



5-2017

# Confidence Intervals for High-Dimensional Linear Regression: Minimax Rates and Adaptivity

Tony Cai  
*University of Pennsylvania*

Zijian Guo  
*University of Pennsylvania*

Follow this and additional works at: [http://repository.upenn.edu/statistics\\_papers](http://repository.upenn.edu/statistics_papers)

 Part of the [Physical Sciences and Mathematics Commons](#)

---

## Recommended Citation

Cai, T., & Guo, Z. (2017). Confidence Intervals for High-Dimensional Linear Regression: Minimax Rates and Adaptivity. *The Annals of Statistics*, 45 (2), 615-646. <http://dx.doi.org/10.1214/16-AOS1461>

This paper is posted at ScholarlyCommons. [http://repository.upenn.edu/statistics\\_papers/73](http://repository.upenn.edu/statistics_papers/73)  
For more information, please contact [repository@pobox.upenn.edu](mailto:repository@pobox.upenn.edu).

---

# Confidence Intervals for High-Dimensional Linear Regression: Minimax Rates and Adaptivity

## **Abstract**

Confidence sets play a fundamental role in statistical inference. In this paper, we consider confidence intervals for high-dimensional linear regression with random design. We first establish the convergence rates of the minimax expected length for confidence intervals in the oracle setting where the sparsity parameter is given. The focus is then on the problem of adaptation to sparsity for the construction of confidence intervals. Ideally, an adaptive confidence interval should have its length automatically adjusted to the sparsity of the unknown regression vector, while maintaining a pre-specified coverage probability. It is shown that such a goal is in general not attainable, except when the sparsity parameter is restricted to a small region over which the confidence intervals have the optimal length of the usual parametric rate. It is further demonstrated that the lack of adaptivity is not due to the conservativeness of the minimax framework, but is fundamentally caused by the difficulty of learning the bias accurately.

## **Keywords**

Adaptivity, confidence interval, coverage probability, expected length, high-dimensional linear regression, minimaxity, sparsity

## **Disciplines**

Physical Sciences and Mathematics

# CONFIDENCE INTERVALS FOR HIGH-DIMENSIONAL LINEAR REGRESSION: MINIMAX RATES AND ADAPTIVITY\*

BY T. TONY CAI, AND ZIJIAN GUO

*University of Pennsylvania*

Confidence sets play a fundamental role in statistical inference. In this paper, we consider confidence intervals for high dimensional linear regression with random design. We first establish the convergence rates of the minimax expected length for confidence intervals in the oracle setting where the sparsity parameter is given. The focus is then on the problem of adaptation to sparsity for the construction of confidence intervals. Ideally, an adaptive confidence interval should have its length automatically adjusted to the sparsity of the unknown regression vector, while maintaining a prespecified coverage probability. It is shown that such a goal is in general not attainable, except when the sparsity parameter is restricted to a small region over which the confidence intervals have the optimal length of the usual parametric rate. It is further demonstrated that the lack of adaptivity is not due to the conservativeness of the minimax framework, but is fundamentally caused by the difficulty of learning the bias accurately.

**1. Introduction.** Driven by a wide range of applications, high-dimensional linear regression, where the dimension  $p$  can be much larger than the sample size  $n$ , has received significant recent attention. The linear model is

$$(1.1) \quad y = X\beta + \epsilon, \quad \epsilon \sim N(0, \sigma^2 I),$$

where  $y \in \mathbb{R}^n$ ,  $X \in \mathbb{R}^{n \times p}$ , and  $\beta \in \mathbb{R}^p$ . Several penalized/constrained  $\ell_1$  minimization methods, including the Lasso [22], Dantzig Selector [11], square-root Lasso [1], and scaled Lasso [21] have been proposed and studied. Under regularity conditions on the design matrix  $X$ , these methods with a suitable choice of the tuning parameter have been shown to achieve the optimal rate of convergence  $k \frac{\log p}{n}$  under the squared error loss over the set of  $k$ -sparse regression coefficient vectors with  $k \leq c \frac{n}{\log p}$  where  $c > 0$  is a constant. That

---

\*The research was supported in part by NSF Grants DMS-1208982 and DMS-1403708, and NIH Grant R01 CA127334..

*MSC 2010 subject classifications:* Primary 62G15; secondary 62C20, 62H35

*Keywords and phrases:* Adaptivity, confidence interval, coverage probability, expected length, high-dimensional linear regression, minimaxity, sparsity.

is, there exists some constant  $C > 0$  such that

$$(1.2) \quad \sup_{\|\beta\|_0 \leq k} \mathbb{P} \left( \|\widehat{\beta} - \beta\|_2^2 > Ck \frac{\log p}{n} \right) = o(1),$$

where  $\|\beta\|_0$  denotes the number of the nonzero coordinates of a vector  $\beta \in \mathbb{R}^p$ . See, for example, [24, 2, 11, 21]. A key feature of the estimation problem is that the optimal rate can be achieved adaptively with respect to the sparsity parameter  $k$ .

Confidence sets play a fundamental role in statistical inference and confidence intervals for high-dimensional linear regression have been actively studied recently with a focus on inference for individual coordinates. But, compared to point estimation, there is still a paucity of methods and fundamental theoretical results on confidence intervals for high-dimensional regression. Zhang and Zhang [25] was the first to introduce the idea of de-biasing for constructing a valid confidence interval for a single coordinate  $\beta_i$ . The confidence interval is centered at a low-dimensional projection estimator obtained through bias correction via score vector using the scaled Lasso as the initial estimator. [14, 15, 23] also used de-biasing for the construction of confidence intervals and [23] established asymptotic efficiency for the proposed estimator. All the aforementioned papers [25, 14, 15, 23] have focused on the ultra-sparse case where the sparsity  $k \ll \frac{\sqrt{n}}{\log p}$  is assumed. Under such a sparsity condition, the expected length of the confidence intervals constructed in [25, 15, 23] is at the parametric rate  $\frac{1}{\sqrt{n}}$  and the procedures do not depend on the specific value of  $k$ .

Compared to point estimation where the sparsity condition  $k \ll \frac{n}{\log p}$  is sufficient for estimation consistency (see equation (1.2)), the condition  $k \ll \frac{\sqrt{n}}{\log p}$  for valid confidence intervals is much stronger. There are several natural questions: What happens in the region where  $\frac{\sqrt{n}}{\log p} \lesssim k \lesssim \frac{n}{\log p}$ ? Is it still possible to construct a valid confidence interval for  $\beta_i$  in this case? Can one construct an adaptive honest confidence interval not depending on  $k$ ?

The goal of the present paper is to address these and other related questions on confidence intervals for high-dimensional linear regression with random design. More specifically, we consider construction of confidence intervals for a linear functional  $T(\beta) = \xi^\top \beta$ , where the loading vector  $\xi \in \mathbb{R}^p$  is given and  $\frac{\max_{i \in \text{supp}(\xi)} |\xi_i|}{\min_{i \in \text{supp}(\xi)} |\xi_i|} \leq \bar{c}$  with  $\bar{c} \geq 1$  being a constant. Based on the sparsity of  $\xi$ , we focus on two specific regimes: the sparse loading regime where  $\|\xi\|_0 \leq Ck$ , with  $C > 0$  being a constant; the dense loading regime where  $\|\xi\|_0$  satisfying (2.7) in Section 2. It will be seen later that for confidence intervals  $T(\beta) = \beta_i$  is a prototypical case for the general functional

$T(\beta) = \xi^\top \beta$  with a sparse loading  $\xi$ , and  $T(\beta) = \sum_{i=1}^p \beta_i$  is a representative case for  $T(\beta) = \xi^\top \beta$  with a dense loading  $\xi$ .

To illustrate the main idea, let us first focus on the two specific functionals  $T(\beta) = \beta_i$  and  $T(\beta) = \sum_{i=1}^p \beta_i$ . We establish the convergence rate of the minimax expected length for confidence intervals in the oracle setting where the sparsity parameter  $k$  is given. It is shown that in this case the minimax expected length is of order  $\frac{1}{\sqrt{n}} + k \frac{\log p}{n}$  for confidence intervals for  $\beta_i$ . An honest confidence interval, which depends on the sparsity  $k$ , is constructed and is shown to be minimax rate optimal. To the best of our knowledge, this is the first construction of confidence intervals in the moderate-sparse region  $\frac{\sqrt{n}}{\log p} \ll k \lesssim \frac{n}{\log p}$ . If the sparsity  $k$  falls into the ultra-sparse region  $k \lesssim \frac{\sqrt{n}}{\log p}$ , the constructed confidence interval is similar to the confidence intervals constructed in [25, 15, 23]. On the other hand, the convergence rate of the minimax expected length of honest confidence intervals for  $\sum_{i=1}^p \beta_i$  in the oracle setting is shown to be  $k \sqrt{\frac{\log p}{n}}$ . A rate-optimal confidence interval that also depends on  $k$  is constructed. It should be noted that this confidence interval is not based on the de-biased estimator.

One drawback of the constructed confidence intervals mentioned above is that they require prior knowledge of the sparsity  $k$ . Such knowledge of sparsity is usually unavailable in applications. A natural question is: Without knowing the sparsity  $k$ , is it possible to construct a confidence interval as good as when the sparsity  $k$  is known? This is a question about adaptive inference, which has been a major goal in nonparametric and high-dimensional statistics. Ideally, an adaptive confidence interval should have its length automatically adjusted to the true sparsity of the unknown regression vector, while maintaining a prespecified coverage probability. We show that, unlike point estimation, such a goal is in general not attainable for confidence intervals. In the case of confidence intervals for  $\beta_i$ , it is impossible to adapt between different sparsity levels, except when the sparsity  $k$  is restricted to the ultra-sparse region  $k \lesssim \frac{\sqrt{n}}{\log p}$ , over which the confidence intervals have the optimal length of the parametric rate  $\frac{1}{\sqrt{n}}$ , which does not depend on  $k$ . In the case of confidence intervals for  $\sum_{i=1}^p \beta_i$ , it is shown that adaptation to the sparsity is not possible at all, even in the ultra-sparse region  $k \lesssim \frac{\sqrt{n}}{\log p}$ .

Minimax theory is often criticized as being too conservative as it focuses on the worst case performance. For confidence intervals for high dimensional linear regression, we establish strong non-adaptivity results which demonstrate that the lack of adaptivity is not due to the conservativeness of the minimax framework. It shows that for any confidence interval with guaranteed coverage probability over the set of  $k$  sparse vectors, its expected length

at any given point in a large subset of the parameter space must be at least of the same order as the minimax expected length. So the confidence interval must be long at a large subset of points in the parameter space, not just at a small number of “unlucky” points. This leads directly to the impossibility of adaptation over different sparsity levels. Fundamentally, the lack of adaptivity is caused by the difficulty in accurately learning the bias of any estimator for high-dimensional linear regression.

We now turn to confidence intervals for general linear functionals. For a linear functional  $\xi^\top \beta$  in the sparse loading regime, the rate of the minimax expected length is  $\|\xi\|_2 \left( \frac{1}{\sqrt{n}} + k \frac{\log p}{n} \right)$ , where  $\|\xi\|_2$  is the vector  $\ell_2$  norm of  $\xi$ . For a linear functional  $\xi^\top \beta$  in the dense loading regime, the rate of the minimax expected length is  $\|\xi\|_\infty k \sqrt{\frac{\log p}{n}}$ , where  $\|\xi\|_\infty$  is the vector  $\ell_\infty$  norm of  $\xi$ . Regarding adaptivity, the phenomena observed in confidence intervals for the two special linear functionals  $T(\beta) = \beta_i$  and  $T(\beta) = \sum_{i=1}^p \beta_i$  extend to the general linear functionals. The case of confidence intervals for  $T(\beta) = \sum_{i=1}^p \xi_i \beta_i$  with a sparse loading  $\xi$  is similar to that of confidence intervals for  $\beta_i$  in the sense that rate-optimal adaptation is impossible except when the sparsity  $k$  is restricted to the ultra-sparse region  $k \lesssim \frac{\sqrt{n}}{\log p}$ . On the other hand, the case for a dense loading  $\xi$  is similar to that of confidence intervals for  $\sum_{i=1}^p \beta_i$ : adaptation to the sparsity  $k$  is not possible at all, even in the ultra-sparse region  $k \lesssim \frac{\sqrt{n}}{\log p}$ .

In addition to the more typical setting in practice where the covariance matrix  $\Sigma$  of the random design and the noise level  $\sigma$  of the linear model are unknown, we also consider the case with the prior knowledge of  $\Sigma = I$  and  $\sigma = \sigma_0$ . It turns out that this case is strikingly different. The minimax rate for the expected length in the sparse loading regime is reduced from  $\|\xi\|_2 \left( \frac{1}{\sqrt{n}} + k \frac{\log p}{n} \right)$  to  $\frac{\|\xi\|_2}{\sqrt{n}}$ , and in particular it does not depend on the sparsity  $k$ . Furthermore, in marked contrast to the case of unknown  $\Sigma$  and  $\sigma$ , adaptation to sparsity is possible over the full range  $k \lesssim \frac{n}{\log p}$ . On the other hand, for linear functionals  $\xi^\top \beta$  with a dense loading  $\xi$ , the minimax rates and impossibility for adaptive confidence intervals do not change even with the prior knowledge of  $\Sigma = I$  and  $\sigma = \sigma_0$ . However, the cost of adaptation is reduced with the prior knowledge.

The rest of the paper is organized as follows: After basic notation is introduced, Section 2 presents a precise formulation for the adaptive confidence interval problem. Section 3 establishes the minimaxity and adaptivity results for a general linear functional  $\xi^\top \beta$  with a sparse loading  $\xi$ . Section 4 focuses on confidence intervals for a general linear functional  $\xi^\top \beta$  with a dense loading  $\xi$ . Section 5 considers the case when there is prior knowledge

of covariance matrix of the random design and the noise level of the linear model. Section 6 discusses connections to other work and further research directions. The proofs of the main results are given in Section 7. More discussion and proofs are presented in the supplement [3].

**2. Formulation for adaptive confidence interval problem.** We present in this section the framework for studying the adaptivity of confidence intervals. We begin with the notation that will be used throughout the paper.

*2.1. Notation.* For a matrix  $X \in \mathbb{R}^{n \times p}$ ,  $X_i$ ,  $X_j$ , and  $X_{i,j}$  denote respectively the  $i$ -th row,  $j$ -th column, and  $(i, j)$  entry of the matrix  $X$ ,  $X_{i,-j}$  denotes the  $i$ -th row of  $X$  excluding the  $j$ -th coordinate, and  $X_{-j}$  denotes the submatrix of  $X$  excluding the  $j$ -th column. Let  $[p] = \{1, 2, \dots, p\}$ . For a subset  $J \subset [p]$ ,  $X_J$  denotes the submatrix of  $X$  consisting of columns  $X_{\cdot j}$  with  $j \in J$  and for a vector  $x \in \mathbb{R}^p$ ,  $x_J$  is the subvector of  $x$  with indices in  $J$  and  $x_{-J}$  is the subvector with indices in  $J^c$ . For a set  $S$ ,  $|S|$  denotes the cardinality of  $S$ . For a vector  $x \in \mathbb{R}^p$ ,  $\text{supp}(x)$  denotes the support of  $x$  and the  $\ell_q$  norm of  $x$  is defined as  $\|x\|_q = (\sum_{i=1}^p |x_i|^q)^{\frac{1}{q}}$  for  $q \geq 1$  with  $\|x\|_0 = |\text{supp}(x)|$  and  $\|x\|_\infty = \max_{1 \leq j \leq p} |x_j|$ . We use  $e_i$  to denote the  $i$ -th standard basis vector in  $\mathbb{R}^p$ . For  $a \in \mathbb{R}$ ,  $a_+ = \max\{a, 0\}$ . We use  $\sum \beta_i$  as a shorthand for  $\sum_{i=1}^p \beta_i$ ,  $\max \|X_{\cdot j}\|_2$  as a shorthand for  $\max_{1 \leq j \leq p} \|X_{\cdot j}\|_2$  and  $\min \|X_{\cdot j}\|_2$  as a shorthand for  $\min_{1 \leq j \leq p} \|X_{\cdot j}\|_2$ . For a matrix  $A$  and  $1 \leq q \leq \infty$ ,  $\|A\|_q = \sup_{\|x\|_q=1} \|Ax\|_q$  is the matrix  $\ell_q$  operator norm. In particular,  $\|A\|_2$  is the spectral norm. For a symmetric matrix  $A$ ,  $\lambda_{\min}(A)$  and  $\lambda_{\max}(A)$  denote respectively the smallest and largest eigenvalue of  $A$ . We use  $c$  and  $C$  to denote generic positive constants that may vary from place to place. For two positive sequences  $a_n$  and  $b_n$ ,  $a_n \lesssim b_n$  means  $a_n \leq Cb_n$  for all  $n$  and  $a_n \gtrsim b_n$  if  $b_n \lesssim a_n$  and  $a_n \asymp b_n$  if  $a_n \lesssim b_n$  and  $b_n \lesssim a_n$ , and  $a_n \ll b_n$  if  $\limsup_{n \rightarrow \infty} \frac{a_n}{b_n} = 0$  and  $a_n \gg b_n$  if  $b_n \ll a_n$ .

*2.2. Framework for adaptivity of confidence intervals.* We shall focus in this paper on the high-dimensional linear model with the Gaussian design,

$$(2.1) \quad y_{n \times 1} = X_{n \times p} \beta_{p \times 1} + \epsilon_{n \times 1}, \quad \epsilon \sim N_n(0, \sigma^2 \mathbf{I}),$$

where the rows of  $X$  satisfy  $X_i \stackrel{\text{i.i.d.}}{\sim} N_p(0, \Sigma)$ ,  $i = 1, \dots, n$ , and are independent of  $\epsilon$ . Both  $\Sigma$  and the noise level  $\sigma$  are unknown. Let  $\Omega = \Sigma^{-1}$  denote the precision matrix. The parameter  $\theta = (\beta, \Omega, \sigma)$  consists of the signal  $\beta$ , the precision matrix  $\Omega$  for the random design, and the noise level  $\sigma$ . The target of interest is the linear functional of  $\beta$ ,  $T(\beta) = \xi^\top \beta$ , where  $\xi \in \mathbb{R}^p$  is a

pre-specified loading vector. The data that we observe is  $Z = (Z_1, \dots, Z_n)^\top$ , where  $Z_i = (y_i, X_i) \in \mathbb{R}^{p+1}$  for  $i = 1, \dots, n$ .

For  $0 < \alpha < 1$  and a given parameter space  $\Theta$  and the linear functional  $T(\beta)$ , denote by  $\mathcal{I}_\alpha(\Theta, T)$  the set of all  $(1 - \alpha)$  level confidence intervals for  $T(\beta)$  over the parameter space  $\Theta$ ,

$$(2.2) \quad \mathcal{I}_\alpha(\Theta, T) = \left\{ \text{CI}_\alpha(T, Z) = [l(Z), u(Z)] : \inf_{\theta \in \Theta} \mathbb{P}_\theta(l(Z) \leq T(\beta) \leq u(Z)) \geq 1 - \alpha \right\}.$$

For any confidence interval  $\text{CI}_\alpha(T, Z) \in \mathcal{I}_\alpha(\Theta, T)$ , the maximum expected length over a parameter space  $\Theta$  is defined as

$$L(\text{CI}_\alpha(T, Z), \Theta, T) = \sup_{\theta \in \Theta} \mathbb{E}_\theta L(\text{CI}_\alpha(T, Z)),$$

where for confidence interval  $\text{CI}_\alpha(T, Z) = [l(Z), u(Z)]$ ,  $L(\text{CI}_\alpha(T, Z)) = u(Z) - l(Z)$  denotes its length. For two parameter spaces  $\Theta_1 \subseteq \Theta$ , we define the benchmark  $L_\alpha^*(\Theta_1, \Theta, T)$  as the infimum of the maximum expected length over  $\Theta_1$  among all  $(1 - \alpha)$ -level confidence intervals over  $\Theta$ ,

$$(2.3) \quad L_\alpha^*(\Theta_1, \Theta, T) = \inf_{\text{CI}_\alpha(T, Z) \in \mathcal{I}_\alpha(\Theta, T)} L(\text{CI}_\alpha(T, Z), \Theta_1, T).$$

We will write  $L_\alpha^*(\Theta, T)$  for  $L_\alpha^*(\Theta, \Theta, T)$ , which is the minimax expected length of confidence intervals over  $\Theta$ .

We should emphasize that  $L_\alpha^*(\Theta_1, \Theta, T)$  is an important quantity that measures the degree of adaptivity over the nested spaces  $\Theta_1 \subset \Theta$ . A confidence interval  $\text{CI}_\alpha(T, Z)$  that is (rate-optimally) adaptive over  $\Theta_1$  and  $\Theta$  should have the optimal expected length performance simultaneously over both  $\Theta_1$  and  $\Theta$  while maintaining a given coverage probability over  $\Theta$ , i.e.,  $\text{CI}_\alpha(T, Z) \in \mathcal{I}_\alpha(\Theta, T)$  such that

$$L(\text{CI}_\alpha(T, Z), \Theta_1, T) \asymp L_\alpha^*(\Theta_1, T) \quad \text{and} \quad L(\text{CI}_\alpha(T, Z), \Theta, T) \asymp L_\alpha^*(\Theta, T).$$

Note that in this case  $L(\text{CI}_\alpha(T, Z), \Theta_1, T) \geq L_\alpha^*(\Theta_1, \Theta, T)$ . So for two parameter spaces  $\Theta_1 \subset \Theta$ , if  $L_\alpha^*(\Theta_1, \Theta, T) \gg L_\alpha^*(\Theta_1, T)$ , then rate-optimal adaptation between  $\Theta_1$  and  $\Theta$  is impossible to achieve.

We consider the following collection of parameter spaces,

$$(2.4) \quad \Theta(k) = \left\{ \theta = (\beta, \Omega, \sigma) : \|\beta\|_0 \leq k, \frac{1}{M_1} \leq \lambda_{\min}(\Omega) \leq \lambda_{\max}(\Omega) \leq M_1, 0 < \sigma \leq M_2 \right\},$$

where  $M_1 > 1$  and  $M_2 > 0$  are positive constants. Basically,  $\Theta(k)$  is the set of all  $k$ -sparse regression vectors.  $\frac{1}{M_1} \leq \lambda_{\min}(\Omega) \leq \lambda_{\max}(\Omega) \leq M_1$  and  $0 < \sigma \leq M_2$  are two mild regularity conditions on the design and noise level.

The main goal of this paper is to address the following two questions:



1. *What is the minimax length  $L_\alpha^*(\Theta(k), \mathbb{T})$  in the oracle setting where the sparsity level  $k$  is known?*
2. *Is it possible to achieve rate-optimal adaptation over different sparsity levels?*

More specifically, for  $k_1 \ll k$ , is it possible to construct a confidence interval  $\text{CI}_\alpha(\mathbb{T}, Z)$  that is adaptive over  $\Theta(k_1)$  and  $\Theta(k)$  in the sense that  $\text{CI}_\alpha(\mathbb{T}, Z) \in \mathcal{I}_\alpha(\Theta(k), \mathbb{T})$  and

$$(2.5) \quad \begin{aligned} L(\text{CI}_\alpha(\mathbb{T}, Z), \Theta(k_1), \mathbb{T}) &\asymp L_\alpha^*(\Theta(k_1), \mathbb{T}), \\ L(\text{CI}_\alpha(\mathbb{T}, Z), \Theta(k), \mathbb{T}) &\asymp L_\alpha^*(\Theta(k), \mathbb{T})? \end{aligned}$$

We will answer these questions by analyzing the two benchmark quantities  $L_\alpha^*(\Theta(k), \mathbb{T})$  and  $L_\alpha^*(\Theta(k_1), \Theta(k), \mathbb{T})$ . Both lower and upper bounds will be established. If (2.5) can be achieved, it means that the confidence interval  $\text{CI}_\alpha(\mathbb{T}, Z)$  can automatically adjust its length to the sparsity level of the true regression vector  $\beta$ . On the other hand, if  $L_\alpha^*(\Theta(k_1), \Theta(k), \mathbb{T}) \gg L_\alpha^*(\Theta(k_1), \mathbb{T})$ , then such a goal is not attainable.

For ease of presentation, we calibrate the sparsity level

$$k \asymp p^\gamma \quad \text{for some } 0 \leq \gamma < \frac{1}{2},$$

and restrict the loading  $\xi$  to the set

$$\xi \in \Xi(q, \bar{c}) = \left\{ \xi \in \mathbb{R}^p : \|\xi\|_0 = q, \xi \neq \mathbf{0} \text{ and } \frac{\max_{j \in \text{supp}(\xi)} |\xi_j|}{\min_{j \in \text{supp}(\xi)} |\xi_j|} \leq \bar{c} \right\},$$

where  $\bar{c} \geq 1$  is a constant. The minimax rate and adaptivity of confidence intervals for the general linear functional  $\xi^\top \beta$  also depends on the sparsity of  $\xi$ . We are particularly interested in the following two regimes:

1. The sparse loading regime:  $\xi \in \Xi(q, \bar{c})$  with

$$(2.6) \quad q \leq Ck.$$

2. The dense loading regime:  $\xi \in \Xi(q, \bar{c})$  with

$$(2.7) \quad q = cp^{\gamma_q} \quad \text{with } 2\gamma < \gamma_q \leq 1.$$

The behavior of the problem is significantly different in these two regimes. We will consider separately the sparse loading regime in Section 3 and the dense loading regime in Section 4.

**3. Minimax rate and adaptivity of confidence intervals for sparse loading linear functionals.** In this section, we establish the rates of convergence for the minimax expected length of confidence intervals for  $\xi^\top \beta$  with a sparse loading  $\xi$  in the oracle setting where the sparsity parameter  $k$  of the regression vector  $\beta$  is given. Both minimax upper and lower bounds are given. Confidence intervals for  $\xi^\top \beta$  are constructed and shown to be minimax rate-optimal in the sparse loading regime. Finally, we establish the possibility of adaptivity for the linear functional  $\xi^\top \beta$  with a sparse loading  $\xi$ .

3.1. *Minimax length of confidence intervals for  $\xi^\top \beta$  in the sparse loading regime.* In this section, we focus on the sparse loading regime defined in (2.6). The following theorem establishes the minimax rates for the expected length of confidence intervals for  $\xi^\top \beta$  in the sparse loading regime.

**THEOREM 1.** *Suppose that  $0 < \alpha < \frac{1}{2}$  and  $k \leq c \min\{p^\gamma, \frac{n}{\log p}\}$  for some constants  $c > 0$  and  $0 \leq \gamma < \frac{1}{2}$ . If  $\xi$  belongs to the sparse loading regime (2.6), the minimax expected length for  $(1 - \alpha)$  level confidence intervals of  $\xi^\top \beta$  over  $\Theta(k)$  satisfies*

$$(3.1) \quad L_\alpha^*(\Theta(k), \xi^\top \beta) \asymp \|\xi\|_2 \left( \frac{1}{\sqrt{n}} + k \frac{\log p}{n} \right).$$

Theorem 1 is established in two separate steps.

1. Minimax upper bound: we construct a confidence interval  $\text{CI}_\alpha^S(\xi^\top \beta, Z)$  such that  $\text{CI}_\alpha^S(\xi^\top \beta, Z) \in \mathcal{I}_\alpha(\Theta(k), \xi^\top \beta)$  and for some constant  $C > 0$

$$(3.2) \quad L(\text{CI}_\alpha^S(\xi^\top \beta, Z), \Theta(k), \xi^\top \beta) \leq C \|\xi\|_2 \left( \frac{1}{\sqrt{n}} + k \frac{\log p}{n} \right).$$

2. Minimax lower bound: we show that for some constant  $c > 0$

$$(3.3) \quad L_\alpha^*(\Theta(k), \xi^\top \beta) \geq c \|\xi\|_2 \left( \frac{1}{\sqrt{n}} + k \frac{\log p}{n} \right).$$

The minimax lower bound is implied by the adaptivity result given in Theorem 2. We now detail the construction of a confidence interval  $\text{CI}_\alpha^S(\xi^\top \beta, Z)$  achieving the minimax rate (3.1) in the sparse loading regime. The interval  $\text{CI}_\alpha^S(\xi^\top \beta, Z)$  is centered at a de-biased scaled Lasso estimator, which generalizes the ideas used in [25, 15, 23]. The construction of the (random) length is different from the aforementioned papers as the asymptotic normality result is not valid once  $k \gtrsim \frac{\sqrt{n}}{\log p}$ .

Let  $\{\widehat{\beta}, \widehat{\sigma}\}$  be the scaled Lasso estimator with  $\lambda_0 = \sqrt{\frac{2.05 \log p}{n}}$ ,

$$(3.4) \quad \{\widehat{\beta}, \widehat{\sigma}\} = \arg \min_{\beta \in \mathbb{R}^p, \sigma \in \mathbb{R}^+} \frac{\|y - X\beta\|_2^2}{2n\sigma} + \frac{\sigma}{2} + \lambda_0 \sum_{j=1}^p \frac{\|X_{\cdot j}\|_2}{\sqrt{n}} |\beta_j|.$$

Define

$$(3.5) \quad \widehat{u} = \arg \min_{u \in \mathbb{R}^p} \left\{ u^\top \widehat{\Sigma} u : \|\widehat{\Sigma} u - \xi\|_\infty \leq \lambda_n \right\},$$

where  $\widehat{\Sigma} = \frac{1}{n} X^\top X$  and  $\lambda_n = 12 \|\xi\|_2 M_1^2 \sqrt{\frac{\log p}{n}}$ . The confidence interval  $\text{CI}_\alpha^S(\xi^\top \beta, Z)$  is centered at the following de-biased estimator

$$(3.6) \quad \widetilde{\mu} = \xi^\top \widehat{\beta} + \widehat{u}^\top \frac{1}{n} X^\top (y - X\widehat{\beta}),$$

where  $\widehat{\beta}$  is the scaled Lasso estimator given in (3.4) and  $\widehat{u}$  is defined in (3.5). Before specifying the length of the confidence interval, we review the following definition of restricted eigenvalue introduced in [2],

$$(3.7) \quad \kappa(X, k, \alpha_0) = \min_{\substack{J_0 \subset \{1, \dots, p\}, \\ |J_0| \leq k}} \min_{\substack{\delta \neq 0, \\ \|\delta_{J_0^c}\|_1 \leq \alpha_0 \|\delta_{J_0}\|_1}} \frac{\|X\delta\|_2}{\sqrt{n} \|\delta_{J_0}\|_2}.$$

Define

$$(3.8) \quad \rho_1(k) = \|\xi\|_2 \widehat{\sigma} \min \left\{ 1.01 \sqrt{\frac{\widehat{u}^\top \widehat{\Sigma} \widehat{u}}{n \|\xi\|_2^2}} z_{\alpha/2} + C_1(X, k) k \frac{\log p}{n}, \log p \left( \frac{1}{\sqrt{n}} + \frac{k \log p}{n} \right) \right\},$$

where  $z_{\alpha/2}$  is the  $\alpha/2$  upper quantile of the standard normal distribution and

$$(3.9) \quad C_1(X, k) = 7000 M_1^2 \frac{\sqrt{n}}{\min \|X_{\cdot j}\|_2} \max \left\{ 1.25, \frac{912 \max \|X_{\cdot j}\|_2^2}{n \kappa^2 \left( X, k, 405 \left( \frac{\max \|X_{\cdot j}\|_2}{\min \|X_{\cdot j}\|_2} \right) \right)} \right\}.$$

Define the event

$$(3.10) \quad A = \{\widehat{\sigma} \leq \log p\}.$$

The confidence interval  $\text{CI}_\alpha^S(\xi^\top \beta, Z)$  for  $\xi^\top \beta$  is defined as

$$(3.11) \quad \text{CI}_\alpha^S(\xi^\top \beta, Z) = \begin{cases} [\widetilde{\mu} - \rho_1(k), \widetilde{\mu} + \rho_1(k)] & \text{on } A \\ \{0\} & \text{on } A^c \end{cases}$$

It will be shown in Section 7 that the confidence interval  $\text{CI}_\alpha^S(\xi^\top \beta, Z)$  has the desired coverage property and achieves the minimax length in (3.1).

REMARK 1. In the special case of  $\xi = e_1$ , the confidence interval defined in (3.11) is similar to the ones based on the de-biased estimators introduced in [25, 15, 23]. The second term  $\widehat{u}^\top \frac{1}{n} X^\top (y - X\widehat{\beta})$  in (3.6) is incorporated to reduce the bias of the scaled Lasso estimator  $\widehat{\beta}$ . The constrained estimator  $\widehat{u}$  defined in (3.5) is a score vector  $u$  such that the variance term  $u^\top \widehat{\Sigma} u$  is minimized and one component of the bias term  $\|\widehat{\Sigma} u - \xi\|_\infty$  is constrained by the tuning parameter  $\lambda_n$ . The tuning parameter  $\lambda_n$  is chosen as  $12\|\xi\|_2 M_1^2 \sqrt{\frac{\log p}{n}}$  such that  $u = \Omega\xi$  lies in the constraint set  $\|\widehat{\Sigma} u - \xi\|_\infty \leq \lambda_n$  in (3.5) with overwhelming probability. For  $C_1(X, k)$  defined in (3.9), it will be shown that it is upper bounded by a constant with overwhelming probability.

3.2. *Adaptivity of confidence intervals for  $\xi^\top \beta$  in the sparse loading regime.* We have constructed a minimax rate-optimal confidence interval for  $\xi^\top \beta$  in the oracle setting where the sparsity  $k$  is assumed to be known. A major drawback of the construction is that it requires prior knowledge of  $k$ , which is typically unavailable in practice. An interesting question is whether it is possible to construct adaptive confidence intervals that have the guaranteed coverage and automatically adjust its length to  $k$ .

We now consider the adaptivity of the confidence intervals for  $\xi^\top \beta$ . In light of the minimax expected length given in Theorem 1, the following theorem provides an answer to the adaptivity question (2.5) for the confidence intervals for  $\xi^\top \beta$  in the sparse loading regime.

THEOREM 2. *Suppose that  $0 < \alpha < \frac{1}{2}$  and  $k_1 \leq k \leq c \min \left\{ p^\gamma, \frac{n}{\log p} \right\}$  for some constants  $c > 0$  and  $0 \leq \gamma < \frac{1}{2}$ . Then*

$$(3.12) \quad L_\alpha^*(\Theta(k_1), \Theta(k), \xi^\top \beta) \geq c_1 \|\xi\|_2 \left( \frac{1}{\sqrt{n}} + k \frac{\log p}{n} \right),$$

for some constant  $c_1 > 0$ .

Note that Theorem 2 implies the minimax lower bound in Theorem 1 by taking  $k_1 = k$ . Theorem 2 rules out the possibility of rate-optimal adaptive confidence intervals beyond the ultra-sparse region. Consider the setting where  $k_1 \ll k$  and  $\frac{\sqrt{n}}{\log p} \ll k \lesssim \frac{n}{\log p}$ . In this case,

$$L_\alpha^*(\Theta(k_1), \Theta(k), \xi^\top \beta) \asymp L_\alpha^*(\Theta(k), \xi^\top \beta) \asymp \|\xi\|_2 k \frac{\log p}{n} \gg L_\alpha^*(\Theta(k_1), \xi^\top \beta).$$

So it is impossible to construct a confidence interval that is adaptive simultaneously over  $\Theta(k_1)$  and  $\Theta(k)$  when  $\frac{\sqrt{n}}{\log p} \ll k \lesssim \frac{n}{\log p}$  and  $k_1 \ll k$ . The only

possible region for adaptation is in the ultra-sparse region  $k \lesssim \frac{\sqrt{n}}{\log p}$ , over which the optimal expected length of confidence intervals is of order  $\frac{1}{\sqrt{n}}$  and in particular does not depend on the specific sparsity level. These facts are illustrated in Figure 1.

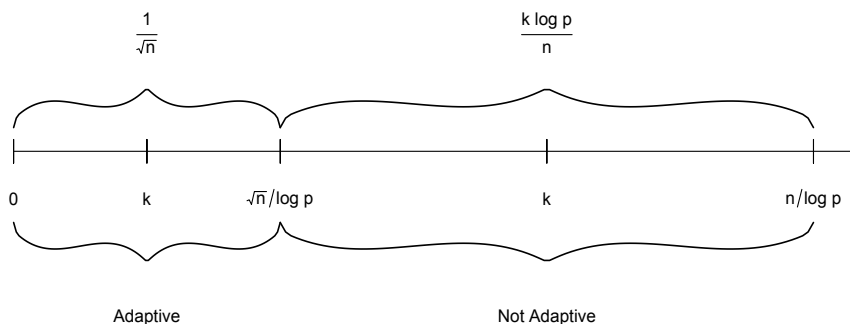


FIG 1. Illustration of adaptivity of confidence intervals for  $\xi^\top \beta$  with a sparse loading  $\xi$ . For adaptation between  $\Theta(k_1)$  and  $\Theta(k)$  with  $k_1 \ll k$ , rate-optimal adaptation is possible if  $k \lesssim \frac{\sqrt{n}}{\log p}$  and impossible otherwise.

So far the analysis is carried out within the minimax framework where the focus is on the performance in the worst case over a large parameter space. The minimax theory is often criticized as being too conservative. In the following, we establish a stronger version of the non-adaptivity result which demonstrates that the lack of adaptivity for confidence intervals is not due to the conservativeness of the minimax framework. The result shows that for any confidence interval  $\text{CI}_\alpha(\xi^\top \beta, Z)$ , under the coverage constraint that  $\text{CI}_\alpha(\xi^\top \beta, Z) \in \mathcal{I}_\alpha(\Theta(k), \xi^\top \beta)$ , its expected length at any given  $\theta^* = (\beta^*, \mathbf{I}, \sigma) \in \Theta(k_1)$  must be of order  $\|\xi\|_2 \left( \frac{1}{\sqrt{n}} + k \frac{\log p}{n} \right)$ . So the confidence interval must be long at a large subset of points in the parameter space, not just at a small number of “unlucky” points.

**THEOREM 3.** *Suppose that  $0 < \alpha < \frac{1}{2}$  and  $k \leq c \min\{p^\gamma, \frac{n}{\log p}\}$  for some constants  $c > 0$  and  $0 \leq \gamma < \frac{1}{2}$ . Let  $k_1 \leq (1 - \zeta_0)k - 1$  and  $q \leq \frac{\zeta_0}{4}k$  for some constant  $0 < \zeta_0 < 1$ . Then for any  $\theta^* = (\beta^*, \mathbf{I}, \sigma) \in \Theta(k_1)$  and  $\xi \in \Xi(q, \bar{c})$ ,*

$$(3.13) \quad \inf_{\text{CI}_\alpha(\xi^\top \beta, Z) \in \mathcal{I}_\alpha(\Theta(k), \xi^\top \beta)} \mathbb{E}_{\theta^*} L(\text{CI}_\alpha(\xi^\top \beta, Z)) \geq c_1 \|\xi\|_2 \left( k \frac{\log p}{n} + \frac{1}{\sqrt{n}} \right) \sigma,$$

for some constant  $c_1 > 0$ .

Note that no supremum is taken over the parameter  $\theta^*$  in (3.13). Theorem

3 illustrates that if a confidence interval  $\text{CI}_\alpha(\xi^\top \beta, Z)$  is “superefficient” at any point  $\theta^* = (\beta^*, \mathbf{I}, \sigma) \in \Theta(k_1)$  in the sense that

$$\mathbb{E}_{\theta^*} L(\text{CI}_\alpha(\xi^\top \beta, Z)) \ll \|\xi\|_2 \left( \frac{1}{\sqrt{n}} + k \frac{\log p}{n} \right) \sigma,$$

then the confidence interval  $\text{CI}_\alpha(\xi^\top \beta, Z)$  can not have the guaranteed coverage over the parameter space  $\Theta(k)$ .

3.3. *Minimax rate and adaptivity of confidence intervals for  $\beta_1$ .* We now turn to the special case  $\mathbb{T}(\beta) = \beta_i$ , which has been the focus of several previous papers [25, 14, 15, 23]. Without loss of generality, we consider  $\beta_1$ , the first coordinate of  $\beta$ , in the following discussion and the results for any other coordinate  $\beta_i$  are the same. The linear functional  $\beta_1$  is the special case of linear functional of sparse loading regime with  $\xi = e_1$ .

Theorem 1 implies that the minimax expected length for  $(1 - \alpha)$  level confidence intervals of  $\beta_1$  over  $\Theta(k)$  satisfies

$$(3.14) \quad L_\alpha^*(\Theta(k), \beta_1) \asymp \frac{1}{\sqrt{n}} + k \frac{\log p}{n}.$$

In the ultra-sparse region with  $k \lesssim \frac{\sqrt{n}}{\log p}$ , the minimax expected length is of order  $\frac{1}{\sqrt{n}}$ . However, when  $k$  falls in the moderate-sparse region  $\frac{\sqrt{n}}{\log p} \ll k \lesssim \frac{n}{\log p}$ , the minimax expected length is of order  $k \frac{\log p}{n}$  and in this case  $k \frac{\log p}{n} \gg \frac{1}{\sqrt{n}}$ . Hence the confidence intervals constructed in [25, 14, 15, 23], which are of parametric length  $\frac{1}{\sqrt{n}}$ , asymptotically have coverage probability going to 0. The condition  $k \lesssim \frac{\sqrt{n}}{\log p}$  is necessary for the parametric rate  $\frac{1}{\sqrt{n}}$ . [23] established asymptotic normality and asymptotic efficiency for a de-biased estimator under the sparsity assumption  $k \ll \frac{\sqrt{n}}{\log p}$ . Similar results have also been given in [19] for a related problem of estimating a single entry of a  $p$ -dimensional precision matrix based on  $n$  i.i.d. samples under the same sparsity condition  $k \ll \frac{\sqrt{n}}{\log p}$ . It was also shown that  $k \ll \frac{\sqrt{n}}{\log p}$  is necessary for the asymptotic normality and asymptotic efficiency results.

The following corollary, as a special case of Theorem 3, illustrates the strong non-adaptivity for confidence intervals of  $\beta_1$  when  $k \gg \frac{\sqrt{n}}{\log p}$ .

COROLLARY 1. *Suppose that  $0 < \alpha < \frac{1}{2}$  and  $k \leq c \min\{p^\gamma, \frac{n}{\log p}\}$  for some constants  $c > 0$  and  $0 \leq \gamma < \frac{1}{2}$ . Let  $k_1 \leq (1 - \zeta_0)k - 1$  for some constant  $0 < \zeta_0 < 1$ . Then for any  $\theta^* = (\beta^*, \mathbf{I}, \sigma) \in \Theta(k_1)$ ,*

$$(3.15) \quad \inf_{\text{CI}_\alpha(\beta_1, Z) \in \mathcal{I}_\alpha(\Theta(k), \beta_1)} \mathbb{E}_{\theta^*} L(\text{CI}_\alpha(\beta_1, Z)) \geq c_1 \left( \frac{1}{\sqrt{n}} + k \frac{\log p}{n} \right) \sigma,$$

for some constant  $c_1 > 0$ .

**4. Minimax rate and adaptivity of confidence intervals for dense loading linear functionals.** We now turn to the setting where the loading  $\xi$  is dense in the sense of (2.7). We will also briefly discuss the special case  $\sum_{i=1}^p \beta_i$  and the computationally feasible confidence intervals.

4.1. *Minimax length of confidence intervals for  $\xi^\top \beta$  in the dense loading regime.* The following theorem establishes the minimax length of confidence intervals of  $\xi^\top \beta$  in the dense loading regime (2.7).

**THEOREM 4.** *Suppose that  $0 < \alpha < \frac{1}{2}$  and  $k \leq c \min\{p^\gamma, \frac{n}{\log p}\}$  for some constants  $c > 0$  and  $0 \leq \gamma < \frac{1}{2}$ . If  $\xi$  belongs to the dense loading regime (2.7), the minimax expected length for  $(1 - \alpha)$  level confidence intervals of  $\xi^\top \beta$  over  $\Theta(k)$  satisfies*

$$(4.1) \quad L_\alpha^*(\Theta(k), \xi^\top \beta) \asymp \|\xi\|_\infty k \sqrt{\frac{\log p}{n}}.$$

Note that the minimax rate in (4.1) is significantly different from the minimax rate  $\|\xi\|_2 (\frac{1}{\sqrt{n}} + k \frac{\log p}{n})$  for the sparse loading case given in Theorem 1. In the following, we construct a confidence interval  $\text{CI}_\alpha^D(\xi^\top \beta, Z)$  achieving the minimax rate (4.1) in the dense loading regime. Define

$$(4.2) \quad C_2(X, k) = 822 \frac{\sqrt{n}}{\min \|X_{\cdot j}\|_2} \max \left\{ 1.25, \frac{912 \max \|X_{\cdot j}\|_2^2}{n \kappa^2 \left( X, k, 405 \left( \frac{\max \|X_{\cdot j}\|_2}{\min \|X_{\cdot j}\|_2} \right) \right)} \right\}.$$

It will be shown that  $C_2(X, k)$  is upper bounded by a constant with overwhelming probability. The confidence interval  $\text{CI}_\alpha^D(\xi^\top \beta, Z)$  is defined to be,

$$(4.3) \quad \text{CI}_\alpha^D(\xi^\top \beta, Z) = \begin{cases} \left[ \xi^\top \hat{\beta} - \|\xi\|_\infty \rho_2(k), \xi^\top \hat{\beta} + \|\xi\|_\infty \rho_2(k) \right] & \text{on } A \\ \{0\} & \text{on } A^c \end{cases}$$

where  $A$  is defined in (3.10) and  $\hat{\beta}$  is the scaled Lasso estimator defined in (3.4) and

$$(4.4) \quad \rho_2(k) = \min \left\{ C_2(X, k) k \sqrt{\frac{\log p}{n}} \hat{\sigma}, \log p \left( k \sqrt{\frac{\log p}{n}} \hat{\sigma} \right) \right\}.$$

The confidence interval constructed in (4.3) will be shown to have the desired coverage property and achieve the minimax length in (4.1). A major

difference between the construction of  $\text{CI}_\alpha^D(\xi^\top\beta, Z)$  and that of  $\text{CI}_\alpha^S(\xi^\top\beta, Z)$  is that  $\text{CI}_\alpha^D(\xi^\top\beta, Z)$  is not centered at a de-biased estimator. If a de-biased estimator is used for the construction of confidence intervals for  $\xi^\top\beta$  with a dense loading, its variance would be too large, much larger than the optimal length  $\|\xi\|_\infty k \sqrt{\frac{\log p}{n}}$ .

4.2. *Adaptivity of confidence intervals for  $\xi^\top\beta$  in the dense loading regime.* In this section, we investigate the possibility of adaptive confidence intervals for  $\xi^\top\beta$  in the dense loading regime. The following theorem leads directly to an answer to the adaptivity question (2.5) for confidence intervals for  $\xi^\top\beta$  in the dense loading regime.

**THEOREM 5.** *Suppose that  $0 < \alpha < \frac{1}{2}$  and  $k_1 \leq k \leq c \min\left\{p^\gamma, \frac{n}{\log p}\right\}$  for some constants  $c > 0$  and  $0 \leq \gamma < \frac{1}{2}$ . Then, for some constant  $c_1 > 0$ ,*

$$(4.5) \quad L_\alpha^*(\Theta(k_1), \Theta(k), \xi^\top\beta) \geq c_1 \|\xi\|_\infty k \sqrt{\frac{\log p}{n}}.$$

Theorem 5 implies the minimax lower bound in Theorem 4 by taking  $k_1 = k$ . If  $k_1 \ll k$ , (4.5) implies

$$(4.6) \quad L_\alpha^*(\Theta(k_1), \Theta(k), \xi^\top\beta) \geq c \|\xi\|_\infty k \sqrt{\frac{\log p}{n}} \gg L_\alpha^*(\Theta(k_1), \xi^\top\beta),$$

which shows that rate-optimal adaptation over two different sparsity levels  $k_1$  and  $k$  is not possible at all for any  $k_1 \ll k$ . In contrast, in the case of the sparse loading regime, Theorem 2 shows that it is possible to construct an adaptive confidence interval in the ultra-sparse region  $k \lesssim \frac{\sqrt{n}}{\log p}$ , although adaptation is not possible in the moderate-sparse region  $\frac{\sqrt{n}}{\log p} \ll k \lesssim \frac{n}{\log p}$ .

Similarly to Theorem 3, the following theorem establishes the strong non-adaptivity results for  $\xi^\top\beta$  in the dense loading regime.

**THEOREM 6.** *Suppose that  $0 < \alpha < \frac{1}{2}$  and  $k \leq c \min\{p^\gamma, \frac{n}{\log p}\}$  for some constants  $c > 0$  and  $0 \leq \gamma < \frac{1}{2}$ . Let  $q$  satisfies (2.7) and  $k_1 \leq (1 - \zeta_0)k - 1$  for some positive constant  $0 < \zeta_0 < 1$ . Then for any  $\theta^* = (\beta^*, \mathbf{I}, \sigma) \in \Theta(k_1)$  and  $\xi \in \Xi(q, \bar{c})$ , there is some constant  $c_1 > 0$  such that*

$$(4.7) \quad \inf_{\text{CI}_\alpha(\xi^\top\beta, Z) \in \mathcal{I}_\alpha(\Theta(k), \xi^\top\beta)} \mathbb{E}_{\theta^*} L(\text{CI}_\alpha(\xi^\top\beta, Z)) \geq c_1 \|\xi\|_\infty k \sqrt{\frac{\log p}{n}} \sigma.$$



4.3. *Minimax length and adaptivity of confidence intervals for  $\sum_{i=1}^p \beta_i$ .* We now turn to the special case of  $T(\beta) = \sum_{i=1}^p \beta_i$ , the sum of all coefficients. Theorem 4 implies that the minimax expected length for  $(1 - \alpha)$  level confidence intervals of  $\sum_{i=1}^p \beta_i$  over  $\Theta(k)$  satisfies

$$(4.8) \quad L_\alpha^* \left( \Theta(k), \sum \beta_i \right) \asymp k \sqrt{\frac{\log p}{n}}.$$

The following impossibility of adaptivity result for confidence intervals for  $\sum_{i=1}^p \beta_i$  is a special case of Theorem 6.

COROLLARY 2. *Suppose that  $0 < \alpha < \frac{1}{2}$  and  $k \leq c \min\{p^\gamma, \frac{n}{\log p}\}$  for some constants  $c > 0$  and  $0 \leq \gamma < \frac{1}{2}$ . Let  $k_1 \leq (1 - \zeta_0)k - 1$  for some constant  $0 < \zeta_0 < 1$ . Then for any  $\theta^* = (\beta^*, \mathbf{I}, \sigma) \in \Theta(k_1)$ ,*

$$(4.9) \quad \inf_{\text{CI}_\alpha(\sum \beta_i, Z) \in \mathcal{I}_\alpha(\Theta(k), \sum \beta_i)} \mathbb{E}_{\theta^*} L \left( \text{CI}_\alpha \left( \sum \beta_i, Z \right) \right) \geq c_1 k \sqrt{\frac{\log p}{n}} \sigma,$$

for some constant  $c_1 > 0$ .

REMARK 2. In the Gaussian sequence model, the problem of estimating the sum of sparse means has been considered in [5, 7] and more recently in [12]. In particular, minimax rate is given in [5] and [12]. The problem of constructing minimax confidence intervals for the sum of sparse normal means was studied in [6].

4.4. *Computationally feasible confidence intervals.* A major drawback of the minimax rate-optimal confidence intervals  $\text{CI}_\alpha^S(\xi^\top \beta, Z)$  given in (3.11) and  $\text{CI}_\alpha^D(\xi^\top \beta, Z)$  given in (4.3) is that they are not computationally feasible as both depend on restricted eigenvalue  $\kappa(X, k, \alpha_0)$ , which is difficult to evaluate. In this section, we assume the prior knowledge of the sparsity  $k$  and discuss how to construct a computationally feasible confidence interval.

The main idea is to replace the term involved with restricted eigenvalue by a computationally feasible lower bound function  $\omega(\Omega, X, k)$  defined by

$$(4.10) \quad \omega(\Omega, X, k) = \left( \frac{1}{4\sqrt{\lambda_{\max}(\Omega)}} - \frac{9 \left( 1 + 405 \frac{\max \|X_{\cdot j}\|_2}{\min \|X_{\cdot j}\|_2} \right)}{\sqrt{\lambda_{\min}(\Omega)}} \sqrt{k \frac{\log p}{n}} \right)_+^2.$$

The lower bound relation is established by Lemma 13 in the supplement [3], which is based on the concentration inequality for Gaussian design in [18]. Except for  $\lambda_{\min}(\Omega)$  and  $\lambda_{\max}(\Omega)$ , all terms in (4.10) are based on the data

$(X, y)$  and the prior knowledge of  $k$ . To construct a data-dependent computationally feasible confidence interval, we make the following assumption,

$$(4.11) \quad \sup_{\Omega \in \mathcal{G}_\Omega} \mathbb{P}_X \left( \max \left\{ \left| \widetilde{\lambda_{\min}}(\Omega) - \lambda_{\min}(\Omega) \right|, \left| \widetilde{\lambda_{\max}}(\Omega) - \lambda_{\max}(\Omega) \right| \right\} \geq C a_{n,p} \right) = o(1),$$

where  $\limsup a_{n,p} = 0$  and  $\mathcal{G}_\Omega$  is a pre-specified parameter space for  $\Omega$  and  $\mathbb{P}_X$  denotes the probability distribution with respect to  $X$ .

REMARK 3. We assume  $\mathcal{G}_\Omega$  is a subspace of the precision matrix defined in (2.4),  $\left\{ \Omega : \frac{1}{M_1} \leq \lambda_{\min}(\Omega) \leq \lambda_{\max}(\Omega) \leq M_1 \right\}$ . By assuming  $\mathcal{G}_\Omega$  is the set of precision matrix of special structure, we can find estimators satisfying (4.11). If  $\mathcal{G}_\Omega$  is assumed to be the set of sparse precision matrix, we can estimate the precision matrix  $\Omega$  by CLIME estimator  $\tilde{\Omega}$  proposed in [4]. Under proper sparsity assumption on  $\Omega$ , the plugin estimator  $(\lambda_{\min}(\tilde{\Omega}), \lambda_{\max}(\tilde{\Omega}))$  satisfies (4.11). Other special structures can also be assumed, for example, the covariance matrix is sparse. We can use the plugin estimator of the estimator proposed in [10].

With  $\widetilde{\lambda_{\min}}(\Omega)$  and  $\widetilde{\lambda_{\max}}(\Omega)$ , we define  $\tilde{\omega}(\Omega, X, k)$  as

$$\tilde{\omega}(\Omega, X, k) = \left( \frac{1}{4\sqrt{\widetilde{\lambda_{\max}}(\Omega)}} - \frac{9 \left( 1 + 405 \frac{\max \|X_{\cdot j}\|_2}{\min \|X_{\cdot j}\|_2} \right)}{\sqrt{\widetilde{\lambda_{\min}}(\Omega)}} \sqrt{k \frac{\log p}{n}} \right)_+^2.$$

and construct computationally feasible confidence intervals by replacing  $\kappa^2 \left( X, k, 405 \left( \frac{\max \|X_{\cdot j}\|_2}{\min \|X_{\cdot j}\|_2} \right) \right)$  in (3.11) and (4.3) with  $\tilde{\omega}(\Omega, X, k)$ .

**5. Confidence intervals for linear functionals with prior knowledge  $\Omega = \mathbf{I}$  and  $\sigma = \sigma_0$ .** We have so far focused on the setting where both the precision matrix  $\Omega$  and the noise level  $\sigma$  are unknown, which is the case in most statistical applications. It is still of theoretical interest to study the problem when  $\Omega$  and  $\sigma$  are known. It is interesting to contrast the results with the ones when  $\Omega$  and  $\sigma$  are unknown. In this case, we consider the setting where it is known a priori that  $\Omega = \mathbf{I}$  and  $\sigma = \sigma_0$  and specify the parameter space as

$$(5.1) \quad \Theta(k, \mathbf{I}, \sigma_0) = \{ \theta = (\beta, \mathbf{I}, \sigma_0) : \|\beta\|_0 \leq k \}.$$

We will discuss separately the minimax rates and adaptivity of confidence intervals for the linear functionals in the sparse loading regime and dense loading regime over the parameter space  $\Theta(k, \mathbf{I}, \sigma_0)$ .

5.1. *Confidence intervals for linear functionals in the sparse loading regime.* The following theorem establishes the minimax rate of confidence intervals for linear functionals in the sparse loading regime when there is prior knowledge that  $\Omega = \mathbf{I}$  and  $\sigma = \sigma_0$ .

**THEOREM 7.** *Suppose that  $0 < \alpha < \frac{1}{2}$  and  $k \leq c \min\{p^\gamma, \frac{n}{\log p}\}$  for some constants  $c > 0$  and  $0 \leq \gamma < \frac{1}{2}$ . If  $\xi$  belongs to the sparse loading regime (2.6), the minimax expected length for  $(1 - \alpha)$  level confidence intervals of  $\xi^\top \beta$  over  $\Theta(k, \mathbf{I}, \sigma_0)$  satisfies*

$$(5.2) \quad L_\alpha^*(\Theta(k, \mathbf{I}, \sigma_0), \xi^\top \beta) \asymp \frac{\|\xi\|_2}{\sqrt{n}}.$$

Compared with the minimax rate  $\frac{\|\xi\|_2}{\sqrt{n}} + \|\xi\|_2 k \frac{\log p}{n}$  for the unknown  $\Omega$  and  $\sigma$  case given in Theorem 1, the minimax rate in (5.2) is significantly different. With the prior knowledge of  $\Omega = \mathbf{I}$  and  $\sigma = \sigma_0$ , the above theorem shows that the minimax expected length of confidence intervals for  $\xi^\top \beta$  is always of parametric rate and in particular does not depend on the sparsity parameter  $k$ . In this case, adaptive confidence intervals for  $\xi^\top \beta$  is possible over the full range  $k \leq c \frac{n}{\log p}$ . A similar result for confidence intervals covering all  $\beta_i$  has been given in a recent paper [16]. The focus of [16] is on individual coordinates, not general linear functionals.

The minimax lower bound of Theorem 7 follows from the parametric lower bound of Theorem 1. As both  $\Omega$  and  $\sigma$  are known, the upper bound analysis is easier than the unknown  $\Omega$  and  $\sigma$  case and is similar to the one given in [16]. For completeness, we detail the construction of a confidence interval achieving the minimax length in (5.2) using the de-biasing method. We first randomly split the samples  $(X, y)$  into two subsamples  $(X^{(1)}, y^{(1)})$  and  $(X^{(2)}, y^{(2)})$  with sample sizes  $n_1$  and  $n_2$ , respectively. Without loss of generality, we assume that  $n$  is even and  $n_1 = n_2 = \frac{n}{2}$ . Let  $\hat{\beta}$  denote the Lasso estimator defined based on the sample  $(X^{(1)}, y^{(1)})$  with the proper tuning parameter  $\lambda = \sqrt{\frac{2.05 \log p}{n_1}} \sigma_0$ ,

$$(5.3) \quad \hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \frac{\|y^{(1)} - X^{(1)}\beta\|_2^2}{2n_1} + \lambda \sum_{j=1}^p \frac{\|X_{\cdot j}^{(1)}\|_2}{\sqrt{n_1}} |\beta_j|.$$

We define the following estimator of  $\xi^\top \beta$ ,

$$(5.4) \quad \bar{\mu} = \xi^\top \hat{\beta} + \frac{1}{n_2} \xi^\top \left(X^{(2)}\right)^\top \left(y^{(2)} - X^{(2)} \hat{\beta}\right).$$

Based on the estimator, we construct the following confidence interval

$$(5.5) \quad \text{CI}_\alpha^{\text{I}}(\xi^\top \beta, Z) = \left[ \bar{\mu} - 1.01 \frac{\|\xi\|_2}{\sqrt{n_2}} z_{\alpha_0/2} \sigma_0, \bar{\mu} + 1.01 \frac{\|\xi\|_2}{\sqrt{n_2}} z_{\alpha_0/2} \sigma_0 \right],$$

where  $\alpha_0 = \gamma_0 \alpha$  with  $0 < \gamma_0 < 1$ . It will be shown in the supplement [3] that the confidence interval proposed in (5.5) has valid coverage and achieves the minimax length in (5.2).

*5.2. Confidence intervals for linear functionals in the dense loading regime.* In marked contrast to the sparse loading regime, the prior knowledge of  $\Omega = \text{I}$  and  $\sigma = \sigma_0$  does not improve the minimax rate in the dense loading regime. That is, Theorem 4 remains true by replacing  $\Theta(k)$  and  $\Theta(k_1)$  with  $\Theta(k, \text{I}, \sigma_0)$  and  $\Theta(k_1, \text{I}, \sigma_0)$ , respectively. However, the cost of adaptation changes when there is prior knowledge of  $\Omega = \text{I}$  and  $\sigma = \sigma_0$ . The following theorem establishes the adaptivity lower bound in the dense loading regime.

**THEOREM 8.** *Suppose that  $0 < \alpha < \frac{1}{2}$  and  $k_1 \leq k \leq c \min \left\{ p^\gamma, \frac{n}{\log p} \right\}$  for some constants  $c > 0$  and  $0 \leq \gamma < \frac{1}{2}$ , then, for some constant  $c_1 > 0$ ,*

$$(5.6) \quad L_\alpha^*(\Theta(k_1, \text{I}, \sigma_0), \Theta(k, \text{I}, \sigma_0), \xi^\top \beta) \geq c_1 \|\xi\|_\infty \sigma_0 \max \left\{ \sqrt{k k_1} \sqrt{\frac{\log p}{n}}, \min \left\{ k \sqrt{\frac{\log p}{n}}, \frac{\sqrt{k}}{n^{\frac{1}{4}}} \right\} \right\}.$$

The lower bound in (5.6) is attainable. For reasons of space, the construction is omitted here. Under the framework (2.5), adaptive confidence intervals are still impossible, since for  $k_1 \ll k$ ,

$$L_\alpha^*(\Theta(k_1, \text{I}, \sigma_0), \Theta(k, \text{I}, \sigma_0), \xi^\top \beta) \gg L_\alpha^*(\Theta(k_1, \text{I}, \sigma_0), \xi^\top \beta).$$

Compared with Theorem 5, we observe that the cost of adaptation is reduced with the prior knowledge of  $\Omega = \text{I}$  and  $\sigma = \sigma_0$ .

**6. Discussion.** In the present paper we studied the minimaxity and adaptivity of confidence intervals for general linear functionals  $\xi^\top \beta$  with a sparse or dense loading  $\xi$  for the setting where  $\Omega$  and  $\sigma$  are unknown as well as the setting with the prior knowledge of  $\Omega = \text{I}$  and  $\sigma = \sigma_0$ . In the more typical case in practice where  $\Omega$  and  $\sigma$  are unknown, the adaptivity results are quite negative: With the exception of the ultra-sparse region for confidence intervals for  $\xi^\top \beta$  with a sparse loading  $\xi$ , it is necessary to know the true sparsity  $k$  in order to have guaranteed coverage probability and rate-optimal expected length. In contrast to estimation, knowledge of

the sparsity  $k$  is crucial to constructing honest confidence intervals. In this sense, the problem of constructing confidence intervals is much harder than the estimation problem.

The case of known  $\Omega = \mathbf{I}$  and  $\sigma = \sigma_0$  is strikingly different. The minimax expected length in the sparse loading regime is of order  $\frac{\|\xi\|_2}{\sqrt{n}}$  and in particular does not depend on  $k$  and adaptivity can be achieved over the full range of sparsity  $k \lesssim \frac{n}{\log p}$ . So in this case, the knowledge of  $\Omega$  and  $\sigma$  is very useful. On the other hand, in the dense loading regime the information on  $\Omega$  and  $\sigma$  is of limited use. In this case, the minimax rate and lack of adaptivity remain unchanged, compared with the unknown  $\Omega$  and  $\sigma$  case, although the cost of adaptation is reduced.

Regarding the construction of confidence intervals, there is a significant difference between the sparse and dense loading regimes. The de-biasing method is useful in the sparse loading regime since such a procedure reduces the bias but does not dramatically increase the variance. However, the de-biasing construction is not applicable to the dense loading regime since the cost of obtaining a near-unbiased estimator is to significantly increase the variance which would lead to an unnecessarily long confidence interval. An interesting open problem is the construction of a confidence interval for  $\xi^\top \beta$  achieving the minimax length where the sparsity  $q$  of the loading  $\xi$  is in the middle regime with  $cp^\gamma \leq q \leq cp^{2\gamma+\varsigma}$  for some  $0 < \varsigma < 1 - 2\gamma$ .

In addition to constructing confidence intervals for linear functionals, another interesting problem is constructing confidence balls for the whole vector  $\beta$ . Such has been considered in [17], where the authors established the impossibility of adaptive confidence balls for sparse linear regression. These problems are connected, but each has its own special features and the behaviors of the problems are different from each other. The connections and differences in adaptivity among various forms of confidence sets have also been observed in nonparametric function estimation problems. See, for example, [6] for adaptive confidence intervals for linear functionals, [13, 9] for adaptive confidence bands, and [8, 20] for adaptive confidence balls.

In the context of nonparametric function estimation, a general adaptation theory for confidence intervals for an arbitrary linear functional was developed in Cai and Low [6] over a collection of convex parameter spaces. It was shown that the key quantity that determines adaptivity is a geometric quantity called the between-class modulus of continuity. The convexity assumption on the parameter space in Cai and Low [6] is crucial for the adaptation theory. In high-dimensional linear regression, the parameter space is highly non-convex. The adaptation theory developed in [6] does not apply to the present setting of high-dimensional linear regression. It would be of

significant interest to develop a general adaptation theory for confidence intervals in such a non-convex setting.

**7. Proofs.** In this section, we prove two main results, Theorem 3 and minimax upper bound of Theorem 1. For reasons of space, the proofs of the other results are given in the supplement [3].

A key technical tool for the proof of the lower bound results is the following lemma which establishes the adaptivity over two nested parameter spaces. Such a formulation has been considered in [6] in the context of adaptive confidence intervals over convex parameter spaces under the Gaussian sequence model. However, the parameter space  $\Theta(k)$  considered in the high dimension setting is highly non-convex. The following lemma can be viewed as a generalization of [6] to the non-convex parameter space, where the lower bound argument requires testing for composite hypotheses.

Suppose that we observe a random variable  $Z$  which has a distribution  $\mathbf{P}_\theta$  where the parameter  $\theta$  belongs to the parameter space  $\mathcal{H}$ . Let  $\text{CI}_\alpha(\mathbf{T}, Z)$  be the confidence interval for the linear functional  $\mathbf{T}(\theta)$  with the guaranteed coverage  $1 - \alpha$  over the parameter space  $\mathcal{H}$ . Let  $\mathcal{H}_0$  and  $\mathcal{H}_1$  be subsets of the parameter space  $\mathcal{H}$  where  $\mathcal{H} = \mathcal{H}_0 \cup \mathcal{H}_1$ . Let  $\pi_{\mathcal{H}_i}$  denote the prior distribution supported on the parameter space  $\mathcal{H}_i$  for  $i = 0, 1$ . Let  $f_{\pi_{\mathcal{H}_i}}(z)$  denote the density function of the marginal distribution of  $Z$  with the prior  $\pi_{\mathcal{H}_i}$  on  $\mathcal{H}_i$  for  $i = 0, 1$ . More specifically,  $f_{\pi_{\mathcal{H}_i}}(z) = \int f_\theta(z) \pi_{\mathcal{H}_i}(\theta) d\theta$ .

Denote by  $\mathbb{P}_{\pi_{\mathcal{H}_i}}$  the marginal distribution of  $Z$  with the prior  $\pi_{\mathcal{H}_i}$  on  $\mathcal{H}_i$  for  $i = 0, 1$ . For any function  $g$ , we write  $\mathbb{E}_{\pi_{\mathcal{H}_0}}(g(Z))$  for the expectation of  $g(Z)$  with respect to the marginal distribution of  $Z$  with the prior  $\pi_{\mathcal{H}_0}$  on  $\mathcal{H}_0$ . We define the  $\chi^2$  distance between two density functions  $f_1$  and  $f_0$  by

$$(7.1) \quad \chi^2(f_1, f_0) = \int \frac{(f_1(z) - f_0(z))^2}{f_0(z)} dz = \int \frac{f_1^2(z)}{f_0(z)} dz - 1$$

and the total variation distance by  $\text{TV}(f_1, f_0) = \int |f_1(z) - f_0(z)| dz$ . It is well known that

$$(7.2) \quad \text{TV}(f_1, f_0) \leq \sqrt{\chi^2(f_1, f_0)}.$$

LEMMA 1. Assume  $\mathbf{T}(\theta) = \mu_0$  for  $\theta \in \mathcal{H}_0$  and  $\mathbf{T}(\theta) = \mu_1$  for  $\theta \in \mathcal{H}_1$  and  $\mathcal{H} = \mathcal{H}_0 \cup \mathcal{H}_1$ . For any  $\text{CI}_\alpha(\mathbf{T}, Z) \in \mathcal{I}_\alpha(\mathbf{T}, \mathcal{H})$ ,

$$(7.3) \quad L(\text{CI}_\alpha(\mathbf{T}, Z), \mathcal{H}) \geq L(\text{CI}_\alpha(\mathbf{T}, Z), \mathcal{H}_0) \geq |\mu_1 - \mu_0| \left( 1 - 2\alpha - \text{TV}\left(f_{\pi_{\mathcal{H}_1}}, f_{\pi_{\mathcal{H}_0}}\right) \right)_+.$$

7.1. *Proof of Lemma 1.* The supremum risk over  $\mathcal{H}_0$  is lower bounded by the Bayesian risk with the prior  $\pi_{\mathcal{H}_0}$  on  $\mathcal{H}_0$ ,

$$(7.4) \quad \sup_{\theta \in \mathcal{H}_0} \mathbb{E}_\theta L(\text{CI}_\alpha(\mathbb{T}, Z)) \geq \int_{\theta} \mathbb{E}_\theta L(\text{CI}_\alpha(\mathbb{T}, Z)) \pi_{\mathcal{H}_0}(\theta) d\theta = \mathbb{E}_{\pi_{\mathcal{H}_0}} L(\text{CI}_\alpha(\mathbb{T}, Z)).$$

By the definition of  $\text{CI}_\alpha(\mathbb{T}, Z) \in \mathcal{I}_\alpha(\mathbb{T}, \mathcal{H})$ , we have

$$(7.5) \quad \mathbb{P}_{\pi_{\mathcal{H}_i}}(\mu_i \in \text{CI}_\alpha(\mathbb{T}, Z)) = \int_{\theta} \mathbb{P}_\theta(\mu_i \in \text{CI}_\alpha(\mathbb{T}, Z)) \pi_{\mathcal{H}_i}(\theta) d\theta \geq 1 - \alpha,$$

for  $i = 0, 1$ . By the following inequality

$$\left| \mathbb{P}_{\pi_{\mathcal{H}_1}}(\mu_1 \in \text{CI}_\alpha(\mathbb{T}, Z)) - \mathbb{P}_{\pi_{\mathcal{H}_0}}(\mu_1 \in \text{CI}_\alpha(\mathbb{T}, Z)) \right| \leq \text{TV}(f_{\pi_{\mathcal{H}_1}}, f_{\pi_{\mathcal{H}_0}}),$$

then we have  $\mathbb{P}_{\pi_{\mathcal{H}_0}}(\mu_1 \in \text{CI}_\alpha(\mathbb{T}, Z)) \geq 1 - \alpha - \text{TV}(f_{\pi_{\mathcal{H}_1}}, f_{\pi_{\mathcal{H}_0}})$ . This together with (7.5) yields  $\mathbb{P}_{\pi_{\mathcal{H}_0}}(\mu_0, \mu_1 \in \text{CI}_\alpha(\mathbb{T}, Z)) \geq 1 - 2\alpha - \text{TV}(f_{\pi_{\mathcal{H}_1}}, f_{\pi_{\mathcal{H}_0}})$ , which leads to  $\mathbb{P}_{\pi_{\mathcal{H}_0}}(L(\text{CI}_\alpha(\mathbb{T}, Z)) \geq |\mu_1 - \mu_0|) \geq 1 - 2\alpha - \text{TV}(f_{\pi_{\mathcal{H}_1}}, f_{\pi_{\mathcal{H}_0}})$ . Hence,  $\mathbb{E}_{\pi_{\mathcal{H}_0}} L(\text{CI}_\alpha(\mathbb{T}, Z)) \geq (\mu_1 - \mu_0)(1 - 2\alpha - \text{TV}(f_{\pi_{\mathcal{H}_1}}, f_{\pi_{\mathcal{H}_0}}))_+$ . The lower bound (7.3) follows from inequality (7.4).

7.2. *Proof of Theorem 3.* The lower bound in (3.13) is involved with a parametric term and a non-parametric term. The proof of the parametric lower bound is postponed to the supplement. In the following, we will prove the non-parametric lower bound

$$(7.6) \quad \inf_{\text{CI}_\alpha(\xi^\top \beta, Z) \in \mathcal{I}_\alpha(\Theta(k), \xi^\top \beta)} \mathbb{E}_{\theta^*} L(\text{CI}_\alpha(\xi^\top \beta, Z)) \geq c_1 \|\xi\|_2 k \frac{\log p}{n} \sigma,$$

for some constant  $c_1 > 0$ . Without loss of generality, we assume  $\text{supp}(\xi) = \{1, \dots, \|\xi\|_0\}$ . We generate the orthogonal matrix  $M \in \mathbb{R}^{\|\xi\|_0 \times \|\xi\|_0}$  such that its first row is  $\frac{1}{\|\xi\|_2} \xi_{\text{supp}(\xi)}$  and define the orthogonal matrix  $Q$  as

$Q = \begin{pmatrix} M & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}$ . We transform both the design matrix  $X$  and the regression vector  $\beta$  and view the linear model (2.1) as  $y = V\psi + \epsilon$ , where  $V = XQ^\top$

and  $\psi = Q\beta$ . The transformed coefficient vector  $\psi^* = Q\beta^* = \begin{pmatrix} M\beta_{\text{supp}(\xi)}^* \\ \beta_{-\text{supp}(\xi)}^* \end{pmatrix}$

is of sparsity at most  $\|\xi\|_0 + k_1$ . The first coefficient  $\psi_1$  of  $\psi$  is  $\frac{1}{\|\xi\|_2} \xi^\top \beta$ . The covariance matrix  $\Psi$  of  $V_1$  is  $Q\Sigma Q^\top$  and its corresponding precision matrix is  $\Gamma = Q\Omega Q^\top$ . To represent the transformed observed data and parameter, we abuse the notation slightly and also use  $Z_i = (y_i, V_i)$  and  $\theta^* = (\psi^*, \mathbf{I}, \sigma)$ . We define the parameter space  $\mathcal{G}(k)$  of  $(\psi, \Gamma, \sigma)$  as

$$(7.7) \quad \mathcal{G}(k) = \{(\psi, \Gamma, \sigma) : \psi = Q\beta, \Gamma = Q\Omega Q^\top \text{ for } (\beta, \Omega, \sigma) \in \Theta(k)\}.$$

For a given  $Q$ , there exists a bijective mapping between  $\Theta(k)$  and  $\mathcal{G}(k)$ . To show that  $(\psi, \Gamma, \sigma) \in \mathcal{G}(k)$ , it is equivalent to show  $(Q^\top \psi, Q^\top \Gamma Q, \sigma) \in \Theta(k)$ . Let  $\mathcal{I}_\alpha(\mathcal{G}(k), \psi_1)$  denote the set of confidence intervals for  $\psi_1 = \frac{1}{\|\xi\|_2} \xi^\top \beta$  with guaranteed coverage over  $\mathcal{G}(k)$ . If  $\text{CI}_\alpha(\psi_1, Z) \in \mathcal{I}_\alpha(\mathcal{G}(k), \psi_1)$ , then  $\|\xi\|_2 \text{CI}_\alpha(\psi_1, Z) \in \mathcal{I}_\alpha(\Theta(k), \xi^\top \beta)$ ; If  $\text{CI}_\alpha(\xi^\top \beta, Z) \in \mathcal{I}_\alpha(\Theta(k), \xi^\top \beta)$ , then  $\frac{1}{\|\xi\|_2} \text{CI}_\alpha(\xi^\top \beta, Z) \in \mathcal{I}_\alpha(\mathcal{G}(k), \psi_1)$ . Because of such one to one correspondence, we have

$$(7.8) \quad \inf_{\text{CI}_\alpha(\xi^\top \beta, Z) \in \mathcal{I}_\alpha(\Theta(k), \xi^\top \beta)} \mathbb{E}_{\theta^*} L(\text{CI}_\alpha(\xi^\top \beta, Z)) = \|\xi\|_2 \inf_{\text{CI}_\alpha(\psi_1, Z) \in \mathcal{I}_\alpha(\mathcal{G}(k), \psi_1)} \mathbb{E}_{\theta^*} L(\text{CI}_\alpha(\psi_1, Z)).$$

By (7.6) and (7.8), we reduce the problem to

$$(7.9) \quad \inf_{\text{CI}_\alpha(\psi_1, Z) \in \mathcal{I}_\alpha(\mathcal{G}(k), \psi_1)} \mathbb{E}_{\theta^*} L(\text{CI}_\alpha(\psi_1, Z)) \geq ck \frac{\log p}{n} \sigma.$$

Under the Gaussian random design model,  $Z_i = (y_i, V_i) \in \mathbb{R}^{p+1}$  follows a joint Gaussian distribution with mean 0. Let  $\Sigma^z$  denotes the covariance matrix of  $Z_i$ . Decompose  $\Sigma^z$  into blocks  $\begin{pmatrix} \Sigma_{yy}^z & (\Sigma_{vy}^z)^\top \\ \Sigma_{vy}^z & \Sigma_{vv}^z \end{pmatrix}$ , where  $\Sigma_{yy}^z$ ,  $\Sigma_{vv}^z$  and  $\Sigma_{vy}^z$  denote the variance of  $y$ , the variance of  $V$  and the covariance of  $y$  and  $V$ , respectively. We define the function  $h : \Sigma^z \rightarrow (\psi, \Gamma, \sigma)$  as  $h(\Sigma^z) = \left( (\Sigma_{vv}^z)^{-1} \Sigma_{vy}^z, (\Sigma_{vv}^z)^{-1}, \Sigma_{yy}^z - (\Sigma_{vy}^z)^\top (\Sigma_{vv}^z)^{-1} \Sigma_{vy}^z \right)$ . The function  $h$  is bijective and its inverse mapping  $h^{-1} : (\psi, \Gamma, \sigma) \rightarrow \Sigma^z$  is

$$h^{-1}((\psi, \Gamma, \sigma)) = \begin{pmatrix} \psi^\top \Gamma^{-1} \psi + \sigma^2 & \psi^\top \Gamma^{-1} \\ \Gamma^{-1} \psi & \Gamma^{-1} \end{pmatrix}.$$

The null space is taken as  $\mathcal{H}_0 = \{(\psi^*, \mathbf{I}, \sigma)\}$  and  $\pi_{\mathcal{H}_0}$  denotes the point mass prior at this point. The proof is divided into three steps:

1. Construct  $\mathcal{H}_1$  and show that  $\mathcal{H}_1 \subset \mathcal{G}(k)$ ;
2. Control the distribution distance  $\text{TV}(f_{\pi_{\mathcal{H}_1}}, f_{\pi_{\mathcal{H}_0}})$ ;
3. Calculate the distance  $\mu_1 - \mu_0$  where  $\mu_0 = \psi_1^*$  and  $\mu_1 = \psi_1$  with  $(\psi, \Gamma, \sigma) \in \mathcal{H}_1$ . We show that  $\mu_1 = \psi_1$  is a fixed constant for all  $(\psi, \Gamma, \sigma) \in \mathcal{H}_1$  and then apply Lemma 1.

**Step 1.** We construct the alternative hypothesis parameter space  $\mathcal{H}_1$ . Let  $\Sigma_0^z$  denote the covariance matrix of  $Z_i$  corresponding to  $(\psi^*, \mathbf{I}, \sigma) \in \mathcal{H}_0$ . Let  $S_1 = \text{supp}(\psi^*) \cup \{1\}$  and  $S = S_1 \setminus \{1\}$ . Let  $k_*$  denote the size of  $S$  and  $p_1$  denote the size of  $S_1^c$  and we have  $k_* \leq k_1 + q$  and  $p_1 \geq p - k_* - 1 \geq cp$ . Without loss of generality, let  $S = \{2, \dots, k_* + 1\}$ . We have the following



expression for the covariance matrix of  $Z_i$  under the null,

$$(7.10) \quad \Sigma_0^z = \begin{pmatrix} \|\psi^*\|_2^2 + \sigma^2 & \psi_1^* & (\psi_S^*)^\top & \mathbf{0}_{1 \times p_1} \\ \psi_1^* & 1 & \mathbf{0}_{1 \times k_*} & \mathbf{0}_{1 \times p_1} \\ \psi_S^* & \mathbf{0}_{k_* \times 1} & \mathbf{I}_{k_* \times k_*} & \mathbf{0}_{k_* \times p_1} \\ \mathbf{0}_{p_1 \times 1} & \mathbf{0}_{p_1 \times 1} & \mathbf{0}_{p_1 \times k_*} & \mathbf{I}_{p_1 \times p_1} \end{pmatrix},$$

To construct  $\mathcal{H}_1$ , we define the following set,

$$(7.11) \quad \ell \left( p_1, \frac{\zeta_0}{2} k, \rho \right) = \left\{ \boldsymbol{\delta} : \boldsymbol{\delta} \in \mathbb{R}^{p_1}, \|\boldsymbol{\delta}\|_0 = \frac{\zeta_0}{2} k, \boldsymbol{\delta}_i \in \{0, \rho\} \text{ for } 1 \leq i \leq p_1 \right\}.$$

Define the parameter space  $\mathcal{F}$  for  $\Sigma^z$  by  $\mathcal{F} = \left\{ \Sigma_{\boldsymbol{\delta}}^z : \boldsymbol{\delta} \in \ell \left( p_1, \frac{\zeta_0}{2} k, \rho \right) \right\}$ , where

$$(7.12) \quad \Sigma_{\boldsymbol{\delta}}^z = \begin{pmatrix} \|\psi^*\|_2^2 + \sigma^2 & \psi_1^* & (\psi_S^*)^\top & \rho_0 \boldsymbol{\delta}^\top \\ \psi_1^* & 1 & \mathbf{0}_{1 \times k_*} & \boldsymbol{\delta}^\top \\ \psi_S^* & \mathbf{0}_{k_* \times 1} & \mathbf{I}_{k_* \times k_*} & \mathbf{0}_{k_* \times p_1} \\ \rho_0 \boldsymbol{\delta} & \boldsymbol{\delta} & \mathbf{0}_{p_1 \times k_*} & \mathbf{I}_{p_1 \times p_1} \end{pmatrix}.$$

Then we construct the alternative hypothesis space  $\mathcal{H}_1$  for  $(\psi, \Gamma, \sigma)$ , which is induced by the mapping  $h$  and the parameter space  $\mathcal{F}$ ,

$$(7.13) \quad \mathcal{H}_1 = \{(\psi, \Gamma, \sigma) : (\psi, \Gamma, \sigma) = h(\Sigma^z) \text{ for } \Sigma^z \in \mathcal{F}\}$$

In the following, we show that  $\mathcal{H}_1 \subset \mathcal{G}(k)$ . It is necessary to identify  $(\psi, \Gamma, \sigma) = h(\Sigma^z)$  for  $\Sigma^z \in \mathcal{F}$  and show  $(Q^\top \psi, Q^\top \Gamma Q, \sigma) \in \Theta(k)$ . Firstly, we identify the expression  $\mathbb{E}(y_i | V_{i,\cdot})$  under the alternative joint distribution (7.12). Assuming  $y_i = V_{i1} \psi_1 + V_{i,S} \psi_S + V_{i,S_1^c} \psi_{S_1^c} + \epsilon'_i$ , we have

$$(7.14) \quad \psi_1 = \frac{-\|\boldsymbol{\delta}\|_2^2 \rho_0 + \psi_1^*}{1 - \|\boldsymbol{\delta}\|_2^2}, \quad \psi_S = \psi_S^*, \quad \psi_{S_1^c} = (\rho_0 - \psi_1) \boldsymbol{\delta},$$

and

$$(7.15) \quad \text{Var}(\epsilon'_i) = \sigma^2 - \frac{\|\boldsymbol{\delta}\|_2^2 (\rho_0 - \psi_1^*)^2}{1 - \|\boldsymbol{\delta}\|_2^2} \leq \sigma^2 \leq M_2.$$

Based on (7.14), the sparsity of  $\psi$  in the alternative hypothesis space is upper bounded by  $1 + |\text{supp}(\psi_S^*)| + |\text{supp}(\boldsymbol{\delta})| \leq \left(1 - \frac{\zeta_0}{4}\right) k$ , and hence the sparsity of the corresponding  $\beta = Q^\top \psi$  is controlled by

$$(7.16) \quad \|\beta\|_0 \leq \left(1 - \frac{\zeta_0}{4}\right) k + q \leq k.$$

Secondly, we show that  $\Omega = Q^\top \Gamma Q$  satisfies the condition  $\frac{1}{M_1} \leq \lambda_{\min}(\Omega) \leq \lambda_{\max}(\Omega) \leq M_1$ . The covariance matrix  $\Psi$  of  $V_{i,\cdot}$  in the alternative hypothesis parameter space is expressed as

$$(7.17) \quad \Psi = \begin{pmatrix} 1 & \mathbf{0}_{1 \times k_*} & \mathbf{0}_{k_* \times p_1} \\ \mathbf{0}_{k_* \times 1} & \mathbf{I}_{k_* \times k_*} & \mathbf{0}_{k_* \times p_1} \\ \mathbf{0}_{p_1 \times 1} & \mathbf{0}_{p_1 \times k_*} & \mathbf{I}_{p_1 \times p_1} \end{pmatrix} + \begin{pmatrix} 0 & \mathbf{0}_{1 \times k_*} & \boldsymbol{\delta}^\top \\ \mathbf{0}_{k_* \times 1} & \mathbf{0}_{k_* \times k_*} & \mathbf{0}_{k_* \times p_1} \\ \boldsymbol{\delta} & \mathbf{0}_{p_1 \times k_*} & \mathbf{0}_{p_1 \times p_1} \end{pmatrix}.$$

Since the second matrix on the above equation is of spectral norm  $\|\boldsymbol{\delta}\|_2$ , Weyl's inequality leads to  $\max\{|\lambda_{\min}(\Psi) - 1|, |\lambda_{\max}(\Psi) - 1|\} \leq \|\boldsymbol{\delta}\|_2$ . When  $\|\boldsymbol{\delta}\|_2$  is chosen such that  $\|\boldsymbol{\delta}\|_2 \leq \min\left\{1 - \frac{1}{M_1}, M_1 - 1\right\}$ , then we have  $\frac{1}{M_1} \leq \lambda_{\min}(\Psi) \leq \lambda_{\max}(\Psi) \leq M_1$ . Since  $\Omega$  and  $\Gamma = Q\Omega Q^\top$  have the same eigenvalues, we have  $\frac{1}{M_1} \leq \lambda_{\min}(\Omega) \leq \lambda_{\max}(\Omega) \leq M_1$ . Combined with (7.15) and (7.16), we show that  $\mathcal{H}_1 \subset \mathcal{G}(k)$ .

**Step 2.** To control  $\text{TV}(f_{\pi_{\mathcal{H}_1}}, f_{\pi_{\mathcal{H}_0}})$ , it is sufficient to control  $\chi^2(f_{\pi_{\mathcal{H}_1}}, f_{\pi_{\mathcal{H}_0}})$  and apply (7.2). Let  $\pi$  denote the uniform prior on  $\boldsymbol{\delta}$  over  $\ell\left(p_1, \frac{\zeta_0}{2}k, \rho\right)$ . Note that this uniform prior  $\pi$  induces a prior distribution  $\pi_{\mathcal{H}_1}$  over the parameter space  $\mathcal{H}_1$ . Let  $\mathbb{E}_{\boldsymbol{\delta}, \tilde{\boldsymbol{\delta}}}$  denote the expectation with respect to the independent random variables  $\boldsymbol{\delta}, \tilde{\boldsymbol{\delta}}$  with uniform prior  $\pi$  over the parameter space  $\ell\left(p_1, \frac{\zeta_0}{2}k, \rho\right)$ . The following lemma controls the  $\chi^2$  distance between the null and the mixture over the alternative distribution.

LEMMA 2. *Let  $f_1 = \left(\sigma^2 + (\psi_1^*)^2 - \rho_0 \psi_1^*\right)$ . Then*

$$(7.18) \quad \chi^2(f_{\pi_{\mathcal{H}_1}}, f_{\pi_{\mathcal{H}_0}}) + 1 = \mathbb{E}_{\boldsymbol{\delta}, \tilde{\boldsymbol{\delta}}} \left(1 - \frac{1}{\sigma^2} (\rho_0 (\rho_0 - \psi_1^*) + f_1) \boldsymbol{\delta}^\top \tilde{\boldsymbol{\delta}}\right)^{-n}.$$

The following lemma is useful in controlling the right hand side of (7.18).

LEMMA 3. *Let  $J$  be a Hypergeometric  $(p, k, k)$  variable with  $\mathbb{P}(J = j) = \frac{\binom{k}{j} \binom{p-k}{k-j}}{\binom{p}{k}}$ , then*

$$(7.19) \quad \mathbb{E} \exp(tJ) \leq e^{\frac{k^2}{p-k}} \left(1 - \frac{k}{p} + \frac{k}{p} \exp(t)\right)^k.$$

Taking  $\rho_0 = \psi_1^* + \sigma$ , we have  $\frac{1}{\sigma^2} (\rho_0 (\rho_0 - \psi_1^*) + f_1) = 2$  and by Lemma 2,

$$\chi^2(f_{\pi_{\mathcal{H}_1}}, f_{\pi_{\mathcal{H}_0}}) + 1 = \mathbb{E}_{\boldsymbol{\delta}, \tilde{\boldsymbol{\delta}}} \left(1 - 2\boldsymbol{\delta}^\top \tilde{\boldsymbol{\delta}}\right)^{-n}.$$

By the inequality  $\frac{1}{1-x} \leq \exp(2x)$  for  $x \in [0, \frac{\log 2}{2}]$ , if  $\delta^\top \tilde{\delta} \leq \frac{\zeta_0}{2} k \rho^2 < \frac{\log 2}{4}$ , then  $(1 - 2\delta^\top \tilde{\delta})^{-n} \leq \exp(4n\delta^\top \tilde{\delta})$ . By Lemma 3, we further have

$$\begin{aligned} \mathbb{E}_{\delta, \tilde{\delta}} \exp(4n\delta^\top \tilde{\delta}) &= \mathbb{E} \exp(4Jn\rho^2) \leq e^{\frac{\zeta_0^2 k^2}{4p_1 - 2\zeta_0 k}} \left(1 - \frac{\zeta_0 k}{2p_1} + \frac{\zeta_0 k}{2p_1} \exp(4n\rho^2)\right)^{\frac{\zeta_0}{2} k} \\ &\leq e^{\frac{\zeta_0^2 k^2}{4p_1 - 2\zeta_0 k}} \left(1 - \frac{\zeta_0 k}{2p_1} + \frac{\zeta_0 k}{2p_1} \sqrt{\frac{4p_1}{\zeta_0^2 k^2}}\right)^{\frac{\zeta_0}{2} k} \leq e^{\frac{c^2 \zeta_0^2 p^{2\gamma}}{4p_1 - 2c\zeta_0 p^\gamma}} \left(1 + \frac{1}{\sqrt{p_1}}\right)^{\frac{c\zeta_0}{2} p^\gamma}, \end{aligned}$$

where the second inequality follows by plugging in  $\rho = \sqrt{\frac{\log \frac{4p_1}{\zeta_0^2 k^2}}{8n}}$  and the last inequality follows by  $k \leq cp^\gamma$ . If  $k \leq c \left\{ \frac{n}{\log p}, p^\gamma \right\}$ , where  $0 \leq \gamma < \frac{1}{2}$  and  $c$  is a sufficient small positive constant, then  $k\rho^2 < \min \left\{ \frac{\log 2}{2\zeta_0}, \left(1 - \frac{1}{M_1}\right)^2, 1 \right\}$  and hence

$$(7.20) \quad \chi^2(f_{\pi_{\mathcal{H}_1}}, f_{\pi_{\mathcal{H}_0}}) \leq \left(\frac{1}{2} - \alpha\right)^2 \quad \text{and} \quad \text{TV}(f_{\pi_{\mathcal{H}_1}}, f_{\pi_{\mathcal{H}_0}}) \leq \frac{1}{2} - \alpha.$$

**Step 3.** We calculate the distance between  $\mu_1$  and  $\mu_0$ . Under  $\mathcal{H}_0$ ,  $\mu_0 = \psi_1^*$ . Under  $\mathcal{H}_1$ ,  $\mu_1 = \psi_1 = \frac{-\|\delta\|_2^2 \rho_0 + \psi_1^*}{1 - \|\delta\|_2^2}$ . For  $\delta \in \ell\left(p_1, \frac{\zeta_0}{2} k, \rho\right)$ ,  $\|\delta\|_2^2 = \frac{\zeta_0}{2} k \rho^2$  and  $\mu_1 = \psi_1 = \frac{-\frac{\zeta_0}{2} k \rho^2 (\psi_1^* + \sigma) + \psi_1^*}{1 - \frac{\zeta_0}{2} k \rho^2}$ . Since  $\rho$  is selected as fixed,  $\mu_1 = \psi_1$  is a fixed constant for  $(\psi, \Omega, \sigma) \in \mathcal{H}_1$ . Note that  $\mu_1 - \mu_0 = \frac{\|\delta\|_2^2 (\psi_1^* - \rho_0)}{1 - \|\delta\|_2^2} = \frac{-\sigma \|\delta\|_2^2}{1 - \|\delta\|_2^2}$ , and it follows that  $|\mu_1 - \mu_0| = \sigma \frac{\|\delta\|_2^2}{1 - \|\delta\|_2^2} \geq ck \frac{\log \frac{4p_1}{\zeta_0^2 k^2}}{n} \sigma$ . Combined with (7.2) and (7.20), Lemma 1 leads to (7.9). By (7.8), we establish (3.13).

**7.3. Proof of upper bound in Theorem 1.** The following proposition establishes the coverage property and the expected length of the constructed confidence interval constructed in (3.11). Such a confidence interval achieves the minimax length in (3.1).

**PROPOSITION 1.** *Suppose that  $k \leq c_* \frac{n}{\log p}$ , where  $c_*$  is a small positive constant, then*

$$(7.21) \quad \liminf_{n, p \rightarrow \infty} \inf_{\theta \in \Theta(k)} \mathbb{P}_\theta(\xi^\top \beta \in \text{CI}_\alpha^S(\xi^\top \beta, Z)) \geq 1 - \alpha,$$

and

$$(7.22) \quad L(\text{CI}_\alpha^S(\xi^\top \beta, Z), \Theta(k)) \leq C \|\xi\|_2 \left( k \frac{\log p}{n} + \frac{1}{\sqrt{n}} \right),$$

for some constant  $C > 0$ .

In the following, we are going to prove Proposition 1. By normalizing the columns of  $X$  and the true sparse vector  $\beta$ , the linear regression model can be expressed as

$$(7.23) \quad y = Wd + \epsilon, \quad \text{with } W = XD, \quad d = D^{-1}\beta \text{ and } \epsilon \sim N(0, \sigma^2 \mathbf{I}),$$

where

$$(7.24) \quad D = \text{diag} \left( \frac{\sqrt{n}}{\|X_{\cdot j}\|_2} \right)_{j \in [p]}$$

denotes the  $p \times p$  diagonal matrix with  $(j, j)$  entry to be  $\frac{\sqrt{n}}{\|X_{\cdot j}\|_2}$ . Take  $\delta_0 = 1.0048$  and  $\eta_0 = 0.01$ , and we have  $\lambda_0 = (1 + \eta_0) \sqrt{\frac{2\delta_0 \log p}{n}}$ . Take  $\epsilon_0 = \frac{2.01}{\eta_0} + 1 = 202$ ,  $\nu_0 = 0.01$ ,  $C_1 = 2.25$ ,  $c_0 = \frac{1}{6}$  and  $C_0 = 3$ . Rather than use the constants directly in the following discussion, we use  $\delta_0, \eta_0, \epsilon_0, \nu_0, C_1, C_0$  and  $c_0$  to represent the above fixed constants in the following discussion. We also assume that  $\frac{\log p}{n} \leq \frac{1}{25}$  and  $\delta_0 \log p > 2$ . Define the  $l_1$  cone invertibility factor ( $CIF_1$ ) as follows,

$$(7.25) \quad CIF_1(\alpha_0, K, W) = \inf \left\{ \frac{|K| \left\| \frac{W^\top W}{n} u \right\|_\infty}{\|u_K\|_1} : \|u_{K^c}\|_1 \leq \alpha_0 \|u_K\|_1, u \neq 0 \right\},$$

where  $K$  is an index set. Define  $\sigma^{ora} = \frac{1}{\sqrt{n}} \|y - X\beta\|_2 = \frac{1}{\sqrt{n}} \|y - Wd\|_2$ ,

$$(7.26) \quad T = \{k : |d_k| \geq \lambda_0 \sigma^{ora}\}, \quad \tau = (1 + \epsilon_0) \lambda_0 \max \left\{ \frac{4}{\sigma^{ora}} \|d_{T^c}\|_1, \frac{8\lambda_0 |T|}{CIF_1(2\epsilon_0 + 1, T, W)} \right\}.$$

To facilitate the proof, we define the following events for the random design  $X$  and the error  $\epsilon$ ,

$$G_1 = \left\{ \frac{2}{5} \frac{1}{\sqrt{M_1}} < \frac{\|X_{\cdot j}\|_2}{\sqrt{n}} < \frac{7}{5} \sqrt{M_1} \text{ for } 1 \leq j \leq p \right\},$$

$$G_2 = \left\{ \left| \frac{(\sigma^{ora})^2}{\sigma^2} - 1 \right| \leq 2 \sqrt{\frac{\log p}{n}} + 2 \frac{\log p}{n} \right\},$$

$$\begin{aligned}
G_3 &= \left\{ \max \left\{ \left| \frac{\xi^\top \widehat{\Sigma} \xi}{\xi^\top \Sigma \xi} - 1 \right|, \left| \frac{u^\top \widehat{\Sigma} u}{\xi^\top \Omega \xi} - 1 \right| \right\} \leq 2\sqrt{\frac{\log p}{n}} + 2\frac{\log p}{n} \right\}, \text{ where } u = \Omega \xi, \\
G_4 &= \left\{ \kappa(X, k, \alpha) \geq \frac{1}{4\sqrt{\lambda_{\max}(\Omega)}} - \frac{9}{\sqrt{\lambda_{\min}(\Omega)}} (1 + \alpha) \sqrt{k \frac{\log p}{n}} \right\}, \\
G_5 &= \left\{ \frac{\|W^T \epsilon\|_\infty}{n} \leq \sigma \sqrt{\frac{2\delta_0 \log p}{n}} \right\}, \\
S_1 &= \left\{ \frac{\|W^T \epsilon\|_\infty}{n} \leq \sigma^{\text{ora}} \lambda_0 \frac{\epsilon_0 - 1}{\epsilon_0 + 1} (1 - \tau) \right\}, \\
S_2 &= \{(1 - \nu_0) \hat{\sigma} \leq \sigma \leq (1 + \nu_0) \hat{\sigma}\}, \\
B_1 &= \left\{ \|\xi^\top \Omega \widehat{\Sigma} - \xi^\top\|_\infty \leq \lambda_n \right\}, \text{ where } \lambda_n = 4C_0 M_1^2 \|\xi\|_2 \sqrt{\frac{\log p}{n}}.
\end{aligned}$$

Define  $G = \cap_{i=1}^5 G_i$  and  $S = \cap_{i=1}^2 S_i$ . The following lemmas control the probability of events  $G$ ,  $S$  and  $B_1$ . The detailed proofs of Lemma 4, 5 and 6 are in the supplement.

LEMMA 4.

$$(7.27) \quad \mathbb{P}_\theta(G) \geq 1 - \frac{6}{p} - 2p^{1-C_1} - \frac{1}{2\sqrt{\pi\delta_0 \log p}} p^{1-\delta_0} - c' \exp(-cn),$$

and

$$(7.28) \quad \mathbb{P}_\theta(B_1) \geq 1 - 2p^{1-c_0} C_0^2,$$

where  $c$  and  $c'$  are universal positive constants. If  $k \leq c_* \frac{n}{\log p}$ , then

$$(7.29) \quad \mathbb{P}_\theta(G \cap S) \geq \mathbb{P}_\theta(G) - 2 \exp\left(-\left(\frac{g_0 + 1 - \sqrt{2g_0 + 1}}{2}\right)n\right) - c'' \frac{1}{\sqrt{\log p}} p^{1-\delta_0},$$

where  $c_*$  and  $c''$  are universal positive constants and  $g_0 = \frac{\nu_0}{2+3\nu_0}$ .

The following lemma establishes a data-dependent upper bound for the term  $\|\widehat{\beta} - \beta\|_1$ .

LEMMA 5. *On the event  $G \cap S$ ,*

$$(7.30) \quad \|\widehat{\beta} - \beta\|_1 \leq (2 + 2\epsilon_0) \frac{\sqrt{n}}{\min \|X_{\cdot j}\|_2} l(Z, k),$$

where

$$(7.31) \quad l(Z, k) = \max \left\{ k\lambda_0\sigma^{ora}, \frac{(2 + 2\epsilon_0) \max \|X_{\cdot j}\|_2^2 \left( \sigma \sqrt{\frac{2\delta_0 \log p}{n}} + \lambda_0 \hat{\sigma} \right) k}{n\kappa^2 \left( X, k, (1 + 2\epsilon_0) \left( \frac{\max \|X_{\cdot j}\|_2}{\min \|X_{\cdot j}\|_2} \right) \right)} \right\}.$$

The following lemma controls the radius of the confidence interval.

LEMMA 6. *On the event  $G \cap S \cap B_1$ , there exists  $p_0$  such that if  $p \geq p_0$ ,*

$$(7.32) \quad \rho_1(k) \leq C \|\xi\|_2 \left( \frac{1}{\sqrt{n}} + \frac{k \log p}{n} \right) \sigma \leq \|\xi\|_2 \log p \left( \frac{1}{\sqrt{n}} + \frac{k \log p}{n} \right) \hat{\sigma},$$

and

$$(7.33) \quad \rho_2(k) \leq Ck \sqrt{\frac{\log p}{n}} \sigma \leq \log p \left( k \sqrt{\frac{\log p}{n}} \hat{\sigma} \right).$$

In the following, we establish the coverage property of the proposed confidence interval. By the definition of  $\tilde{\mu}$  in (3.6), we have

$$(7.34) \quad \tilde{\mu} - \xi^\top \beta = \frac{1}{n} \hat{u}^\top X^\top \epsilon + \left( \xi^\top - \hat{u}^\top \hat{\Sigma} \right) \left( \hat{\beta} - \beta \right).$$

We now construct a confidence interval for the variance term  $\frac{1}{n} \hat{u}^\top X^\top \epsilon$  by normal distribution and a high probability upper bound for the bias term  $\left( \xi^\top - \hat{u}^\top \hat{\Sigma} \right) \left( \hat{\beta} - \beta \right)$ . Since  $\epsilon$  is independent of  $X$  and  $\hat{u}$  and  $\hat{\Sigma}$  is a function of  $X$ , we have  $\frac{1}{n} \hat{u}^\top X^\top \epsilon \mid X \sim N \left( 0, \sigma^2 \frac{\hat{u}^\top \hat{\Sigma} \hat{u}}{n} \right)$ , and

$$\mathbb{P}_{\epsilon \mid X} \left( \frac{1}{n} \hat{u}^\top X^\top \epsilon \in \left( -\sqrt{\frac{\hat{u}^\top \hat{\Sigma} \hat{u}}{n}} \sigma z_{\alpha/2}, \sqrt{\frac{\hat{u}^\top \hat{\Sigma} \hat{u}}{n}} \sigma z_{\alpha/2} \right) \mid X \right) = 1 - \alpha.$$

By (7.34), we have  $\mathbb{P}_{\epsilon \mid X} \left( \xi^\top \beta \in \text{CI}_0(Z, k) \mid X \right) = 1 - \alpha$ , where

$$\begin{aligned} \text{CI}_0(Z, k) = & \left[ \tilde{\mu} - \left( \xi^\top - \hat{u}^\top \hat{\Sigma} \right) \left( \hat{\beta} - \beta \right) - \sqrt{\frac{\hat{u}^\top \hat{\Sigma} \hat{u}}{n}} \sigma z_{\alpha/2}, \right. \\ & \left. \tilde{\mu} - \left( \xi^\top - \hat{u}^\top \hat{\Sigma} \right) \left( \hat{\beta} - \beta \right) + \sqrt{\frac{\hat{u}^\top \hat{\Sigma} \hat{u}}{n}} \sigma z_{\alpha/2} \right]. \end{aligned}$$

Integrating with respect to  $X$ , we have

$$(7.35) \quad \mathbb{P}_\theta(\xi^\top \beta \in \text{CI}_0(Z, k)) = \int \mathbb{P}_{\epsilon|x}(\xi^\top \beta \in \text{CI}_0(Z, k)|x) f(x) dx = 1 - \alpha.$$

Since  $\left| (\xi^\top - \widehat{u}^\top \widehat{\Sigma}) (\widehat{\beta} - \beta) \right| \leq \|\xi^\top - \widehat{u}^\top \widehat{\Sigma}\|_\infty \|\widehat{\beta} - \beta\|_1$ , on the event  $S \cap G$ , Lemma 5 and the constraint in (3.5) lead to

$$(7.36) \quad \|\xi^\top - \widehat{u}^\top \widehat{\Sigma}\|_\infty \|\widehat{\beta} - \beta\|_1 \leq \lambda_n (2 + 2\epsilon_0) \frac{\sqrt{n}}{\min \|X_{\cdot j}\|_2} l(Z, k),$$

where  $l(Z, k)$  is defined in (7.31). On the event  $G \cap S$ , we also have  $\sigma \leq (1 + \nu_0) \hat{\sigma}$  and  $\sigma^{ora} \leq (1 + \nu_0) \sqrt{1 + 2\sqrt{\frac{\log p}{n}} + 2\frac{\log p}{n}} \hat{\sigma}$ . We define the following confidence interval to facilitate the discussion,  $\text{CI}_1(Z, k) = [\tilde{\mu} - l_k, \tilde{\mu} + l_k]$ , where  $l_k = (1 + \nu_0) \sqrt{\frac{\widehat{u}^\top \widehat{\Sigma} \widehat{u}}{n}} z_{\alpha/2} \hat{\sigma} + C_1(X, k) \|\xi\|_2 k \frac{\log p}{n} \hat{\sigma}$ . On the event  $G \cap S$ , we have

$$(7.37) \quad \text{CI}_0(Z, k) \subset \text{CI}_1(Z, k).$$

On the event  $S_2$ , if  $p \geq \exp(2M_2)$ , then  $\hat{\sigma} \leq \frac{1}{1-\nu_0} \sigma \leq \frac{1}{1-\nu_0} M_2 < \log p$ . Hence, the event  $A$  holds and  $\text{CI}_\alpha^S(\xi^\top \beta, Z) = [\tilde{\mu} - \rho_1(k), \tilde{\mu} + \rho_1(k)]$ . By Lemma 6, on the event  $G \cap S \cap B_1$ , if  $p \geq \max\{p_0, \exp(2M_2)\}$ , we have  $\rho_1(k) = l_k$ , and hence

$$(7.38) \quad \text{CI}_1(Z, k) = \text{CI}_\alpha^S(\xi^\top \beta, Z).$$

We have the following bound on the coverage probability,

$$\begin{aligned} & \mathbb{P}_\theta(\{\xi^\top \beta \in \text{CI}_\alpha^S(\xi^\top \beta, Z)\}) \geq \mathbb{P}_\theta(\{\xi^\top \beta \in \text{CI}_0(Z, k)\} \cap S \cap G \cap B_1) \\ & \geq \mathbb{P}_\theta(\{\xi^\top \beta \in \text{CI}_0(Z, k)\}) - \mathbb{P}_\theta((S \cap G \cap B_1)^c) = 1 - \alpha - \mathbb{P}_\theta((S \cap G \cap B_1)^c) \\ & = \mathbb{P}_\theta(S \cap G \cap B_1) - \alpha, \end{aligned}$$

where the first inequality follows from (7.37) and (7.38) and the first equality follows from (7.35). Combined with Lemma 4, we establish (7.21). We control the expected length as follows,

$$(7.39) \quad \begin{aligned} & \mathbb{E}_\theta L(\text{CI}_\alpha^S(\xi^\top \beta, Z)) = \mathbb{E}_\theta L(\text{CI}_\alpha^S(\xi^\top \beta, Z)) \mathbf{1}_A \\ & = \mathbb{E}_\theta L(\text{CI}_\alpha^S(\xi^\top \beta, Z)) \mathbf{1}_{A \cap (S \cap G \cap B_1)} + \mathbb{E}_\theta L(\text{CI}_\alpha^S(\xi^\top \beta, Z)) \mathbf{1}_{A \cap (S \cap G \cap B_1)^c} \\ & \leq C \|\xi\|_2 \left( k \frac{\log p}{n} + \frac{1}{\sqrt{n}} \right) \sigma + \|\xi\|_2 (\log p)^2 \left( \frac{1}{\sqrt{n}} + \frac{k \log p}{n} \right) \mathbb{P}_\theta((S \cap G \cap B_1)^c) \\ & \leq C \|\xi\|_2 \left( k \frac{\log p}{n} + \frac{1}{\sqrt{n}} \right) \left( \sigma + C \left( p^{1-\min\{\delta_0, C_1, c_0 C_0^2\}} + c' \exp(-cn) \right) (\log p)^2 \right), \end{aligned}$$

where the first inequality follows from (7.32) and second inequality follows from Lemma 4. If  $\frac{\log p}{n} \leq c$ , then  $\left(p^{1-\min\{\delta_0, C_1, c_0 C_0^2\}} + c' \exp(-cn)\right) (\log p)^2 \rightarrow 0$ , and hence  $\mathbb{E}_\theta L(\text{CI}_\alpha^S(\xi^\top \beta, Z)) \leq C \|\xi\|_2 \left(k \frac{\log p}{n} + \frac{1}{\sqrt{n}}\right) M_2$ .

**Acknowledgements.** The authors thank Zhao Ren for the discussion on the confidence intervals for linear functionals with sparse loadings.

## SUPPLEMENTARY MATERIAL

### Supplement to “Confidence Intervals for High-Dimensional Linear Regression: Minimax Rates and Adaptivity”.

(<http://www-stat.wharton.upenn.edu/~tcai/paper/CI-Reg-Supplement.pdf>).

Detailed proofs of the adaptivity lower bound and minimax upper bound for confidence intervals of the linear functional  $\xi^\top \beta$  with a dense loading  $\xi$  are given. The minimax rates and adaptivity of confidence intervals of the linear functional  $\xi^\top \beta$  are established when there is prior knowledge that  $\Omega = I$  and  $\sigma = \sigma_0$ . Extra propositions and technical lemmas are also proved in the supplement.

### References.

- [1] Alexandre Belloni, Victor Chernozhukov, and Lie Wang. Square-root lasso: pivotal recovery of sparse signals via conic programming. *Biometrika*, 98(4):791–806, 2011.
- [2] Peter J Bickel, Ya’acov Ritov, and Alexandre B Tsybakov. Simultaneous analysis of lasso and dantzig selector. *The Annals of Statistics*, 37(4):1705–1732, 2009.
- [3] T Tony Cai and Zijian Guo. Supplement to “confidence intervals for high-dimensional linear regression: Minimax rates and adaptivity”. 2015.
- [4] T Tony Cai, Weidong Liu, and Xi Luo. A constrained  $\ell_1$  minimization approach to sparse precision matrix estimation. *Journal of the American Statistical Association*, 106(494):594–607, 2011.
- [5] T Tony Cai and Mark G Low. Minimax estimation of linear functionals over non-convex parameter spaces. *The Annals of statistics*, 32(2):552–576, 2004.
- [6] T Tony Cai and Mark G Low. An adaptation theory for nonparametric confidence intervals. *The Annals of statistics*, 32(5):1805–1840, 2005.
- [7] T Tony Cai and Mark G Low. On adaptive estimation of linear functionals. *The Annals of Statistics*, 33(5):2311–2343, 2005.
- [8] T Tony Cai and Mark G Low. Adaptive confidence balls. *The Annals of Statistics*, 34(1):202–228, 2006.
- [9] T Tony Cai, Mark G Low, and Zongming Ma. Adaptive confidence bands for nonparametric regression functions. *Journal of the American Statistical Association*, 109:1054–1070, 2014.
- [10] T Tony Cai and Harrison H Zhou. Optimal rates of convergence for sparse covariance matrix estimation. *The Annals of Statistics*, 40(5):2389–2420, 2012.
- [11] Emmanuel Candès and Terence Tao. The dantzig selector: statistical estimation when  $p$  is much larger than  $n$ . *The Annals of Statistics*, 35(6):2313–2351, 2007.



- [12] Olivier Collier, Latitia Comminges, and Alexandre B. Tsybakov. Minimax estimation of linear and quadratic functionals on sparsity classes. *arXiv preprint arXiv:1502.00665*, 2015.
- [13] Marc Hoffmann and Richard Nickl. On adaptive inference and confidence bands. *The Annals of Statistics*, 39(5):2383–2409, 2011.
- [14] Adel Javanmard and Alessandro Montanari. Hypothesis testing in high-dimensional regression under the gaussian random design model: Asymptotic theory. *Information Theory, IEEE Transactions on*, 60(10):6522–6554, 2014.
- [15] Adel Javanmard and Andrea Montanari. Confidence intervals and hypothesis testing for high-dimensional regression. *The Journal of Machine Learning Research*, 15(1):2869–2909, 2014.
- [16] Adel Javanmard and Andrea Montanari. De-biasing the lasso: Optimal sample size for gaussian designs. *arXiv preprint arXiv:1508.02757*, 2015.
- [17] Richard Nickl and Sara van de Geer. Confidence sets in sparse regression. *The Annals of Statistics*, 41(6):2852–2876, 2013.
- [18] Garvesh Raskutti, Martin J Wainwright, and Bin Yu. Restricted eigenvalue properties for correlated gaussian designs. *The Journal of Machine Learning Research*, 11:2241–2259, 2010.
- [19] Zhao Ren, Tingni Sun, Cun-Hui Zhang, and Harrison H Zhou. Asymptotic normality and optimalities in estimation of large gaussian graphical model. *arXiv preprint arXiv:1309.6024*, 2013.
- [20] James Robins and Aad Van Der Vaart. Adaptive nonparametric confidence sets. *The Annals of Statistics*, 34(1):229–253, 2006.
- [21] Tingni Sun and Cun-Hui Zhang. Scaled sparse linear regression. *Biometrika*, 101(2):269–284, 2012.
- [22] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1):267–288, 1996.
- [23] Sara van de Geer, Peter Bühlmann, Yaacov Ritov, and Ruben Dezeure. On asymptotically optimal confidence regions and tests for high-dimensional models. *The Annals of Statistics*, 42(3):1166–1202, 2014.
- [24] Nicolas Verzelen. Minimax risks for sparse regressions: Ultra-high dimensional phenomena. *Electronic Journal of Statistics*, 6:38–90, 2012.
- [25] Cun-Hui Zhang and Stephanie S Zhang. Confidence intervals for low dimensional parameters in high dimensional linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(1):217–242, 2014.

DEPARTMENT OF STATISTICS  
THE WHARTON SCHOOL  
UNIVERSITY OF PENNSYLVANIA  
PHILADELPHIA, PENNSYLVANIA 19104  
USA  
E-MAIL: [tcai@wharton.upenn.edu](mailto:tcai@wharton.upenn.edu)  
[zijguo@wharton.upenn.edu](mailto:zijguo@wharton.upenn.edu)  
URL: <http://www-stat.wharton.upenn.edu/~tcai/>