}, \alpha) &= \bar{n}L(\hat{\beta}_{j_0}^{DS}, \alpha) + o_p(1) \\ &\sim \bar{n}L(\hat{\beta}_{j_0}^{l+}, \alpha) + \frac{\phi_1^{2/3} z_{\alpha/3}}{k^{4/3} c_1^6} \left(\frac{s_n}{n} \sqrt{\boldsymbol{\xi}_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)}, j_0}} \sqrt{\frac{3}{3}} + \frac{1}{n} \sum_{t \leq s_n}^{n-2} \sqrt{\boldsymbol{\xi}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}} \sqrt{\frac{3}{3}} \right) + o_p(1).\end{aligned}$$

Theorem 3.2 implies that the proposed CI is asymptotically shorter than those based on the de-sparsified Lasso and the decorrelated score statistic. The difference depends on the L_3 norm of $\boldsymbol{\xi}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}$, which measures the partial dependence between X_{1,j_0} and $\mathbf{X}_{1,(\widehat{\mathcal{M}}_{j_0}^{(t)} \cup j_0)^c}$, after adjusted by $\mathbf{X}_{1, \widehat{\mathcal{M}}_{j_0}^{(t)}}$. For linear regression models, we have $\boldsymbol{\xi}_{\mathcal{M},j_0} = 0$ when X_{1,j_0} is independent of other predictors. However, $\sqrt{\boldsymbol{\xi}_{\mathcal{M},j_0}} \sqrt{\frac{3}{3}}$ can be positive when X_{1,j_0} is partially correlated with $\mathbf{X}_{1,(\mathcal{M} \cup j_0)^c}$ given $\mathbf{X}_{1,\mathcal{M}}$.

Although our method yields narrower CI on average, its validity relies on certain minimal-signal-strength conditions on $\boldsymbol{\beta}_{1, \mathbb{I}_{j_0}}$, as discussed in Section 2.2. This is a potential disadvantage of our method. Moreover, our procedure can be more time consuming than the existing methods, as it requires to recursively estimate the support set based on different data subsets. A variant of our method is proposed in Section 3.3 to reduce the computational cost.

3.2.2 Comparison with the oracle method

We compare the proposed CI with the CI of the oracle method. The oracle knew the set \mathcal{M}_{j_0} ahead of time. It estimates β_{1,j_0} by $\hat{\beta}_{j_0}^{oracle}$ defined as

$$(\hat{\beta}_{j_0}^{oracle}, \hat{\beta}_{\mathcal{M}_{j_0}}^{oracle}) = \arg \min_{\beta_{j_0}, \beta_{\mathcal{M}_{j_0}}} \frac{1}{n} \sum_{i \in \mathcal{I}_2} \left(b(X_{i,j_0} \beta_{j_0} + \mathbf{X}_{i,\mathcal{M}_{j_0}}^T \beta_{\mathcal{M}_{j_0}}) - Y_i(X_{i,j_0} \beta_{j_0} + \mathbf{X}_{i,\mathcal{M}_{j_0}}^T \beta_{\mathcal{M}_{j_0}}) \right).$$

Let

$$\hat{\Sigma}^{oracle} = \frac{1}{n} \sum_{i \in \mathcal{I}_2} \mathbf{X}_i b''(X_{i,j_0} \hat{\beta}_{j_0}^{oracle} + \mathbf{X}_{i,\mathcal{M}_{j_0}}^T \hat{\beta}_{\mathcal{M}_{j_0}}^{oracle}) \mathbf{X}_i.$$

The asymptotic variance of $\bar{n} \hat{\beta}_{j_0}^{oracle}$ can be consistently estimated by

$$\hat{\phi} \mathbf{e}_{2,|\mathcal{M}_{j_0}|}^T \left(\begin{array}{cc} \hat{\Sigma}_{j_0,j_0}^{oracle} & \hat{\Sigma}_{j_0,\mathcal{M}_{j_0}}^{oracle} \\ \hat{\Sigma}_{\mathcal{M}_{j_0},j_0}^{oracle} & \hat{\Sigma}_{\mathcal{M}_{j_0},\mathcal{M}_{j_0}}^{oracle} \end{array} \right)^{-2} \mathbf{e}_{2,|\mathcal{M}_{j_0}|} = \hat{\phi} \left\{ \hat{\Sigma}_{j_0,j_0}^{oracle} - \hat{\Sigma}_{j_0,\mathcal{M}_{j_0}}^{oracle} \left(\hat{\Sigma}_{\mathcal{M}_{j_0},\mathcal{M}_{j_0}}^{oracle} \right)^{-2} \hat{\Sigma}_{\mathcal{M}_{j_0},j_0}^{oracle} \right\}^{-2},$$

where the equality follows by the matrix inversion formula (see Lemma A.5). Let

$$\hat{\sigma}_{\mathcal{M}_{j_0},j_0}^{oracle} = \sqrt{\hat{\Sigma}_{j_0,j_0}^{oracle} - \hat{\Sigma}_{j_0,\mathcal{M}_{j_0}}^{oracle} \left(\hat{\Sigma}_{\mathcal{M}_{j_0},\mathcal{M}_{j_0}}^{oracle} \right)^{-2} \hat{\Sigma}_{\mathcal{M}_{j_0},j_0}^{oracle}}.$$

The corresponding confidence interval is given by

$$\left[\hat{\beta}_{j_0}^{oracle} - z_{\frac{\alpha}{2}} \hat{\phi}^{2/3} / (\bar{n} \hat{\sigma}_{\mathcal{M}_{j_0},j_0}^{oracle}), \hat{\beta}_{j_0}^{oracle} + z_{\frac{\alpha}{2}} \hat{\phi}^{2/3} / (\bar{n} \hat{\sigma}_{\mathcal{M}_{j_0},j_0}^{oracle}) \right], \quad (22)$$

where $\hat{\phi}$ is some constant estimator for ϕ_1 .

Under certain conditions, the length of (22) satisfies

$$\bar{n} \mathbb{L}(\hat{\beta}_{j_0}^{oracle}, \alpha) = \frac{2z_{\alpha/3} \phi_1^{2/3}}{\sigma_{\mathcal{M}_{j_0},j_0}} + o_p(1). \quad (23)$$

Theorem 3.3. Assume (16) and (23) hold. Assume $s_n \in \mathcal{E}$, $n \in \mathcal{E}$, and there

exists some $\alpha_1 > 1$ such that

$$Pr\left(\lfloor j_0 = \widehat{j_0}^{n+} \right) \sim 1 - O\left(\frac{1}{n^{\alpha_0}}\right). \quad (24)$$

Then for any $0 < \alpha < 1$, $l \sim 1$, we have

$$\bar{n}L(\hat{\beta}_{j_0}^{l+}, \alpha) = \bar{n}L(\hat{\beta}_{j_0}^{oracle}, \alpha) + o_p(1).$$

Condition (24) in Theorem 3.3 requires the variable selection procedure to be consistent. Under this condition, we prove the ‘‘oracle’’ property of our method, which means that the length of the proposed CI is asymptotically equivalent to the CI of the oracle method.

3.3 Computationally efficient procedure

The proposed estimation procedure in Section 3.1 requires to estimate $\lfloor j_0$ ($n - s_n + 1$) times. This will be time consuming for large n . To address this concern, for a given integer $S > 1$, we can compute $\widehat{j_0}^{t+}$ approximately $(n - s_n)/S$ times based on the sub-dataset \mathcal{M} for $t = s_n, s_n + S, s_n + 2S, \dots, s_n + \lfloor (n - s_n)/S \rfloor S$ where $\lfloor z \rfloor$ denotes the largest integer smaller than or equal to z . For any $s_n < t < n$, define

$$\widehat{j_0}^{t+} = \widehat{j_0}^{s_n + l_1 S},$$

for some nonnegative integer l_1 such that $s_n + l_1 S \geq t < s_n + (l_1 + 1)S$. The resulting estimator $\hat{\beta}_{j_0}^{l+}$ is computed by (11). The corresponding CI can be similarly derived as in (12).

4 Numerical examples

4.1 Linear regression

In this section, we conduct some simulation studies to examine the performance of the proposed CI in high dimensional linear regression models. Suppose that $\{\mathbf{X}_i, Y_i\}$, $i =$

$1, \dots, n$ is a sample from the following model:

$$Y_i = \mathbf{X}_i^T \boldsymbol{\beta}_1 + \varepsilon_i, \quad (25)$$

where $\varepsilon_i \rightarrow N(0, 1)$, $\mathbf{X}_i \rightarrow N(0, \boldsymbol{\Sigma})$.

Consider the following four settings: (A) $n = 100$, $\beta_{1,2} = \beta_{1,3} = 1.0$ and $\beta_{1,j} = 0$ for $j > 2$; (B) $n = 100$, $\beta_{1,2} = \beta_{1,3} = 2.0$ and $\beta_{1,j} = 0$ for $j > 2$; (C) $n = 200$, $\beta_{1,2} = 2.0$, $\beta_{1,3} = 2.0$ and $\beta_{1,j} = 0$ for $j > 2$; (D) $n = 200$, $\beta_{1,2} = \beta_{1,3} = \beta_{1,4} = \beta_{1,5} = \beta_{1,6} = 1.0$ and $\beta_{1,j} > 0$ for $j > 5$. For each setting, we set $p = 1000$, and consider two different covariance matrices $\boldsymbol{\Sigma}$, corresponding to $\boldsymbol{\Sigma} = \mathbf{I}$ and $\boldsymbol{\Sigma} = \{0.5^{|i-j|}\}_{i,j \in \{2, \dots, p\}}$. This yields a total of 8 scenarios. For the first three settings, the objective is to construct 95% two-sided CIs for $\beta_{1,3}$ and $\beta_{1,4}$. For the last setting, we aim to construct 95% two-sided CIs for $\beta_{1,4}$, $\beta_{1,5}$, $\beta_{1,6}$ and $\beta_{1,=}$. Comparison is made among the following CIs:

- (i) The proposed CI in (9), labeled by ROSE in Tables 1 and 2;
- (ii) The CI constructed by the de-sparsified Lasso (DLASSO) method;
- (iii) The CI constructed by the Bootstrap Lasso+Partial Ridge (BLPR) method (Liu et al., 2017);
- (iv) The CI constructed by the simple sample-splitting (S3) method.

To calculate the CI in (9), we set $s_n = \lfloor 2n/\log(n) \rfloor$. Such a choice of s_n satisfies the conditions in Theorem 2.1. The set $\lfloor \cdot \rfloor_{j_0}$ is estimated by ISIS. The estimation procedure is implemented by the R package `SIS` (Saldana and Feng, 2016). To compute the initial estimator $\tilde{\boldsymbol{\beta}}$, we first apply ISIS based on all observations and then fit a penalized linear regression model using the R package `ncvreg` (Breheny and Huang, 2011) with SCAD penalty function for the variables picked by ISIS. The variance estimator $\hat{\sigma}$ is computed by refitted cross-validation. We implement the CI in (ii) by the R package `hdi` (Dezeure et al., 2015). BLPR estimates β_{1,j_0} by the Lasso+Partial Ridge (LPR) estimator. More specifically, it first uses the Lasso to select important predictors and then refit the model using partial ridge regression based on the selected variables. The corresponding CI for β_{1,j_0} is constructed by bootstrapping the LPR estimator. We implement the BLPR method by the R package `HDICI`. To compute the CI in (iv), we randomly split the samples into two

equal halves, use ISIS to estimate the support of control variables and construct the CI based on the remaining second half of the data. In Table 1 and 2, we report the empirical coverage probability (ECP) and average length (AL) of these CIs. Results are averaged over 500 simulations.

It can be seen from Table 1 that ECPs of our procedure and the S3 method are close to the nominal level in all cases. However, CIs constructed by the S3 method are approximately $\sqrt{2}$ time wide than our proposed method, according to Table 2. As commented in the introduction, this is because S3 only uses half of the samples to evaluate β_{1,j_0} .

Under the settings where $\Sigma = \{0.5^{|i-j|}\}_{i,j}$, ECPs of the DLASSO method are significantly smaller than the nominal level. For example, in Setting (A) and (B), ECPs of the DLASSO method are smaller than 90% when $\Sigma = \{0.5^{|i-j|}\}_{i,j}$. Under the settings where $\Sigma = \mathbf{I}_p$, CIs constructed by the DLASSO method have approximately nominal coverage probabilities. However, we note these CIs are wider than the proposed CIs in all cases. Take Setting (D) as an example. When $\Sigma = \mathbf{I}_p$, ALs of the DLASSO method are approximately 10% larger than the proposed method.

We note that BLPR yields very narrow CIs for zero parameters. For nonzero parameters however, the CIs based on the BLPR method are much wider than the proposed CIs in all cases. Moreover, under the settings where $\Sigma = \mathbf{I}_p$, ECPs of the BLPR method are significantly smaller than the nominal level for nearly all nonzero parameters.

4.2 Logistic regression

We generate $\{\mathbf{X}_i, Y_i\}_{i \in \{2, \dots, n\}}$ from the following logistic regression model

$$\text{logit}\{\Pr(Y_i = 1|\mathbf{X}_i)\} = \mathbf{X}_i^T \boldsymbol{\beta}_1,$$

where $\text{logit}(z) = \log\{z/(1-z)\}$ for $0 < z < 1$.

We consider two settings: (A) $n = 500$, $\beta_{1,2} = 2.0$, $\beta_{1,3} = -2.0$ and $\beta_{1,j} = 0$ for $j > 2$; (B) $n = 600$, $\beta_{1,2} = \beta_{1,3} = \beta_{1,4} = \beta_{1,5} = \beta_{1,6} = 1.0$, $\beta_{1,j} = 0$ for $j > 6$. As in Section 4.1, we set $p = 1000$ and consider two different covariance matrices, $\Sigma = \mathbf{I}$ and

Table 1: ECP (%) of the CIs with standard errors in parenthesis

Setting (A)		ROSE	DLASSO	BLPR	S3
$\Sigma = \mathbf{I}_p$	β_3	93.0 (1.1)	94.0 (1.1)	83.0 (1.7)	94.0 (1.1)
	β_4	96.4 (0.8)	96.0 (0.9)	97.4 (0.7)	95.2 (1.0)
$\Sigma = \}0.5^{ i-j } _{i,j}$	β_3	93.6 (1.1)	89.0 (1.4)	92.0 (1.2)	94.6 (1.0)
	β_4	94.6 (1.0)	86.0 (1.6)	95.4 (0.9)	93.4 (1.1)
Setting (B)		ROSE	DLASSO	BLPR	S3
$\Sigma = \mathbf{I}_p$	β_3	94.0 (1.1)	94.0 (1.1)	87.0 (1.5)	94.4 (1.0)
	β_4	96.8 (0.8)	96.0 (0.9)	97.4 (0.7)	95.2 (1.0)
$\Sigma = \}0.5^{ i-j } _{i,j}$	β_3	93.8 (1.1)	89.0 (1.4)	93.4 (1.1)	94.4 (1.0)
	β_4	95.6 (0.9)	85.6 (1.6)	96.8 (0.8)	95.2 (1.0)
Setting (C)		ROSE	DLASSO	BLPR	S3
$\Sigma = \mathbf{I}_p$	β_3	95.6 (0.9)	95.8 (0.9)	94.6 (1.0)	93.8 (1.1)
	β_4	94.8 (1.0)	95.2 (1.0)	96.0 (0.9)	96.6 (0.8)
$\Sigma = \}0.5^{ i-j } _{i,j}$	β_3	94.8 (1.0)	76.4 (1.9)	90.8 (4.0)	96.4 (0.8)
	β_4	94.0 (1.1)	92.0 (1.2)	95.6 (0.9)	93.6 (1.1)
Setting (D)		ROSE	DLASSO	BLPR	S3
$\Sigma = \mathbf{I}_p$	β_4	94.8 (1.0)	94.0 (1.1)	92.6 (1.2)	95.2 (1.0)
	β_5	93.8 (1.1)	93.4 (1.1)	91.0 (1.3)	93.8 (1.1)
	β_6	96.2 (0.9)	95.4 (0.9)	91.0 (1.3)	95.2 (1.0)
	$\beta_=\$	94.4 (1.0)	95.2 (1.0)	95.0 (1.0)	95.6 (0.9)
$\Sigma = \}0.5^{ i-j } _{i,j}$	β_4	96.0 (0.9)	81.2 (1.7)	93.2 (1.1)	94.6 (1.0)
	β_5	93.4 (1.1)	82.6 (1.7)	94.8 (1.0)	93.8 (1.1)
	β_6	94.6 (1.0)	91.0 (1.2)	93.6 (1.1)	96.4 (0.8)
	$\beta_=\$	93.8 (1.1)	91.6 (1.3)	95.0 (1.0)	95.4 (0.9)

Table 2: AL of the CIs with standard errors in parenthesis (numbers reported in the table are multiplied by 100)

Setting (A)		ROSE	DLASSO	BLPR	S3
$\Sigma = \mathbf{I}_p$	β_3	42.2 (0.3)	45.4 (0.3)	88.5 (1.4)	63.1 (0.5)
	β_4	42.4 (0.3)	45.1 (0.2)	2.6 (0.1)	64.1 (0.5)
$\Sigma = \}0.5^{ i-j } \}_{i,j}$	β_3	48.0 (0.3)	47.5 (0.2)	106.2 (1.9)	72.4 (0.5)
	β_4	49.1 (0.3)	47.8 (0.2)	5.4 (0.3)	74.8 (0.6)
Setting (B)		ROSE	DLASSO	BLPR	S3
$\Sigma = \mathbf{I}_p$	β_3	40.6 (0.2)	45.4 (0.2)	154.8 (3.4)	60.2 (0.5)
	β_4	40.7 (0.2)	45.1 (0.2)	3.4 (0.1)	60.5 (0.4)
$\Sigma = \}0.5^{ i-j } \}_{i,j}$	β_3	46.7 (0.2)	47.6 (0.2)	193.5 (4.4)	69.2 (0.5)
	β_4	47.6 (0.3)	47.8 (0.2)	6.8 (0.3)	70.9 (0.6)
Setting (C)		ROSE	DLASSO	BLPR	S3
$\Sigma = \mathbf{I}_p$	β_3	27.9 (0.1)	30.1 (0.1)	142.8 (3.3)	40.0 (0.2)
	β_4	28.0 (0.1)	30.2 (0.1)	4.7 (0.1)	40.5 (0.2)
$\Sigma = \}0.5^{ i-j } \}_{i,j}$	β_3	32.2 (0.1)	34.5 (0.1)	161.6 (2.8)	46.3 (0.2)
	β_4	32.4 (0.1)	34.7 (0.1)	9.8 (0.3)	46.6 (0.2)
Setting (D)		ROSE	DLASSO	BLPR	S3
$\Sigma = \mathbf{I}_p$	β_4	28.4 (0.1)	31.4 (0.1)	65.6 (1.6)	42.7 (0.2)
	β_5	28.5 (0.1)	31.4 (0.1)	65.6 (1.6)	42.6 (0.2)
	β_6	28.3 (0.1)	31.3 (0.1)	63.9 (1.6)	42.3 (0.2)
	β_7	28.5 (0.1)	31.4 (0.1)	3.8 (0.1)	42.7 (0.2)
$\Sigma = \}0.5^{ i-j } \}_{i,j}$	β_4	37.0 (0.1)	34.0 (0.1)	55.1 (1.3)	55.6 (0.3)
	β_5	36.9 (0.1)	33.9 (0.1)	68.9 (1.4)	55.7 (0.3)
	β_6	33.2 (0.1)	33.9 (0.1)	84.5 (1.6)	50.3 (0.2)
	β_7	33.1 (0.1)	33.8 (0.1)	6.5 (0.2)	50.1 (0.3)

Table 3: ECP and AL of the CIs, with standard errors in parenthesis

Setting (A)		ROSE		DLASSO		S3	
Σ		ECP(%)	AL*100	ECP(%)	AL*100	ECP(%)	AL*100
I_p	β_3	95.8 (0.9)	80.5 (0.3)	13.0 (1.5)	53.4 (0.2)	93.6 (1.1)	118.8 (0.8)
	β_4	94.8 (1.0)	51.2 (0.1)	97.2 (0.7)	45.2 (0.1)	93.8 (1.1)	75.5 (0.3)
$\}0.5^{ i-j } _{i,j}$	β_3	95.8 (0.9)	77.1 (0.3)	26.6 (2.0)	52.9 (0.1)	95.2 (1.0)	113.0 (0.6)
	β_4	94.8 (1.0)	52.9 (0.1)	95.4 (0.9)	47.5 (0.1)	95.8 (0.9)	76.9 (0.2)
Setting (B)		ROSE		DLASSO		S3	
I_p	β_4	94.0 (1.1)	51.1 (0.1)	30.8 (1.1)	39.9 (0.1)	95.2 (1.0)	76.2 (0.3)
	β_5	93.2 (1.1)	50.9 (0.1)	26.0 (1.1)	39.8 (0.1)	95.6 (0.9)	75.7 (0.3)
	β_6	96.2 (0.9)	50.9 (0.1)	30.6 (0.9)	39.8 (0.1)	94.8 (1.0)	76.2 (0.3)
	$\beta_{=}$	93.4 (1.0)	43.9 (0.1)	96.4 (1.0)	38.1 (0.1)	92.4 (1.2)	64.9 (0.2)
$\}0.5^{ i-j } _{i,j}$	β_4	95.6 (0.9)	71.5 (0.2)	88.2 (1.4)	55.7 (0.2)	95.0 (1.0)	107.5 (0.5)
	β_5	93.2 (1.1)	71.5 (0.2)	84.6 (1.6)	54.9 (0.2)	93.8 (1.1)	108.0 (0.5)
	β_6	93.8 (1.1)	65.5 (0.2)	67.6 (2.1)	53.1 (0.1)	93.6 (1.1)	99.1 (0.5)
	$\beta_{=}$	94.0 (1.1)	58.5 (0.2)	95.0 (1.0)	50.8 (0.1)	94.0 (1.1)	88.2 (0.4)

$\Sigma = \}0.5^{|i-j|} |_{i,j[2,\dots,p}$. The objective is to construct two-sided CIs for $\beta_{1,3}, \beta_{1,4}$ in Setting (A) and $\beta_{1,4}, \beta_{1,5}, \beta_{1,6}, \beta_{1,=}$ in Setting (B).

To implement the proposed CI in (12), we set $s_n = \lfloor 2n/\log(n) \rfloor$ and $l = 5$. We use the R package `SIS` and estimate β_{j_0} by ISIS. The initial estimator $\tilde{\beta}$ is computed by fitting a penalized logistic regression model with SCAD penalty function for the variables picked by ISIS. We implement the penalized logistic regression by the R package `ncvreg`. For Setting (B), we update $\widehat{\beta}_{j_0}^{t+}$ using the method discussed in Section 3.3 with $S = 2$.

We further compare the proposed CI with the CI constructed by the DLASSO method and the S3 method. In Table 3, we report the ECP and AL of the proposed CI and the CIs constructed by DLASSO and S3. It can be seen that DLASSO performs poorly for nonzero parameters. On the contrary, ECPs of the proposed CIs are close to the nominal level in almost all cases. In addition, our CIs are much narrower than those based on the S3 method in all cases.

4.3 Real data analysis

We apply the proposed methods to a real dataset riboflavin (vitamin B2) production in *Bacillus subtilis*. This dataset is provided by DSM (Kaiseraugst, Switzerland) and is pub-

licly available in the R package `hdi`. It consists of a response variable which is the logarithm of the riboflavin production rate and 4088 predictors measuring the logarithm of the expression level of 4088 genes. There are a total of 71 observations. We model this data with a linear regression model, center the response and standardize all the covariates before analysis. To identify genes that are significantly associated with the response, we construct CIs for each individual coefficient and apply Bonferroni’s method for multiple adjustment. We compare the proposed method with the de-sparsified Lasso method and implement both methods as discussed in Section 4.1. At the 5% significance level, the proposed method finds three important genes (the 1588th, 3154th and 4004th) while the de-sparsified Lasso procedure claims no variables are significant.

5 Discussion

5.1 Statistical inference via online estimation

In this paper, we develop an online estimation procedure for high-dimensional statistical inference, to account for model selection uncertainty in subsequent inferences. Such an online inference method can be applied to some other non-regular problems as well. Variations of this approach has been used by Luedtke and van der Laan (2016) to provide a CI for the mean outcome under a non-unique optimal treatment regime, and by Luedtke and van der Laan (2017) to construct a CI for the maximal absolute correlation between responses and covariates.

5.2 Multi-dimensional extensions

We focus on constructing CIs for a single regression coefficient in GLMs. The proposed procedure can be naturally extended to form confidence regions for multi-dimensional parameters as well. Let \mathbb{J}_1 be an arbitrary subset of \mathbb{I} with $|\mathbb{J}_1| > 1$. The confidence region for β_{1,\mathbb{J}_0} can be constructed as follows.

Let $\mathbb{J}_0 = \{j \in \mathbb{I} : \beta_{1,j} \neq 0\}$. We first estimate \mathbb{J}_0 by some model selection procedure based on the sub-dataset $\mathcal{M} = \{(\mathbf{X}_2, Y_2), \dots, (\mathbf{X}_t, Y_t)\}$ for $t = s_n, \dots, n - 1$

and $\{(\mathbf{X}_{s_n 0 2}, Y_{s_n 0 2}), \dots, (\mathbf{X}_n, Y_n)\}$. Denoted by $\widehat{\mathbb{J}}_0^{s_n+}$, $\widehat{\mathbb{J}}_0^{s_n 0 2+}$, \dots , $\widehat{\mathbb{J}}_0^{n-2+}$ and $\widehat{\mathbb{J}}_0^{-s_n+}$ the corresponding estimators. We calculate $\widehat{\Sigma}$ as in (10) based on some consistent initial estimator $\widetilde{\beta}$ and compute

$$\widehat{\omega}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(t)}, \mathbb{J}_0} = \widehat{\Sigma}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(t)}, \widehat{\mathcal{M}}_{\mathbb{J}_0}^{(t)}}^{-2} \widehat{\Sigma}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(t)}, \mathbb{J}_0}, \quad \widehat{\sigma}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(t)}, \mathbb{J}_0} = \left(\widehat{\Sigma}_{\mathbb{J}_0, \mathbb{J}_0} \quad \widehat{\Sigma}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(t)}, \mathbb{J}_0}^T \widehat{\omega}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(t)}, \mathbb{J}_0} \right)^{2/3},$$

for $t = s_n, \dots, n-1$ and

$$\widehat{\omega}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(-s_n)}, \mathbb{J}_0} = \widehat{\Sigma}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(-s_n)}, \widehat{\mathcal{M}}_{\mathbb{J}_0}^{(-s_n)}}^{-2} \widehat{\Sigma}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(-s_n)}, \mathbb{J}_0}, \quad \widehat{\sigma}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(-s_n)}, \mathbb{J}_0} = \left(\widehat{\Sigma}_{\mathbb{J}_0, \mathbb{J}_0} \quad \widehat{\Sigma}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(-s_n)}, \mathbb{J}_0}^T \widehat{\omega}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(-s_n)}, \mathbb{J}_0} \right)^{2/3}.$$

Consider the following score equation:

$$\begin{aligned} & \sum_{t=1}^{s_n-2} \widehat{\sigma}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(-s_n)}, \mathbb{J}_0}^{-2} \widehat{\mathbf{Z}}_{t0 2, \mathbb{J}_0}^T \left\{ Y_{t0 2} - \mu \left(\mathbf{X}_{t0 2, \mathbb{J}_0} \beta_{1, \mathbb{J}_0} + \mathbf{X}_{t0 2, \widehat{\mathcal{M}}_{\mathbb{J}_0}^{(-s_n)}} \widetilde{\beta}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(-s_n)}} \right) \right\} \\ & + \sum_{t=s_n}^{n-2} \widehat{\sigma}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(t)}, \mathbb{J}_0}^{-2} \widehat{\mathbf{Z}}_{t0 2, \mathbb{J}_0}^T \left\{ Y_{t0 2} - \mu \left(\mathbf{X}_{t0 2, \mathbb{J}_0} \beta_{1, \mathbb{J}_0} + \mathbf{X}_{t0 2, \widehat{\mathcal{M}}_{\mathbb{J}_0}^{(t)}} \widetilde{\beta}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(t)}} \right) \right\} = 0, \end{aligned}$$

where $\widehat{\mathbf{Z}}_{t0 2, \mathbb{J}_0} = \mathbf{X}_{t0 2, \mathbb{J}_0} \quad \widehat{\omega}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(t)}, \mathbb{J}_0}^T \mathbf{X}_{t0 2, \widehat{\mathcal{M}}_{\mathbb{J}_0}^{(t)}}$ for $t = s_n, \dots, n-1$ and $\widehat{\mathbf{Z}}_{t0 2, \mathbb{J}_0} = \mathbf{X}_{t0 2, \mathbb{J}_0} \quad \widehat{\omega}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(-s_n)}, \mathbb{J}_0}^T \mathbf{X}_{t0 2, \widehat{\mathcal{M}}_{\mathbb{J}_0}^{(-s_n)}}$ for $t = 0, \dots, s_n-1$. The estimator $\widehat{\beta}_{\mathcal{J}_0}$ can be computed by solving the score equation via Newton's method with initial value $\widetilde{\beta}_{\mathcal{J}_0}$. The corresponding $1-\alpha$ 100% confidence region is given by

$$\left\{ \beta_{\mathbb{J}_0} / \mathbb{R}^{|\mathbb{J}_0|} : n(\beta_{\mathbb{J}_0} - \widehat{\beta}_{\mathbb{J}_0})^T (\Gamma_n^*)^T \Gamma_n^* (\beta_{\mathbb{J}_0} - \widehat{\beta}_{\mathbb{J}_0}) \widehat{\phi} \geq \chi_\alpha^3(\|\mathbb{J}_1\|) \right\}, \quad (26)$$

where $\widehat{\phi}$ denotes some constant estimator for ϕ_1 , $\chi_\alpha^3(\|\mathbb{J}_1\|)$ is the upper α -quantile of a central χ^3 distribution with $\|\mathbb{J}_1\|$ degrees of freedom, and

$$\begin{aligned} \Gamma_n^* &= \frac{1}{n} \sum_{t=1}^{s_n-2} \widehat{\sigma}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(-s_n)}, \mathbb{J}_0}^{-2} \widehat{\mathbf{Z}}_{t0 2, \mathbb{J}_0}^T b'' \left(\mathbf{X}_{t0 2, \mathbb{J}_0} \widehat{\beta}_{1, \mathbb{J}_0} + \mathbf{X}_{t0 2, \widehat{\mathcal{M}}_{\mathbb{J}_0}^{(-s_n)}} \widetilde{\beta}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(-s_n)}} \right) \\ &+ \frac{1}{n} \sum_{t=s_n}^{n-2} \widehat{\sigma}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(t)}, \mathbb{J}_0}^{-2} \widehat{\mathbf{Z}}_{t0 2, \mathbb{J}_0}^T b'' \left(\mathbf{X}_{t0 2, \mathbb{J}_0} \widehat{\beta}_{1, \mathbb{J}_0} + \mathbf{X}_{t0 2, \widehat{\mathcal{M}}_{\mathbb{J}_0}^{(t)}} \widetilde{\beta}_{\widehat{\mathcal{M}}_{\mathbb{J}_0}^{(t)}} \right). \end{aligned}$$

To guarantee the validity of (26), the number of elements in \mathbb{J}_1 needs to be much smaller than n . It would be interesting to construct confidence regions for the entire regression coefficient vector β_1 based on some multiple comparison procedures. However, this is beyond the scope of the current paper.

5.3 Extension to generic penalized M-estimators

The proposed method can also be extended beyond the class of GLMs to a general framework with a convex loss function. Specifically, given a high-dimensional random vector \mathbf{U}_1 , define

$$\beta_1 = \arg \min_{\beta \in \mathbb{R}^p} \mathbb{E} \ell(\mathbf{U}_1, \beta),$$

for some convex loss function ℓ . An initial estimator for β_1 can be computed by minimizing

$$\tilde{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \left(\frac{1}{n} \sum_{i=2}^n \ell(\mathbf{U}_i, \beta) + \sum_{j=2}^p \rho_\lambda(\|\beta_j\|) \right), \quad (27)$$

where $\mathbf{U}_2, \dots, \mathbf{U}_n$ are i.i.d random vectors generated according as \mathbf{U}_1 , and $\rho_\lambda(\cdot)$ denotes some penalty function. In addition to estimating the regression coefficients in GLMs, such a generic framework includes some other important applications such as estimation of the precision matrix in Gaussian graphical models (as illustrated in Section 2.1.4 of Ning and Liu, 2017).

Here, we focus on constructing the CI for a univariate parameter β_{1,j_0} . Let $\widehat{[\]}_{j_0}^{t+}$ denote the estimated support of the control variables based on $\}\mathbf{U}_i\}_{i=2}^t$ for $t = s_n, s_n + 1, \dots, n - 1$ and $\widehat{[\]}_{j_0}^{-s_n+}$ the estimated support based on $\}\mathbf{U}_i\}_{i=[s_n, 0, 2}^n$. Define

$$\widehat{\Sigma} = \frac{1}{n} \sum_{i=2}^n \frac{\partial^3}{\partial \beta \partial \beta^T} \ell(\mathbf{U}_i, \tilde{\beta}),$$

where $\tilde{\beta}$ corresponds to the initial estimator in (27). Given $\widehat{\Sigma}$, we compute $\widehat{\omega}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}$, and $\widehat{\omega}_{\widehat{\mathcal{M}}_{j_0}^{(-s_n)}, j_0}$ as in Section 3.1. For any model $\llbracket \cdot \rrbracket \leq \mathbb{J}_{j_0}$, suppose we have some consistent

estimator $\hat{\sigma}_{\mathcal{M},j_0}^3$ for

$$\sigma_{\mathcal{M},j_0}^3 = \mathbb{E} \left(\frac{\partial \ell(\mathbf{U}_1, \boldsymbol{\beta}_1)}{\partial \beta_{j_0}} \quad \boldsymbol{\omega}_{\mathcal{M},j_0}^T \frac{\partial \ell(\mathbf{U}_1, \boldsymbol{\beta}_1)}{\partial \boldsymbol{\beta}_{\mathcal{M}}} \right)^3.$$

For any $c \in \mathbb{R}$, $\lfloor \cdot \rfloor \subseteq \mathbb{I}_{j_0}$ and $\boldsymbol{\alpha} \in \mathbb{R}^{|\mathcal{M}|}$, we define a p -dimensional vector $\boldsymbol{\theta} = \mathbf{h}(c, \lfloor \cdot \rfloor, \boldsymbol{\alpha})$ such that $\theta_{j_0} = c$, $\boldsymbol{\theta}_{\mathcal{M}} = \boldsymbol{\alpha}$ and $\boldsymbol{\theta}_{\mathcal{M}^c - \{j_0\}} = \mathbf{0}$. Let $\widehat{\lfloor \cdot \rfloor}_{j_0}^{t+} = \widehat{\lfloor \cdot \rfloor}_{j_0}^{-s_n+}$ for $t = 0, 1, \dots, s_n - 1$ and $\widehat{\beta}_{j_0}^{l+} = \widetilde{\beta}_{j_0}$, we update $\widehat{\beta}_{j_0}$ as

$$\widehat{\beta}_{j_0}^{l+} = \widehat{\beta}_{j_0}^{l-2+} \frac{\sum_{t \in \mathbb{I}_{j_0}} \frac{1}{n \widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}} \left(\frac{\partial \ell(\mathbf{U}_{t0}, \mathbf{h}(\widehat{\beta}_{j_0}^{l-2+}, \widehat{\lfloor \cdot \rfloor}_{j_0}^{t+}, \widetilde{\boldsymbol{\beta}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}}))}{\partial \beta_{j_0}} \quad \widehat{\boldsymbol{\omega}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}^T \frac{\partial \ell(\mathbf{U}_{t0}, \mathbf{h}(\widehat{\beta}_{j_0}^{l-2+}, \widehat{\lfloor \cdot \rfloor}_{j_0}^{t+}, \widetilde{\boldsymbol{\beta}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}}))}{\partial \boldsymbol{\beta}_{\widehat{\mathcal{M}}_{j_0}^{(t)}}} \right)}{\sum_{t \in \mathbb{I}_{j_0}} \frac{1}{n \widehat{\sigma}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}} \left(\frac{\partial^3 \ell(\mathbf{U}_{t0}, \mathbf{h}(\widehat{\beta}_{j_0}^{l-2+}, \widehat{\lfloor \cdot \rfloor}_{j_0}^{t+}, \widetilde{\boldsymbol{\beta}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}}))}{\partial \beta_{j_0}^3} \quad \widehat{\boldsymbol{\omega}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}, j_0}^T \frac{\partial^3 \ell(\mathbf{U}_{t0}, \mathbf{h}(\widehat{\beta}_{j_0}^{l-2+}, \widehat{\lfloor \cdot \rfloor}_{j_0}^{t+}, \widetilde{\boldsymbol{\beta}}_{\widehat{\mathcal{M}}_{j_0}^{(t)}}))}{\partial \beta_{j_0} \partial \boldsymbol{\beta}_{\widehat{\mathcal{M}}_{j_0}^{(t)}}} \right)},$$

$\underbrace{\hspace{15em}}_{n^{*(l-1)}}$

for $l = 1, 2, \dots$. The corresponding CI for β_{1,j_0} is given by

$$\widehat{\beta}_{j_0}^{l+} \subseteq \frac{z_{\frac{\alpha}{2}}}{n \Gamma_n^{*(l-2+)}}.$$

In Section C of the supplementary article, we sketch a few lines to show that the above CI achieves nominal coverage under certain conditions.

5.4 Doubly-robust procedure

The proposed ROSE algorithm constructs the score equation for β_{1,j_0} by recursively estimating the support of control variables. As commented in (2.2), such a procedure requires certain minimal-signal-strength conditions on $\boldsymbol{\beta}_{1, \mathbb{I}_{j_0}}$.

We now introduce a variant of our method that is valid even when the minimal-signal-strength conditions fail. At the t -th iteration, instead of estimating $\lfloor \cdot \rfloor_{j_0}$ only, we might apply another variable selection procedure to estimate the support of $\boldsymbol{\omega}_{\mathbb{I}_{j_0}, j_0}$ based on \mathcal{M} and set $\widehat{\lfloor \cdot \rfloor}_{j_0}^{t+}$ to be a union of the two sets of important variables selected. The result-

ing CI is doubly-robust in the sense that it achieves nominal coverage as long as either $\beta_{1, \mathbb{I}_{j_0}}$ satisfies certain minimal-signal-strength conditions, or the ℓ_3 norm of weak signals in $\beta_{1, \mathbb{I}_{j_0}}$ and $\omega_{\mathbb{I}_{j_0}, j_0}$ is $o(n^{-2/5})$. The latter condition allows the existence of weak signals in $\beta_{1, \mathbb{I}_{j_0}}$. It automatically holds when variables with signals larger than or proportional to $(n/\log \log n)^{-2/5}(s^*)^{-2/3}$ can be consistently identified by the model selection procedure. In addition, it is considerably weaker than the zonal assumption (Bühlmann and Mandozzi, 2014) that requires the strength of weak signals to be $o(n^{-2/3})$. More detailed discussions are given in Section B.1.3 of the supplementary article.

References

- Bickel, P. J. (1975). One-step Huber estimates in the linear model. *J. Amer. Statist. Assoc.* *70*, 428–434.
- Bickel, P. J., Y. Ritov, and A. B. Tsybakov (2009). Simultaneous analysis of lasso and Dantzig selector. *Ann. Statist.* *37*(4), 1705–1732.
- Breheny, P. and J. Huang (2011). Coordinate descent algorithms for nonconvex penalized regression, with applications to biological feature selection. *Ann. Appl. Stat.* *5*(1), 232–253.
- Bühlmann, P. and J. Mandozzi (2014). High-dimensional variable screening and bias in subsequent inference, with an empirical comparison. *Comput. Statist.* *29*(3-4), 407–430.
- Candes, E. and T. Tao (2007). The dantzig selector: statistical estimation when p is much larger than n (with discussions). *Ann. Stat.* *35*, 2313–2404.
- Csörgö, M. (1968). On the strong law of large numbers and the central limit theorem for martingales. *Trans. Amer. Math. Soc.* *131*, 259–275.
- Dezeure, R., P. Bühlmann, L. Meier, and N. Meinshausen (2015). High-dimensional inference: confidence intervals, p -values and **r**-software **hdi**. *Statist. Sci.* *30*(4), 533–558.

- Fan, J., S. Guo, and N. Hao (2012). Variance estimation using refitted cross-validation in ultrahigh dimensional regression. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* *74*(1), 37–65.
- Fan, J. and R. Li (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Amer. Statist. Assoc.* *96*(456), 1348–1360.
- Fan, J. and J. Lv (2008). Sure independence screening for ultrahigh dimensional feature space. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* *70*(5), 849–911.
- Javanmard, A. and A. Montanari (2014). Confidence intervals and hypothesis testing for high-dimensional regression. *J. Mach. Learn. Res.* *15*, 2869–2909.
- Lee, J. D., D. L. Sun, Y. Sun, and J. E. Taylor (2016). Exact post-selection inference, with application to the lasso. *Ann. Statist.* *44*(3), 907–927.
- Liu, H., X. Xu, and J. J. Li (2017). A bootstrap lasso+ partial ridge method to construct confidence intervals for parameters in high-dimensional sparse linear models. *arXiv preprint arXiv:1706.02150*.
- Liu, H. and B. Yu (2013). Asymptotic properties of Lasso+mLS and Lasso+Ridge in sparse high-dimensional linear regression. *Electron. J. Stat.* *7*, 3124–3169.
- Luedtke, A. R. and M. J. van der Laan (2016). Statistical inference for the mean outcome under a possibly non-unique optimal treatment strategy. *Ann. Statist.* *44*(2), 713–742.
- Luedtke, A. R. and M. J. van der Laan (2017). Parametric-rate inference for one-sided differentiable parameters. *J. Amer. Statist. Assoc.* (just-accepted).
- McCullagh and Nelder (1989). *Generalized Linear Models*. Chapman and Hall.
- Ning, Y. and H. Liu (2017). A general theory of hypothesis tests and confidence regions for sparse high dimensional models. *Ann. Statist.* *45*(1), 158–195.
- Saldana, D. F. and Y. Feng (2016). Sis: An r package for sure independence screening in ultrahigh dimensional statistical models. *Journal of Statistical Software*, to appear.

- Schifano, E. D., J. Wu, C. Wang, J. Yan, and M.-H. Chen (2016). Online updating of statistical inference in the big data setting. *Technometrics* 58(3), 393–403.
- Sun, T. and C.-H. Zhang (2013). Sparse matrix inversion with scaled lasso. *J. Mach. Learn. Res.* 14, 3385–3418.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* 58, 267–288.
- van de Geer, S., P. Bühlmann, Y. Ritov, and R. Dezeure (2014). On asymptotically optimal confidence regions and tests for high-dimensional models. *Ann. Statist.* 42(3), 1166–1202.
- Wang, C., M.-H. Chen, E. Schifano, J. Wu, and J. Yan (2016). Statistical methods and computing for big data. *Stat. Interface* 9(4), 399–414.
- Wasserman, L. and K. Roeder (2009). High-dimensional variable selection. *Ann. Statist.* 37(5A), 2178–2201.
- Zhang, C.-H. (2010). Nearly unbiased variable selection under minimax concave penalty. *Ann. Statist.* 38(2), 894–942.
- Zhang, C.-H. and S. S. Zhang (2014). Confidence intervals for low dimensional parameters in high dimensional linear models. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* 76(1), 217–242.