

## GEE ANALYSIS OF CLUSTERED BINARY DATA WITH DIVERGING NUMBER OF COVARIATES

BY LAN WANG<sup>1</sup>

*University of Minnesota*

Clustered binary data with a large number of covariates have become increasingly common in many scientific disciplines. This paper develops an asymptotic theory for generalized estimating equations (GEE) analysis of clustered binary data when the number of covariates grows to infinity with the number of clusters. In this “large  $n$ , diverging  $p$ ” framework, we provide appropriate regularity conditions and establish the existence, consistency and asymptotic normality of the GEE estimator. Furthermore, we prove that the sandwich variance formula remains valid. Even when the working correlation matrix is misspecified, the use of the sandwich variance formula leads to an asymptotically valid confidence interval and Wald test for an estimable linear combination of the unknown parameters. The accuracy of the asymptotic approximation is examined via numerical simulations. We also discuss the “diverging  $p$ ” asymptotic theory for general GEE. The results in this paper extend the recent elegant work of Xie and Yang [*Ann. Statist.* **31** (2003) 310–347] and Balan and Schiopu-Kratina [*Ann. Statist.* **32** (2005) 522–541] in the “fixed  $p$ ” setting.

**1. Introduction.** A fundamental problem in statistical analysis is to characterize the effects of a set of covariates  $X_1, \dots, X_p$  on a response variable  $Y$  based on a sample of size  $n$ . Recently, there has been considerable interest in investigating this problem in the so-called “large  $n$ , diverging  $p$ ” asymptotic framework, where the dimension of the covariates increases to infinity with the sample size. This setup allows statisticians to adopt a more complex statistical model as more abundant data become available, and thus to reduce the modeling bias.

The “large  $n$ , diverging  $p$ ” framework can be traced back to the earlier pioneering work on M-estimators with a diverging number of parameter; see Huber (1973), Portnoy (1984, 1985, 1988), Mammen (1989), Welsh (1989), Bai and Wu (1994), He and Shao (2000) and the references therein. With the advent of high-dimensional data in many scientific areas, statistical theory developed in this new framework has become crucial for guiding practical data analysis with high-dimensional covariates, which relies heavily on asymptotic theory to justify its validity. By allowing the covariates’ dimension to increase with the sample size, Fan

---

Received December 2009; revised July 2010.

<sup>1</sup>Supported by NSF Grant DMS-1007603.

*AMS 2000 subject classifications.* Primary 62F12; secondary 62J12.

*Key words and phrases.* Clustered binary data, generalized estimating equations (GEE), high-dimensional covariates, sandwich variance formula.

and Peng (2004) studied nonconcave penalized likelihood; Lam and Fan (2008) investigated profile-kernel likelihood inference with generalized varying coefficient partially linear models; Huang, Horowitz and Ma (2008) explored bridge estimators in linear regression; Hjort, McKeague and Van Keilegom (2009) and Chen, Peng and Qin (2009) studied the effects of data dimension on empirical likelihood; Zou and Zhang (2009) studied the adaptive elastic net, Zhu and Zhu (2009) investigated parameter estimation in a semiparametric regression model with highly correlated predictors. In the aforementioned literature, the number of covariates  $p$  grows to infinity at a polynomial rate  $o(n^\alpha)$  for some  $0 < \alpha < 1$ . In particular, most of these papers provide necessary conditions under which classical asymptotic theories remain valid for  $\alpha$  in the range  $[\frac{1}{3}, \frac{1}{2}]$ .

A different line of research considers the case where  $p$  can be much larger than  $n$  and even grow at an exponential rate of  $n$ , in which case the sparsity assumption and other more stringent regularity conditions are generally required to investigate the large-sample properties. Furthermore, it is worth noting that much work has also been devoted to classification and multiple hypotheses testing problems with high-dimensional covariates, but these problems are different in nature from what is discussed in this paper. We refer to the review papers of Donoho (2000), Fan and Li (2006) and Fan and Lv (2010) for more comprehensive references on high-dimensional data analysis.

When the research focus is on modeling the relationship between  $Y$  and a high-dimensional vector of covariates, the existing literature in the “large  $n$ , diverging  $p$ ” setting has been largely restricted to independent data. In many modern data sets, in addition to the large dimensionality of covariates, complexity also arises when the responses are correlated due to repeated measures or clustered design. One representative example is the Framingham Heart Study, where the researchers are interested in linking common risk factors to the occurrence of cardiovascular diseases. In this study, many variables, such as age, smoking status, cholesterol level and blood pressure, were recorded for the participants during their clinic visits over the years to describe their physical characteristics and lifestyles. Another example is the Chicago Longitudinal Study in social science, which investigated the educational and social development of about 1500 low income, minority youths in the Chicago area. The study collected a large amount of information on many variables that measure children’s early antisocial behavior, individual-level attributes of the child, family attributes and social characteristics of both the child and the family, among others. In some other examples of clustered data, the number of variables measured for each individual or experimental unit may not be many, but when one considers various interaction effects, the actual number of predictors in the statistical model can still be large and better fits the “large  $p$ ” setup.

The intrinsic complexity of clustered data raises challenging issues for statistical analysis, especially for correlated non-Gaussian data where it is difficult to specify the full likelihood. In this paper, we establish the asymptotic properties of

generalized estimating equations (GEE), a semiparametric procedure widely used in practice for clustered data analysis, while allowing the covariate dimension to grow to infinity with the sample size.

The GEE procedure was introduced in a seminal paper of Liang and Zeger (1986) as a useful extension of generalized linear models [McCullagh and Nelder (1989)] to correlated data. Instead of specifying the full likelihood, it only requires the knowledge of the first two marginal moments and a working correlation matrix. Thus, it is particularly effective for modeling clustered binary or count data. A key advantage of the GEE approach is that it yields a consistent estimator (in the classical “large  $n$ , fixed  $p$ ” setup), even if the working correlation structure is misspecified. The GEE estimator is also asymptotically efficient if the correlation structure is indeed correctly specified. The original paper of Liang and Zeger focused mostly on the methodology development. Li (1997) adopted a min-max approach to study the consistency of GEE. A more complete and systematic large-sample theory for GEE, including consistency and asymptotic normality, was elegantly established by Xie and Yang (2003). Balan and Schiopu-Kratina (2005) also rigorously studied a closely related pseudo-likelihood framework for GEE. However, these papers all assume that  $p$  is fixed and that the number of clusters  $n$  goes to infinity. Xie and Yang (2003) also considered the case where the cluster size (number of observations within each cluster) is itself large, which corresponds to a large number of time points in the longitudinal setting.

This paper examines the effect of high-dimensional covariates on the GEE estimator in the “large  $n$ , diverging  $p$ ” setup, where  $p = p_n$  is a function of the sample size  $n$ . We focus on clustered binary data because binary response (e.g., disease status) is ubiquitous in many scientific applications and because of the relative transparency of technical derivation. We also discuss the related theory for general GEE in Section 5.1 The main technical challenges come from the high dimensionality of the covariates, the dependence among observations within each cluster and the nuisance parameters in the working correlation matrix. We provide a self-contained derivation and extend earlier theory in the literature on M-estimation with a large number of parameters, which is not tailored for clustered data and generally has not considered nuisance parameters.

We aim to answer the following essential questions. To what extent can the asymptotic results derived in the classical asymptotic framework for GEE still be deemed trustworthy when the number of covariates is large? How large can  $p_n$  be (relative to  $n$ )? The main findings in this paper reveal that under reasonable conditions, the GEE estimator  $\hat{\beta}_n$  is  $\sqrt{p_n/n}$ -consistent when  $p_n^2/n \rightarrow 0$  and that an arbitrary linear combination  $\alpha_n^T (\hat{\beta}_n - \beta_{n0})$  is asymptotically normal when  $p_n^3/n \rightarrow 0$ , where  $\beta_{n0}$  is the true parameter value. These findings resonate with those in the literature for independent data in the “large  $p$ ” setting. Moreover, we also verify that the desirable robustness property against working correlation matrix misspecification still holds and that both the sandwich variance formula and the large-sample

Wald test still remain valid in this new context. Understanding these fundamental questions is essential to justifying asymptotic statistical inference based on GEE for analyzing real-world clustered data containing many covariates, such as the validity of the confidence intervals provided by the GEE package in R, SAS and other statistical software packages.

The rest of the paper is organized as follows. In Section 2, we provide a brief review of the GEE procedure for analyzing clustered binary data. Section 3 establishes the consistency and asymptotic normality of the GEE estimator, the consistency of the sandwich variance formula and the validity of the large-sample Wald test in the “large  $n$ , diverging  $p$ ” framework. Section 4 examines the asymptotic results via numerical simulations. Section 5 discusses general GEE and related problems.

**2. Generalized estimating equations.** For the  $j$ th observation of the  $i$ th cluster, we observe a binary response variable  $Y_{ij}$  and a  $p_n$ -dimensional vector of covariates  $\mathbf{X}_{ij}$ ,  $i = 1, \dots, n$  and  $j = 1, \dots, m_i$ . Observations from different clusters are independent, but those from the same clusters are correlated. Let  $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{im_i})^T$  denote the vector of responses for the  $i$ th cluster and let  $\mathbf{X}_i = (\mathbf{X}_{i1}, \dots, \mathbf{X}_{im_i})^T$  be the associated  $m_i \times p_n$  matrix of covariates.

The marginal regression approach of GEE assumes that  $E(Y_{ij}|\mathbf{X}_{ij}) = \pi_{ij}$  and  $\text{Var}(Y_{ij}|\mathbf{X}_{ij}) = \pi_{ij}(1 - \pi_{ij})$ , where a dispersion parameter may be added in the marginal variance function if overdispersion is suspected to be present. Furthermore, it relates the covariates to the marginal mean by specifying that

$$(2.1) \quad \text{logit}(\pi_{ij}) = \mathbf{X}_{ij}^T \boldsymbol{\beta}_n,$$

where  $\text{logit}(\pi_{ij}) = \log\left(\frac{\pi_{ij}}{1-\pi_{ij}}\right)$  is the link function and  $\boldsymbol{\beta}_n$  is a  $p_n$ -dimensional vector of parameters. The true unknown parameter value is denoted by  $\boldsymbol{\beta}_{n0}$ .

Let  $\boldsymbol{\pi}_i(\boldsymbol{\beta}_n) = (\pi_{i1}(\boldsymbol{\beta}_n), \dots, \pi_{im_i}(\boldsymbol{\beta}_n))^T$ , where  $\pi_{ij}(\boldsymbol{\beta}_n) = \exp(\mathbf{X}_{ij}^T \boldsymbol{\beta}_n) / [1 + \exp(\mathbf{X}_{ij}^T \boldsymbol{\beta}_n)]$ . Further, let  $\mathbf{A}_i(\boldsymbol{\beta}_n)$  be the  $m_i \times m_i$  diagonal matrix with the  $j$ th diagonal element  $\mathbf{A}_{ij}(\boldsymbol{\beta}_n) = \pi_{ij}(\boldsymbol{\beta}_n)(1 - \pi_{ij}(\boldsymbol{\beta}_n))$ ,  $j = 1, \dots, m_i$ . In what follows, we assume  $m_i = m < \infty$ , for simplicity. Liang and Zeger (1986) suggested to estimate  $\boldsymbol{\beta}_{n0}$  by solving the following generalized estimating equation in  $\boldsymbol{\beta}_n$ :

$$(2.2) \quad \sum_{i=1}^n \mathbf{X}_i^T \mathbf{A}_i(\boldsymbol{\beta}_n) \mathbf{V}_i^{-1} (\mathbf{Y}_i - \boldsymbol{\pi}_i(\boldsymbol{\beta}_n)) = 0,$$

where  $\mathbf{V}_i$  is a working covariance matrix.

**3. Asymptotic properties when  $p_n \rightarrow \infty$ .**

3.1. *GEE estimator with estimated working correlation matrix.* In applications, the true correlation matrix of  $\mathbf{Y}_i$ , denoted by  $\mathbf{R}_0$ , is unknown. The working covariance matrix is often specified via a working correlation matrix  $\mathbf{R}(\boldsymbol{\tau})$ :  $\mathbf{V}_i = \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_n) \mathbf{R}(\boldsymbol{\tau}) \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_n)$ , where  $\boldsymbol{\tau}$  is a finite-dimensional parameter. Com-

monly used working correlation structures include AR-1, compound symmetry and unstructured working correlation, among others. Note that, in practice, the working correlation matrix is chosen to be independent of the covariates, for simplicity. However, for correlated non-normal data, the range of correlation generally depends on the univariate marginals. Thus,  $\mathbf{R}(\boldsymbol{\tau})$  should be understood as a weight matrix [Chaganty and Joe (2004)]. Chaganty and Joe demonstrated that GEE with an appropriately chosen working correlation matrix does have good efficiency when compared with a proper likelihood model.

Given a working correlation structure,  $\boldsymbol{\tau}$  is often estimated using a residual-based moment method, which requires an initial consistent estimator of  $\boldsymbol{\beta}_{n0}$ . We use  $\widehat{\mathbf{R}}$  to denote the resulting estimated working correlation matrix, with the subscript “ $n$ ” suppressed. Following (2.2), we formally define the *GEE estimator*  $\widehat{\boldsymbol{\beta}}_n$  as the solution of

$$(3.1) \quad \mathbf{S}_n(\boldsymbol{\beta}_n) = \sum_{i=1}^n \mathbf{X}_i^T \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_n) \widehat{\mathbf{R}}^{-1} \mathbf{A}_i^{-1/2}(\boldsymbol{\beta}_n) (\mathbf{Y}_i - \boldsymbol{\pi}_i(\boldsymbol{\beta}_n)) = 0.$$

To solve for  $\widehat{\boldsymbol{\beta}}_n$ , we can iterate between a modified Fisher scoring algorithm for  $\boldsymbol{\beta}_n$  and the moment estimation for  $\boldsymbol{\tau}$ . In the following, we provide examples of an initial consistent estimator and an estimated working correlation matrix.

**EXAMPLE 1** (Initial estimator for  $\boldsymbol{\beta}_{n0}$  when  $p_n \rightarrow \infty$ ). A simple way to obtain an initial estimator for  $\boldsymbol{\beta}_{n0}$  is to solve the generalized estimating equations under the working independence assumption

$$(3.2) \quad \widetilde{\mathbf{S}}_n(\boldsymbol{\beta}_n) = \sum_{i=1}^n \mathbf{X}_i^T (\mathbf{Y}_i - \boldsymbol{\pi}_i(\boldsymbol{\beta}_n)) = 0.$$

Under conditions (A1)–(A3) in Section 3.2, we can show that if  $p_n^2/n \rightarrow 0$  as  $n \rightarrow \infty$ , then the independence estimating equations in (3.2) have a root  $\widetilde{\boldsymbol{\beta}}_n$  such that

$$(3.3) \quad \|\widetilde{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_{n0}\| = O_p(\sqrt{p_n/n}),$$

where  $\|\cdot\|$  denotes the Euclidean norm of a vector. A detailed derivation of (3.3) is given in the [Appendix](#).

**EXAMPLE 2** (Estimated working correlation matrix when  $p_n \rightarrow \infty$ ). In Balan and Schiopu-Kratina (2005), it was suggested to use

$$\widehat{\mathbf{R}} = \frac{1}{n} \sum_{i=1}^n \mathbf{A}_i^{-1/2}(\widetilde{\boldsymbol{\beta}}_n) (\mathbf{Y}_i - \boldsymbol{\pi}_i(\widetilde{\boldsymbol{\beta}}_n)) (\mathbf{Y}_i - \boldsymbol{\pi}_i(\widetilde{\boldsymbol{\beta}}_n))^T \mathbf{A}_i^{-1/2}(\widetilde{\boldsymbol{\beta}}_n),$$

where  $\widetilde{\boldsymbol{\beta}}_n$  is a preliminary  $\sqrt{n/p_n}$ -consistent estimator of  $\boldsymbol{\beta}_{n0}$ , such as the one discussed in Example 1. This provides a moment estimator of the unstructured

working correlation matrix. Assuming conditions (A1)–(A3) of Section 3.2, we can prove that if  $p_n^2/n \rightarrow 0$  as  $n \rightarrow \infty$ , then

$$(3.4) \quad \|\widehat{\mathbf{R}}^{-1} - \mathbf{R}_0^{-1}\| = O_p(\sqrt{p_n/n}),$$

where  $\mathbf{R}_0$  denotes the true common correlation matrix. Here, and throughout the paper, for a matrix  $\mathbf{B}$ ,  $\|\mathbf{B}\| = [\text{Tr}(\mathbf{B}\mathbf{B}^T)]^{1/2}$  denotes its Frobenius norm. A detailed derivation of (3.4) is given in the supplementary article [Wang (2010)].

3.2. *Existence and consistency.* In Fan and Peng (2004), Lam and Fan (2008) and Huang, Horowitz and Ma (2008), the estimator is defined as the minimizer of a certain objective function. We use alternative techniques here to establish the existence and consistency of the GEE estimator, which involve the roots of estimating equations. The approach we adopt here is also different from that of Xie and Yang (2003) and Balan and Schiopu-Kratina (2005), both of which rely on properties of injective functions.

We directly verify the following condition:  $\forall \varepsilon > 0$ , there exists a constant  $\Delta > 0$  such that for all  $n$  sufficiently large,

$$(3.5) \quad P\left(\sup_{\|\beta_n - \beta_{n0}\| = \Delta\sqrt{p_n/n}} (\beta_n - \beta_{n0})^T \mathbf{S}_n(\beta_n) < 0\right) \geq 1 - \varepsilon.$$

Condition (3.5) is sufficient to ensure the existence of a sequence of roots  $\widehat{\beta}_n$  of the equation  $\mathbf{S}_n(\beta_n) = 0$  such that  $\|\widehat{\beta}_n - \beta_{n0}\| = O_p(\sqrt{p_n/n})$ . This approach follows from Theorem 6.3.4 of Ortega and Rheinboldt (1970). In Portnoy (1984), this technique was applied to establish the existence and consistency of an M-estimator for i.i.d. data; in a different setting, it was used by Wang et al. (2010) to study a partial linear single-index model. This leads to a more straightforward and elegant proof of weak consistency. On the other hand, the method relying on injective functions [Xie and Yang (2003); Balan and Schiopu-Kratina (2005)] can also be used to prove strong consistency.

To prove consistency and asymptotic normality, we need the following general regularity conditions:

(A1)  $\sup_{i,j} \|\mathbf{X}_{ij}\| = O(\sqrt{p_n})$ ;

(A2) the unknown parameter  $\beta_n$  belongs to a compact subset  $\mathcal{B} \subseteq R^{p_n}$ , the true parameter value  $\beta_{n0}$  lies in the interior of  $\mathcal{B}$  and there exist two positive constants,  $b_1$  and  $b_2$ , such that  $0 < b_1 \leq \pi_{ij}(\beta_{n0}) \leq b_2 < 1, \forall i, j$ ;

(A3) there exist two positive constants,  $b_3$  and  $b_4$ , such that

$$b_3 \leq \lambda_{\min}\left(n^{-1} \sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i\right) \leq \lambda_{\max}\left(n^{-1} \sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i\right) \leq b_4,$$

where  $\lambda_{\min}$  (resp.  $\lambda_{\max}$ ) denotes the minimum (resp. maximum) eigenvalue of a matrix;

(A4) the common true correlation matrix  $\mathbf{R}_0$  has eigenvalues bounded away from zero and  $+\infty$ ; the estimated working correlation matrix  $\widehat{\mathbf{R}}$  satisfies  $\|\widehat{\mathbf{R}}^{-1} - \overline{\mathbf{R}}^{-1}\| = O_p(\sqrt{p_n/n})$ , where  $\overline{\mathbf{R}}$  is a constant positive definite matrix with eigenvalues bounded away from zero and  $+\infty$ ; we do not require  $\overline{\mathbf{R}}$  to be the true correlation matrix  $\mathbf{R}_0$ .

REMARK 1. Condition (A1) is a common assumption in the literature on M-estimators with diverging dimension. For example, it is the same as assumption (3.9) of Portnoy (1985) and it is implied by conditions (C.9) and (C.10) of Welsh (1989). This condition holds almost surely under some weak moment conditions for  $X_{ij}$  from spherically symmetric distributions [see, e.g., the discussions in He and Shao (2000)]. When  $m = 1$  (i.e., each cluster has only one observation), condition (A3) is also popularly adopted in the literature on high-dimensional regression for independent data. It can be shown that condition (A3) is implied by the following slightly stronger condition: there exist two positive constants,  $c_1 \leq c_2$ , such that  $\forall 1 \leq j \leq m$ ,

$$c_1 \leq \lambda_{\min} \left( n^{-1} \sum_{i=1}^n \mathbf{X}_{ij} \mathbf{X}_{ij}^T \right) \leq \lambda_{\max} \left( n^{-1} \sum_{i=1}^n \mathbf{X}_{ij} \mathbf{X}_{ij}^T \right) \leq c_2.$$

Finally, condition (A4) is a direct extension of a similar assumption in the “fixed  $p$ ” case. Liang and Zeger (1986) assumes that the estimator of the working correlation matrix parameter  $\widehat{\boldsymbol{\tau}}$  satisfies  $\sqrt{n}(\widehat{\boldsymbol{\tau}} - \boldsymbol{\tau}_0) = O_p(1)$  for some  $\boldsymbol{\tau}_0$ . Assumption (C2) of Chen and Jin (2006) is of similar nature, while Xie and Yang (2003) assumes the nuisance parameter  $\boldsymbol{\tau}$  to be completely known. Note that Example 2 in Section 3.1 guarantees that (A4) is satisfied when a nonparametric moment estimator is used for the working correlation matrix, in which case  $\overline{\mathbf{R}} = \mathbf{R}_0$ .

We use notation similar to that in Xie and Yang (2003) and Balan and Schiopu-Kratina (2005). Consider the following estimating equation:

$$\overline{\mathbf{S}}_n(\boldsymbol{\beta}_n) = \sum_{i=1}^n \mathbf{X}_i^T \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_n) \overline{\mathbf{R}}^{-1} \mathbf{A}_i^{-1/2}(\boldsymbol{\beta}_n) (\mathbf{Y}_i - \boldsymbol{\pi}_i(\boldsymbol{\beta}_n)).$$

If we let  $\overline{\mathbf{M}}_n(\boldsymbol{\beta}_n)$  denote the covariance matrix of  $\overline{\mathbf{S}}_n(\boldsymbol{\beta}_n)$ , then

$$\overline{\mathbf{M}}_n(\boldsymbol{\beta}_n) = \sum_{i=1}^n \mathbf{X}_i^T \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_n) \overline{\mathbf{R}}^{-1} \mathbf{R}_0 \overline{\mathbf{R}}^{-1} \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_n) \mathbf{X}_i.$$

To prove the consistency, the essential idea is to approximate  $\mathbf{S}_n(\boldsymbol{\beta}_n)$  by  $\overline{\mathbf{S}}_n(\boldsymbol{\beta}_n)$ , whose moments are easier to evaluate. Lemma 3.1 below establishes the accuracy of this approximation, which also plays an important role in deriving the asymptotic normality in Section 3.3.

LEMMA 3.1. Assume conditions (A1)–(A4). If  $n^{-1}p_n^2 = o(1)$ , then

$$\|\mathbf{S}_n(\boldsymbol{\beta}_{n0}) - \bar{\mathbf{S}}_n(\boldsymbol{\beta}_{n0})\| = O_p(p_n).$$

To facilitate the Taylor expansion of the estimating function  $\mathbf{S}_n(\boldsymbol{\beta}_n)$ , we also use  $\bar{\mathbf{D}}_n(\boldsymbol{\beta}_n) = -\frac{\partial}{\partial \boldsymbol{\beta}_n^T} \bar{\mathbf{S}}_n(\boldsymbol{\beta}_n)$  to approximate the negative gradient function  $\mathbf{D}_n(\boldsymbol{\beta}_n) = -\frac{\partial}{\partial \boldsymbol{\beta}_n^T} \mathbf{S}_n(\boldsymbol{\beta}_n)$ . Lemma 3.2 below provides a useful representation of  $\bar{\mathbf{D}}_n(\boldsymbol{\beta}_n)$ , based on which, Lemma 3.3 establishes the approximation of gradient functions.

LEMMA 3.2.

$$(3.6) \quad \bar{\mathbf{D}}_n(\boldsymbol{\beta}_n) = \bar{\mathbf{H}}_n(\boldsymbol{\beta}_n) + \bar{\mathbf{E}}_n(\boldsymbol{\beta}_n) + \bar{\mathbf{G}}_n(\boldsymbol{\beta}_n),$$

where

$$\begin{aligned} \bar{\mathbf{H}}_n(\boldsymbol{\beta}_n) &= \sum_{i=1}^n \mathbf{X}_i^T \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_n) \bar{\mathbf{R}}^{-1} \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_n) \mathbf{X}_i, \\ \bar{\mathbf{E}}_n(\boldsymbol{\beta}_n) &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m (1 - 2\pi_{ij}(\boldsymbol{\beta}_n)) \varepsilon_{ij}(\boldsymbol{\beta}_n) \mathbf{X}_i^T \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_n) \bar{\mathbf{R}}^{-1} \mathbf{e}_j \mathbf{e}_j^T \mathbf{X}_i, \\ \bar{\mathbf{G}}_n(\boldsymbol{\beta}_n) &= -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m (1 - 2\pi_{ij}(\boldsymbol{\beta}_n)) \mathbf{A}_{ij}^{1/2}(\boldsymbol{\beta}_n) \mathbf{X}_i \mathbf{X}_i^T \mathbf{e}_j \bar{\mathbf{R}}^{-1} \boldsymbol{\varepsilon}_i(\boldsymbol{\beta}_n), \end{aligned}$$

where  $\varepsilon_{ij}(\boldsymbol{\beta}_n) = \mathbf{A}_{ij}^{-1/2}(\boldsymbol{\beta}_n)(Y_{ij} - \pi_{ij}(\boldsymbol{\beta}_n))$ ,  $\boldsymbol{\varepsilon}_i(\boldsymbol{\beta}_n) = \mathbf{A}_i^{-1/2}(\boldsymbol{\beta}_n)(\mathbf{Y}_i - \boldsymbol{\pi}(\boldsymbol{\beta}_n))$  and  $\mathbf{e}_j$  denotes a unit vector of length  $m$  whose  $j$ th entry is 1 and all other entries of which are 0.

LEMMA 3.3. Assume conditions (A1)–(A4). If  $n^{-1}p_n^2 = o(1)$ , then  $\forall \Delta > 0$ , for  $\mathbf{b}_n \in R^{p_n}$ , we have

$$\sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} \sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T [\mathbf{D}_n(\boldsymbol{\beta}_n) - \bar{\mathbf{D}}_n(\boldsymbol{\beta}_n)] \mathbf{b}_n| = O_p(\sqrt{np_n}).$$

REMARK 2. The matrix  $\mathbf{D}_n(\boldsymbol{\beta}_n) - \bar{\mathbf{D}}_n(\boldsymbol{\beta}_n)$  is symmetric. The above lemma immediately implies that

$$\begin{aligned} \sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} |\lambda_{\min}[\mathbf{D}_n(\boldsymbol{\beta}_n) - \bar{\mathbf{D}}_n(\boldsymbol{\beta}_n)]| &= O_p(\sqrt{np_n}), \\ \sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} |\lambda_{\max}[\mathbf{D}_n(\boldsymbol{\beta}_n) - \bar{\mathbf{D}}_n(\boldsymbol{\beta}_n)]| &= O_p(\sqrt{np_n}). \end{aligned}$$

Furthermore, we can use the leading term  $\bar{\mathbf{H}}_n(\boldsymbol{\beta}_n)$  in (3.6) to approximate the negative gradient function  $\bar{\mathbf{D}}_n(\boldsymbol{\beta}_n)$ . This result is given by Lemma 3.4 below. Lemma 3.5 further establishes an equicontinuity result for  $\bar{\mathbf{H}}_n(\boldsymbol{\beta}_n)$ .



LEMMA 3.4. Assume conditions (A1)–(A4). If  $n^{-1}p_n^2 = o(1)$ , then  $\forall \Delta > 0$ , for  $\mathbf{b}_n \in R^{p_n}$ , we have

$$\sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} \sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T [\bar{\mathbf{D}}_n(\boldsymbol{\beta}_n) - \bar{\mathbf{H}}_n(\boldsymbol{\beta}_n)] \mathbf{b}_n| = O_p(\sqrt{n} p_n).$$

LEMMA 3.5. Assume conditions (A1)–(A4). If  $n^{-1}p_n^2 = o(1)$ , then  $\forall \Delta > 0$ , for  $\mathbf{b}_n \in R^{p_n}$ , we have

$$\sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} \sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T [\bar{\mathbf{H}}_n(\boldsymbol{\beta}_n) - \bar{\mathbf{H}}_n(\boldsymbol{\beta}_{n0})] \mathbf{b}_n| = O_p(\sqrt{n} p_n).$$

The proofs of Lemmas 3.1–3.4 are given in the Appendix; the proof of Lemma 3.5 is given in the supplementary article [Wang (2010)]. The following theorem ensures the existence and consistency of the GEE estimator when  $p_n \rightarrow \infty$ .

THEOREM 3.6 (Existence and consistency). Assume conditions (A1)–(A4) and that  $n^{-1}p_n^2 = o(1)$ . Then,  $\mathbf{S}_n(\boldsymbol{\beta}_n) = 0$  has a root  $\hat{\boldsymbol{\beta}}_n$  such that

$$\|\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_{n0}\| = O_p(\sqrt{p_n/n}).$$

PROOF. We will prove that (3.5) holds. This requires us to evaluate the sign of  $(\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0})^T \mathbf{S}_n(\boldsymbol{\beta}_n)$  on  $\{\boldsymbol{\beta}_n : \|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| = \Delta \sqrt{p_n/n}\}$ . Note that

$$\begin{aligned} &(\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0})^T \mathbf{S}_n(\boldsymbol{\beta}_n) \\ &= (\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0})^T \mathbf{S}_n(\boldsymbol{\beta}_{n0}) - (\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0})^T \mathbf{D}_n(\boldsymbol{\beta}_n^*)(\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}) \\ &\triangleq I_{n1} + I_{n2}, \end{aligned}$$

where  $\boldsymbol{\beta}_n^*$  lies between  $\boldsymbol{\beta}_n$  and  $\boldsymbol{\beta}_{n0}$ , that is,  $\boldsymbol{\beta}_n^* = t\boldsymbol{\beta}_n + (1-t)\boldsymbol{\beta}_{n0}$  for some  $0 < t < 1$ . Next, we write

$$\begin{aligned} I_{n1} &= (\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0})^T \bar{\mathbf{S}}_n(\boldsymbol{\beta}_{n0}) + (\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0})^T [\mathbf{S}_n(\boldsymbol{\beta}_{n0}) - \bar{\mathbf{S}}_n(\boldsymbol{\beta}_{n0})] \\ &\triangleq I_{n11} + I_{n12}. \end{aligned}$$

We have  $|I_{n11}| \leq \|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \cdot \|\bar{\mathbf{S}}_n(\boldsymbol{\beta}_{n0})\| = \Delta \sqrt{p_n/n} \|\bar{\mathbf{S}}_n(\boldsymbol{\beta}_{n0})\|$  by the Cauchy–Schwarz inequality. Furthermore,

$$\begin{aligned} &E[\|\bar{\mathbf{S}}_n(\boldsymbol{\beta}_{n0})\|^2] \\ &= E \left\{ \sum_{i=1}^n \boldsymbol{\varepsilon}_i^T(\boldsymbol{\beta}_{n0}) \bar{\mathbf{R}}^{-1} \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0}) \mathbf{X}_i \mathbf{X}_i^T \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0}) \bar{\mathbf{R}}^{-1} \boldsymbol{\varepsilon}_i(\boldsymbol{\beta}_{n0}) \right\} \\ &\leq \sum_{i=1}^n \lambda_{\max}(\mathbf{X}_i \mathbf{X}_i^T) \lambda_{\max}(\mathbf{A}_i(\boldsymbol{\beta}_{n0})) \lambda_{\max}(\bar{\mathbf{R}}^{-2}) E[\boldsymbol{\varepsilon}_i^T(\boldsymbol{\beta}_{n0}) \boldsymbol{\varepsilon}_i(\boldsymbol{\beta}_{n0})] \\ &\leq C \operatorname{Tr} \left( \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T \right) = C \sum_{i=1}^n \sum_{j=1}^m \mathbf{X}_{ij}^T \mathbf{X}_{ij} = O(np_n). \end{aligned}$$

Here, and throughout the paper, we use  $C$  to denote a generic positive constant which may vary from line to line. Thus,  $\|\bar{\mathbf{S}}_n(\boldsymbol{\beta}_{n0})\| = O_p(\sqrt{np_n})$ . This implies that  $|I_{n11}| = \Delta O_p(p_n)$ . For  $I_{n12}$ , we have

$$|I_{n12}| \leq \|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \cdot \|\mathbf{S}_n(\boldsymbol{\beta}_{n0}) - \bar{\mathbf{S}}_n(\boldsymbol{\beta}_{n0})\| = \Delta\sqrt{p_n/n}O_p(p_n) = \Delta o_p(p_n),$$

by Lemma 3.1. Hence,  $|I_{n1}| = \Delta O_p(p_n)$ . In what follows, we evaluate  $I_{n2}$ :

$$\begin{aligned} I_{n2} &= -(\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0})^T \bar{\mathbf{D}}_n(\boldsymbol{\beta}_n^*)(\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}) \\ &\quad - (\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0})^T [\mathbf{D}_n(\boldsymbol{\beta}_n^*) - \bar{\mathbf{D}}_n(\boldsymbol{\beta}_n^*)](\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}) \\ &\triangleq I_{n21} + I_{n22}. \end{aligned}$$

First, note that

$$\begin{aligned} |I_{n22}| &\leq \max(|\lambda_{\max}(\mathbf{D}_n(\boldsymbol{\beta}_n^*) - \bar{\mathbf{D}}_n(\boldsymbol{\beta}_n^*))|, |\lambda_{\min}(\mathbf{D}_n(\boldsymbol{\beta}_n^*) - \bar{\mathbf{D}}_n(\boldsymbol{\beta}_n^*))|) \\ &\quad \times \|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\|^2 \\ &= O_p(\sqrt{np_n})\Delta^2 \frac{p_n}{n} = \Delta^2 o_p(p_n), \end{aligned}$$

by Lemma 3.3. On the other hand,

$$\begin{aligned} I_{n21} &= -(\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0})^T \bar{\mathbf{H}}_n(\boldsymbol{\beta}_{n0})(\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}) \\ &\quad - (\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0})^T [\bar{\mathbf{H}}_n(\boldsymbol{\beta}_n^*) - \bar{\mathbf{H}}_n(\boldsymbol{\beta}_{n0})](\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}) \\ &\quad - (\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0})^T [\bar{\mathbf{D}}_n(\boldsymbol{\beta}_n^*) - \bar{\mathbf{H}}_n(\boldsymbol{\beta}_n^*)](\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}) \\ &\triangleq I_{n21}^a + I_{n21}^b + I_{n21}^c. \end{aligned}$$

From Lemma 3.5, we have  $I_{n21}^b = \Delta^2 o_p(p_n)$ ; from Lemma 3.4, we have  $I_{n21}^c = \Delta^2 o_p(p_n)$ . Finally, we evaluate  $I_{n21}^a$ . We have

$$\begin{aligned} I_{n21}^a &= -(\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0})^T \left[ \sum_{i=1}^n \mathbf{X}_i^T \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0}) \bar{\mathbf{R}}^{-1} \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0}) \mathbf{X}_i \right] (\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}) \\ &\leq -\lambda_{\min}(\bar{\mathbf{R}}^{-1}) \min_i \lambda_{\min}(\mathbf{A}_i(\boldsymbol{\beta}_{n0})) \lambda_{\min} \left( \sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i \right) \|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\|^2 \\ &\leq -C\Delta^2 p_n, \end{aligned}$$

by (A3). Thus,  $(\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0})^T \mathbf{S}_n(\boldsymbol{\beta}_n)$  on  $\{\boldsymbol{\beta}_n : \|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| = \Delta\sqrt{p_n/n}\}$  is asymptotically dominated in probability by  $I_{n11} + I_{n21}^a$ , which is negative for  $\Delta$  large enough.  $\square$

3.3. *Asymptotic normality of the GEE estimator.* The asymptotic distribution of the GEE estimator  $\widehat{\beta}_n$  is closely related to that of the ideal estimating function  $\overline{\mathbf{S}}_n(\beta_{n0})$ . When appropriately normalized,  $\overline{\mathbf{S}}_n(\beta_{n0})$  has an asymptotic normal distribution, as shown by the following lemma.

LEMMA 3.7. *Assume conditions (A1)–(A4). If  $n^{-1}p_n^3 = o(1)$ , then  $\forall \alpha_n \in R^{p_n}$  such that  $\|\alpha_n\| = 1$ , we have*

$$\alpha_n^T \overline{\mathbf{M}}_n^{-1/2}(\beta_{n0}) \overline{\mathbf{S}}_n(\beta_{n0}) \rightarrow N(0, 1) \quad \text{in distribution.}$$

To prove Lemma 3.7, we write  $\alpha_n^T \overline{\mathbf{M}}_n^{-1/2}(\beta_{n0}) \overline{\mathbf{S}}_n(\beta_{n0})$  as a sum of independent random variables and then check the Lindeberg–Feller condition for the central limit theorem. The detailed proof is given in the [Appendix](#). The following theorem ensures the asymptotic normality of the GEE estimator when  $n^{-1}p_n^3 = o(1)$ .

THEOREM 3.8 (Asymptotic normality). *Assume conditions (A1)–(A4). If  $n^{-1}p_n^3 = o(1)$ , then  $\forall \alpha_n \in R^{p_n}$  such that  $\|\alpha_n\| = 1$ , we have*

$$\alpha_n^T \overline{\mathbf{M}}_n^{-1/2}(\beta_{n0}) \overline{\mathbf{H}}_n(\beta_{n0})(\widehat{\beta}_n - \beta_{n0}) \rightarrow N(0, 1)$$

*in distribution.*

PROOF. We have

$$\begin{aligned} & \alpha_n^T \overline{\mathbf{M}}_n^{-1/2}(\beta_{n0}) \overline{\mathbf{S}}_n(\beta_{n0}) \\ &= \alpha_n^T \overline{\mathbf{M}}_n^{-1/2}(\beta_{n0}) \mathbf{S}_n(\beta_{n0}) + \alpha_n^T \overline{\mathbf{M}}_n^{-1/2}(\beta_{n0}) [\overline{\mathbf{S}}_n(\beta_{n0}) - \mathbf{S}_n(\beta_{n0})] \\ &= \alpha_n^T \overline{\mathbf{M}}_n^{-1/2}(\beta_{n0}) \mathbf{D}_n(\beta_n^*)(\widehat{\beta}_n - \beta_{n0}) \\ & \quad + \alpha_n^T \overline{\mathbf{M}}_n^{-1/2}(\beta_{n0}) [\overline{\mathbf{S}}_n(\beta_{n0}) - \mathbf{S}_n(\beta_{n0})] \\ &= \alpha_n^T \overline{\mathbf{M}}_n^{-1/2}(\beta_{n0}) \overline{\mathbf{H}}_n(\beta_{n0})(\widehat{\beta}_n - \beta_{n0}) \\ & \quad + \alpha_n^T \overline{\mathbf{M}}_n^{-1/2}(\beta_{n0}) [\mathbf{D}_n(\beta_n^*) - \overline{\mathbf{H}}_n(\beta_{n0})](\widehat{\beta}_n - \beta_{n0}) \\ & \quad + \alpha_n^T \overline{\mathbf{M}}_n^{-1/2}(\beta_{n0}) [\overline{\mathbf{S}}_n(\beta_{n0}) - \mathbf{S}_n(\beta_{n0})], \end{aligned}$$

where, to obtain the second equality, we note that  $\mathbf{S}_n(\widehat{\beta}_n) = 0$  and thus, by a Taylor expansion,  $\mathbf{S}_n(\beta_{n0}) = \mathbf{D}_n(\beta_n^*)(\widehat{\beta}_n - \beta_{n0})$  for some  $\beta_n^*$  between  $\widehat{\beta}_n$  and  $\beta_{n0}$ . By Lemma 3.7,  $\alpha_n^T \overline{\mathbf{M}}_n^{-1/2}(\beta_{n0}) \overline{\mathbf{S}}_n(\beta_{n0}) \rightarrow N(0, 1)$ . Therefore, to prove the theorem, it is sufficient to verify that  $\forall \Delta > 0$ ,

$$\begin{aligned} & \sup_{\|\beta_n - \beta_{n0}\| \leq \Delta \sqrt{p_n/n}} |\alpha_n^T \overline{\mathbf{M}}_n^{-1/2}(\beta_{n0}) [\mathbf{D}_n(\beta_n) - \overline{\mathbf{H}}_n(\beta_{n0})](\widehat{\beta}_n - \beta_{n0})| \\ (3.7) \quad & = o_p(1) \end{aligned}$$

and

$$(3.8) \quad |\boldsymbol{\alpha}_n^T \bar{\mathbf{M}}_n^{-1/2}(\boldsymbol{\beta}_{n0}) [\bar{\mathbf{S}}_n(\boldsymbol{\beta}_{n0}) - \mathbf{S}_n(\boldsymbol{\beta}_{n0})]| = o_p(1).$$

We prove (3.8) first. Note that

$$\begin{aligned} & [ \boldsymbol{\alpha}_n^T \bar{\mathbf{M}}_n^{-1/2}(\boldsymbol{\beta}_{n0}) [\bar{\mathbf{S}}_n(\boldsymbol{\beta}_{n0}) - \mathbf{S}_n(\boldsymbol{\beta}_{n0})] ]^2 \\ &= \boldsymbol{\alpha}_n^T \bar{\mathbf{M}}_n^{-1/2}(\boldsymbol{\beta}_{n0}) [\bar{\mathbf{S}}_n(\boldsymbol{\beta}_{n0}) - \mathbf{S}_n(\boldsymbol{\beta}_{n0})] [\bar{\mathbf{S}}_n(\boldsymbol{\beta}_{n0}) - \mathbf{S}_n(\boldsymbol{\beta}_{n0})]^T \bar{\mathbf{M}}_n^{-1/2}(\boldsymbol{\beta}_{n0}) \boldsymbol{\alpha}_n \\ &\leq \lambda_{\max}(\bar{\mathbf{M}}_n^{-1}(\boldsymbol{\beta}_{n0})) \lambda_{\max}([\bar{\mathbf{S}}_n(\boldsymbol{\beta}_{n0}) - \mathbf{S}_n(\boldsymbol{\beta}_{n0})] [\bar{\mathbf{S}}_n(\boldsymbol{\beta}_{n0}) - \mathbf{S}_n(\boldsymbol{\beta}_{n0})]^T) \\ &\leq \frac{\|\bar{\mathbf{S}}_n(\boldsymbol{\beta}_{n0}) - \mathbf{S}_n(\boldsymbol{\beta}_{n0})\|^2}{\lambda_{\min}(\bar{\mathbf{M}}_n(\boldsymbol{\beta}_{n0}))} \\ &\leq \frac{\|\bar{\mathbf{S}}_n(\boldsymbol{\beta}_{n0}) - \mathbf{S}_n(\boldsymbol{\beta}_{n0})\|^2}{C \lambda_{\min}(\sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i)} \\ &= O_p(p_n^2/n) = o_p(1), \end{aligned}$$

by Lemma 3.1 and the fact that

$$(3.9) \quad \lambda_{\min}(\bar{\mathbf{M}}_n(\boldsymbol{\beta}_{n0})) \geq C \lambda_{\min}\left(\sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i\right).$$

A justification of (3.9) is given in the proof of Lemma 3.7 in the [Appendix](#). Thus, (3.8) holds. Next, we prove (3.7). We have

$$\begin{aligned} & \sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} |\boldsymbol{\alpha}_n^T \bar{\mathbf{M}}_n^{-1/2}(\boldsymbol{\beta}_{n0}) [\mathbf{D}_n(\boldsymbol{\beta}_n) - \bar{\mathbf{H}}_n(\boldsymbol{\beta}_{n0})] (\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_{n0})| \\ &\leq \sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} |\boldsymbol{\alpha}_n^T \bar{\mathbf{M}}_n^{-1/2}(\boldsymbol{\beta}_{n0}) [\mathbf{D}_n(\boldsymbol{\beta}_n) - \bar{\mathbf{D}}_n(\boldsymbol{\beta}_n)] (\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_{n0})| \\ &\quad + \sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} |\boldsymbol{\alpha}_n^T \bar{\mathbf{M}}_n^{-1/2}(\boldsymbol{\beta}_{n0}) [\bar{\mathbf{D}}_n(\boldsymbol{\beta}_n) - \bar{\mathbf{H}}_n(\boldsymbol{\beta}_n)] (\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_{n0})| \\ &\quad + \sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} |\boldsymbol{\alpha}_n^T \bar{\mathbf{M}}_n^{-1/2}(\boldsymbol{\beta}_{n0}) [\bar{\mathbf{H}}_n(\boldsymbol{\beta}_n) - \bar{\mathbf{H}}_n(\boldsymbol{\beta}_{n0})] (\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_{n0})| \\ &\triangleq I_{n1} + I_{n2} + I_{n3}. \end{aligned}$$

By the Cauchy–Schwarz inequality and Remark 2, we have

$$\begin{aligned} I_{n1} \leq & \sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} [ \boldsymbol{\alpha}_n^T \bar{\mathbf{M}}_n^{-1/2}(\boldsymbol{\beta}_{n0}) (\mathbf{D}_n(\boldsymbol{\beta}_n) - \bar{\mathbf{D}}_n(\boldsymbol{\beta}_n))^2 \\ & \times \bar{\mathbf{M}}_n^{-1/2}(\boldsymbol{\beta}_{n0}) \boldsymbol{\alpha}_n ]^{1/2} \|\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_{n0}\| \end{aligned}$$

$$\begin{aligned} &\leq \sup_{\|\beta_n - \beta_{n0}\| \leq \Delta \sqrt{p_n/n}} \max(|\lambda_{\min}(\mathbf{D}_n(\beta_n) - \overline{\mathbf{D}}_n(\beta_n))|, \\ &\quad |\lambda_{\max}(\mathbf{D}_n(\beta_n) - \overline{\mathbf{D}}_n(\beta_n))|) \\ &\quad \times \lambda_{\min}^{-1/2}(\overline{\mathbf{M}}_n(\beta_{n0})) O_p(p_n^{1/2} n^{-1/2}) \\ &= O_p(\sqrt{n} p_n) O(n^{-1/2}) O_p(n^{-1/2} p_n^{1/2}) = O_P(n^{-1/2} p_n^{3/2}) = o_p(1). \end{aligned}$$

Hence,  $I_{n1} = o_p(1)$ . By the same argument and Lemmas 3.4 and 3.5, we also have  $I_{n2} = o_p(1)$  and  $I_{n3} = o_p(1)$ . This proves (3.7).  $\square$

REMARK 3. Note that the condition  $n^{-1} p_n^3 = o(1)$  is the same as that of Huber (1973) for an M-estimator with independent data and diverging number of parameters. It is weaker than the condition  $n^{-1} p_n^5 = o(1)$  in Fan and Peng (2004) and Lam and Fan (2008) for asymptotic normality.

REMARK 4. Combining the result of Theorem 3.8 with the Cramér–Wold device, it is easy to see that for any  $l \times p_n$  matrix  $\mathbf{B}_n$  with  $l$  fixed and such that  $\mathbf{B}_n \mathbf{B}_n^T \rightarrow \mathbf{F}$ , a positive definite matrix, we have

$$\mathbf{B}_n \Sigma_n^{-1/2}(\beta_{n0})(\widehat{\beta}_n - \beta_{n0}) \rightarrow \mathbf{N}_l(\mathbf{0}, \mathbf{F}),$$

where

$$\Sigma_n = \overline{\mathbf{H}}_n^{-1}(\beta_{n0}) \overline{\mathbf{M}}_n(\beta_{n0}) \overline{\mathbf{H}}_n^{-1}(\beta_{n0}).$$

Now, take  $\mathbf{B}_n = (\mathbf{L}_n \Sigma_n \mathbf{L}_n^T)^{-1/2} \mathbf{L}_n \Sigma_n^{1/2}$ , where  $\mathbf{L}_n$  is an  $l \times p_n$  matrix such that  $\mathbf{L}_n \Sigma_n \mathbf{L}_n^T$  is invertible. Then,  $\mathbf{B}_n \mathbf{B}_n^T = \mathbf{I}_l$  and we have the following corollary which gives the asymptotic distribution of  $\mathbf{L}_n(\widehat{\beta}_n - \beta_{n0})$ .

COROLLARY 3.9. Under the same conditions as in Theorem 3.8, if  $n^{-1} p_n^3 = o(1)$ , then

$$(\mathbf{L}_n \Sigma_n \mathbf{L}_n^T)^{-1/2} \mathbf{L}_n(\widehat{\beta}_n - \beta_{n0}) \rightarrow \mathbf{N}_l(\mathbf{0}, \mathbf{I}_l)$$

in distribution.

3.4. Sandwich covariance formula and large-sample Wald test. Theorem 3.8 and Corollary 3.9 suggest that the covariance matrix of  $\widehat{\beta}_n$  is approximately  $\Sigma_n$ . To estimate  $\Sigma_n$ , Liang and Zeger (1986) proposed, in the “fixed  $p$ ” setup, the following well-known sandwich covariance matrix estimator:

$$\widehat{\Sigma}_n = \mathbf{H}_n^{-1}(\widehat{\beta}_n) \widehat{\mathbf{M}}_n(\widehat{\beta}_n) \mathbf{H}_n^{-1}(\widehat{\beta}_n),$$

where  $\mathbf{H}_n(\beta_n)$  is defined similarly as  $\overline{\mathbf{H}}_n(\beta_n)$ , but with  $\overline{\mathbf{R}}$  replaced by  $\widehat{\mathbf{R}}$ ;  $\widehat{\mathbf{M}}_n(\beta_n)$  is defined similarly as  $\overline{\mathbf{M}}_n(\beta_n)$ , except that  $\overline{\mathbf{R}}$  is replaced by  $\widehat{\mathbf{R}}$  and the unknown true correlation matrix  $\mathbf{R}_0$  is replaced by  $\boldsymbol{\varepsilon}_i(\beta_n) \boldsymbol{\varepsilon}_i^T(\beta_n)$ , with  $\boldsymbol{\varepsilon}_i(\beta_n)$  defined in

Lemma 3.2. Based on Corollary 3.9 and the sandwich covariance matrix estimator, an asymptotic  $(1 - \alpha)\%$  confidence interval ( $0 < \alpha < 1$ ) for  $\beta_j$  is

$$(3.10) \quad \widehat{\beta}_j \pm z_{\alpha/2} \mathbf{u}_j^T \widehat{\Sigma}_n \mathbf{u}_j,$$

where  $z_{\alpha/2}$  denotes the upper  $\frac{\alpha}{2}$  quantile of the standard normal distribution and  $\mathbf{u}_j$  is the unit vector of length  $p_n$  with the  $j$ th element equal to 1 and all the other elements equal to 0.

The sandwich covariance formula plays an important role in GEE methodology. In the “fixed  $p$ ” setup, it is known that the sandwich covariance matrix estimator provides a consistent estimator for the variance of the GEE estimator, even when the working correlation matrix is misspecified. The following theorem shows that this appealing property is still valid when  $p_n$  converges to  $\infty$  at an appropriate rate.

The proofs of Theorem 3.10 and Corollary 3.11 below are given in the Appendix.

THEOREM 3.10. Assume conditions (A1)–(A4) and that  $n^{-1} p_n^3 = o(1)$ . Then,

$$\mathbf{C}_n \widehat{\Sigma}_n \mathbf{C}_n^T - \mathbf{C}_n \Sigma_n \mathbf{C}_n^T = o_p(n^{-1}),$$

where  $\mathbf{C}_n$  is any  $l \times q_n$  matrix such that  $l$  is fixed and  $\mathbf{C}_n \mathbf{C}_n^T = \mathbf{G}$  with  $\mathbf{G}$  being an  $l \times l$  positive definite matrix.

REMARK 5. It is worth pointing out a subtle issue that is sometimes overlooked in the existing literature on high-dimensional analysis of independent data. In order to justify the validity of the asymptotic confidence interval or large-sample test for estimable contrast, it is necessary to show that the convergence rate in Theorem 3.10 is  $o_p(n^{-1})$ . Note that the estimable contrast is asymptotically normal with convergence rate  $O_p(n^{1/2})$ ; see, for example, Corollary 2.1 in He and Shao (2000) for the case of an M-estimator based on independent data. In the literature, sometimes only the  $o_p(1)$  rate is provided, which is not adequate, but can be fixed.

Next, we consider the large-sample Wald test for testing the following linear hypothesis:

$$H_0 : \mathbf{L}_n \beta_{n0} = \mathbf{0} \quad \text{vs.} \quad H_1 : \mathbf{L}_n \beta_{n0} \neq \mathbf{0},$$

where  $\mathbf{L}_n$  is an  $l \times p_n$  matrix with  $l$  fixed and  $\mathbf{L}_n \mathbf{L}_n^T = \mathbf{I}_l$ . The Wald test statistic is defined as

$$W_n = (\mathbf{L}_n \widehat{\beta}_n)^T (\mathbf{L}_n \widehat{\Sigma}_n \mathbf{L}_n^T)^{-1} (\mathbf{L}_n \widehat{\beta}_n).$$

The corollary below shows that the Wald test remains valid, even when the number of covariates diverges with the sample size.

COROLLARY 3.11. Assume conditions (A1)–(A4). If  $n^{-1}p_n^3 = o(1)$ , then  $W_n \rightarrow \chi_l^2$  in distribution under  $H_0$ , where  $\chi_l^2$  denotes the  $\chi^2$  distribution with  $l$  degrees of freedom.

REMARK 6. For testing a high-dimensional hypothesis  $H_0: \beta_n = \beta_{n0}^*$  versus  $H_1: \beta_n \neq \beta_{n0}^*$ , it can be shown that

$$(3.11) \quad \frac{(\widehat{\beta}_n - \beta_{n0}^*)^T \widehat{\Sigma}_n^{-1} (\widehat{\beta}_n - \beta_{n0}^*) - p_n}{\sqrt{2p_n}} \rightarrow N(0, 1)$$

in distribution under  $H_0$ , under some regularity conditions. A proof of this result is given in the supplementary article [Wang (2010)].

**4. Numerical studies.** We consider the following model for the marginal expectation of  $Y_{ij}$ ,  $i = 1, \dots, n$ , given  $\mathbf{X}_{ij}$ ,

$$(4.1) \quad \text{logit}(\pi_{ij}) = X_{ij}^T \beta_{n0}, \quad j = 1, 2, 3,$$

where  $\beta_{n0}$  is a  $p_n$ -dimensional vector of parameters with  $p_n = \lfloor 2.5n^{1/3} \rfloor$ , with  $\lfloor q \rfloor$  denoting the the largest integer not greater than  $q$ . In this example,  $\beta_{n0}^T = (0.4 \cdot \mathbf{1}_k^T, -0.3 \cdot \mathbf{1}_k^T, 0.2 \cdot \mathbf{1}_k^T, -0.1 \cdot \mathbf{1}_{p_n-3k}^T)$ , where  $\mathbf{1}_k$  denotes a  $k$ -dimensional vector of 1's and  $k = \lfloor p_n/4 \rfloor$ . In addition,  $X_{ij} = (x_{ij1}, \dots, x_{ijp_n})^T$  has a multivariate normal distribution with mean zero, marginal variance 0.2 and an AR-1 correlation matrix with autocorrelation coefficient 0.5. The binary response vector for each cluster has the above marginal mean and an exchangeable (also called compound symmetry or CS) correlation structure with correlation coefficient 0.5. Such correlated binary data are generated using Bahadur's representation [see, e.g., Fitzmaurice (1995)].

Since, for different sample sizes, the parameter dimension is different, we measure the accuracy of estimation by the *simulated average mean square error*, which is obtained by averaging  $\|\widehat{\beta}_n - \beta_{n0}\|^2/p_n$  over 500 simulated samples. Table 1 reports simulation results using four different working correlation structures: independence working correlation matrix (IN), unstructured working correlation matrix (UN), compound symmetry working correlation matrix (CS) and the first order autocorrelation working correlation matrix (AR-1), for sample sizes  $n = 500, 1000, 2000$  and  $3000$ . Table 1 demonstrates that when the covariate dimension grows at an appropriate rate with the sample size, the accuracy of GEE estimator is satisfactory. We also observe that when the true correlation matrix (CS in this case) is adopted, the estimator is more efficient.

We next examine the accuracy of the sandwich variance formula. The standard deviations of the estimated coefficients over 500 simulations are averaged and regarded as the true standard error (SD). Table 2 compares SD with the standard error obtained from the sandwich variance formula (SD2) when the unstructured working correlation matrix is used for estimating  $\widehat{\beta}_k, \widehat{\beta}_{2k}, \widehat{\beta}_{3k}$  and  $\widehat{\beta}_{p_n}$ . We observe that the sandwich variance formula works remarkably well. Similar phenomena are

TABLE 1  
*The simulated average mean squared error ( $\times 10$ ) for estimating  $\beta_{n0}$  using four different working correlation structures*

$n$	$p_n$	Working correlation structure			
		IN	UN	CS	AR-1
500	19	0.265	0.156	0.154	0.179
1000	24	0.141	0.103	0.100	0.111
2000	31	0.090	0.074	0.071	0.075
3000	36	0.070	0.065	0.063	0.065

also observed for estimating other regression coefficients and with other working correlation structures, but, for reasons of brevity, these are not reported.

Finally, we investigate hypothesis testing based on the large-sample Wald test. We consider model (4.1) with  $n = 1000$ ,  $p_n = 24$  and  $\beta_{n0}^T = (0.4 \cdot \mathbf{1}_6^T, -0.3 \cdot \mathbf{1}_6^T, 0.2 \cdot \mathbf{1}_6^T, -0.1 \cdot \mathbf{1}_2^T, 0, 0, 0, 0)$ . The left panel of Figure 1 depicts the density of the Wald test under the null hypothesis  $H_0: \beta_{21} = \beta_{22} = \beta_{23} = \beta_{24} = 0$  and compares it with the density curve of the  $\chi_4^2$  distribution. It demonstrates that the  $\chi^2$  approximation given in Corollary 3.11 is accurate. The right panel of Figure 1 gives the normal Q–Q plot for the Wald test statistic under the null hypothesis  $\beta_n = \beta_{n0}$  and it shows that the null distribution can be approximated well by a normal distribution for testing a higher-dimensional alternative, as discussed in Remark 6.

**5. Discussions.**

5.1. *Extension to general GEE.* Although the focus of the paper is on clustered binary data, the approaches and techniques can be extended to general GEE. For general GEE, the decomposition of  $\bar{\mathbf{D}}_n(\beta_n)$  given in Lemma 3.2 has a more complex expression, and the potential unboundedness of  $Y_{ij}$  makes the derivation

TABLE 2  
*Standard deviation (SD) and estimated standard deviation (SD2) using the sandwich variance formula*

$n$	$p_n$	$\hat{\beta}_k$		$\hat{\beta}_{2k}$		$\hat{\beta}_{3k}$		$\hat{\beta}_{p_n}$	
		SD	SD2	SD	SD2	SD	SD2	SD	SD2
500	19	0.126	0.111	0.114	0.110	0.117	0.111	0.089	0.098
1000	24	0.082	0.083	0.079	0.083	0.085	0.083	0.072	0.074
2000	31	0.073	0.060	0.063	0.060	0.065	0.060	0.051	0.053
3000	36	0.060	0.051	0.049	0.051	0.052	0.051	0.051	0.045



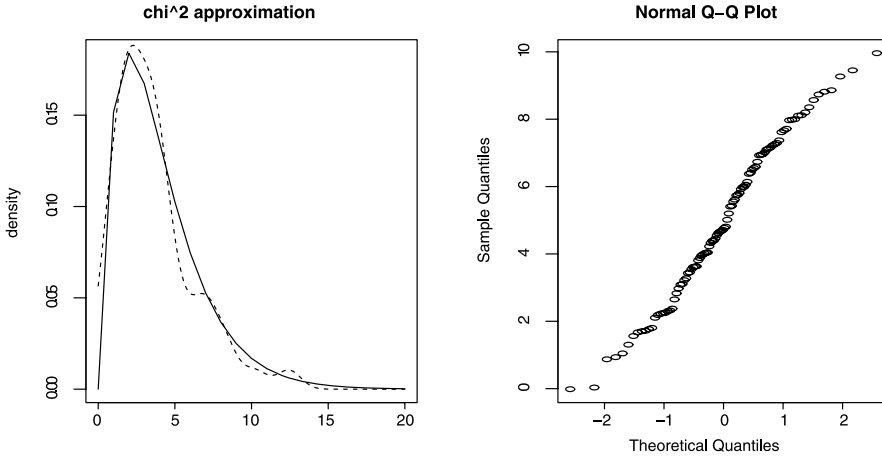


FIG. 1. The left panel gives the estimated null density of the large-sample Wald test (dashed curve) and the density of the chi-square distribution with four degrees of freedom (solid curve) for testing  $H_0 : \beta_{21} = \beta_{22} = \beta_{23} = \beta_{24} = 0$ . The right panel gives the normal Q–Q plot of the Wald test statistic under the null hypothesis  $\beta_n = \beta_{n0}$ .

of various probability bounds and asymptotic equivalence more delicate. Below, we give a brief discussion of the large- $p$  asymptotics for general GEE.

Assume that the first two marginal moments of  $Y_{ij}$  are  $\mu_{ij}(\beta_n) := E_{\beta_n}(Y_{ij}) = \mu(\theta_{ij})$  and  $\sigma_{ij}^2(\beta_n) := \text{Var}_{\beta_n}(Y_{ij}) = \dot{\mu}(\theta_{ij})$ , where  $\theta_{ij} = \mathbf{X}_{ij}^T \beta_n$ . These moment assumptions would follow when the marginal response variable has a canonical exponential family distribution with scaling parameter 1. Let  $\mathbf{A}_i(\beta_n) = \text{diag}(\sigma_{i1}^2(\beta_n), \dots, \sigma_{im}^2(\beta_n))$  and  $\boldsymbol{\mu}_i(\beta_n) = (\mu_{i1}(\beta_n), \dots, \mu_{im}(\beta_n))^T$ . The GEE estimator  $\hat{\beta}_n$  is the solution of

$$(5.1) \quad \sum_{i=1}^n \mathbf{X}_i^T \mathbf{A}_i^{1/2}(\beta_n) \hat{\mathbf{R}}^{-1} \mathbf{A}_i^{-1/2}(\beta_n) (\mathbf{Y}_i - \boldsymbol{\mu}_i(\beta_n)) = 0.$$

In addition to assumptions (A1)–(A4) in Section 3.2, we adopt two additional conditions:

(A5) there exists a finite constant  $M_1 > 0$  such that  $E(\|\mathbf{A}_i^{-1/2}(\beta_n)(\mathbf{Y}_i - \boldsymbol{\mu}_i(\beta_n))\|^{2+\delta}) \leq M_1$  for all  $i$  and some  $\delta > 0$ ;

(A6) if  $B_n = \{\beta_n : \|\beta_n - \beta_{n0}\| \leq \Delta \sqrt{p_n/n}\}$ , then  $\dot{\mu}(\mathbf{X}_{ij}^T \beta_n)$ ,  $1 \leq i \leq n$ ,  $1 \leq j \leq m$ , are uniformly bounded away from 0 and  $\infty$  on  $B_n$ ;  $\ddot{\mu}(\mathbf{X}_{ij}^T \beta_n)$  and  $\mu^{(3)}(\mathbf{X}_{ij}^T \beta_n)$ ,  $1 \leq i \leq n$ ,  $1 \leq j \leq m$ , are uniformly bounded by a finite positive constant  $M_2$  on  $B_n$ .

REMARK 7. Condition (A5) is similar to the condition in Lemma 2 of Xie and Yang (2003) and condition ( $\tilde{N}_\delta$ ) in Balan and Schiopu-Kratina (2005). Condition (A6) requires  $\mu_{ij}^{(k)}(\mathbf{X}_{ij}^T \boldsymbol{\beta}_n)$ ,  $k = 1, 2, 3$ , to be uniformly bounded when  $\boldsymbol{\beta}_n$  is in a local neighborhood around  $\boldsymbol{\beta}_{n0}$ . This condition is generally satisfied for GEE. For example, when the marginal model follows a Poisson distribution,  $\mu(t) = \exp(t)$ , thus  $\mu_{ij}^{(k)}(\mathbf{X}_{ij}^T \boldsymbol{\beta}_n) = \exp(\mathbf{X}_{ij}^T \boldsymbol{\beta}_n)$ ,  $k = 1, 2, 3$ , are uniformly bounded on  $B_n$ .

THEOREM 5.1. Assume conditions (A1)–(A6) and that  $n^{-1} p_n^2 = o(1)$ . The generalized estimating equation (5.1) then has a root  $\hat{\boldsymbol{\beta}}_n$  such that  $\|\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_{n0}\| = O_p(\sqrt{p_n/n})$ . Furthermore, if  $n^{-1} p_n^3 = o(1)$ , then  $\forall \boldsymbol{\alpha}_n \in R^{p_n}$  such that  $\|\boldsymbol{\alpha}_n\| = 1$ ,

$$\boldsymbol{\alpha}_n^T \bar{\mathbf{M}}_n^{-1/2}(\boldsymbol{\beta}_{n0}) \bar{\mathbf{H}}_n(\boldsymbol{\beta}_{n0})(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_{n0}) \rightarrow N(0, 1)$$

in distribution, where  $\bar{\mathbf{M}}_n^{-1/2}(\boldsymbol{\beta}_{n0})$  and  $\bar{\mathbf{H}}_n(\boldsymbol{\beta}_{n0})$  have the same expressions as in Section 3.2.

A sketch of the proof of Theorem 5.1 is given in the supplementary article [Wang (2010)].

5.2. *Related problems.* In some scenarios, a “large  $n$ , diverging  $m$ ” asymptotic framework, where  $p$  is either fixed or also diverges at an appropriate rate, may be more appropriate. This corresponds to a real situation where the cluster size is itself large. For example, in a longitudinal study, doctors take measurements on the patients during each visit. Each patient forms a cluster. The cluster size is large if the number of visits is large. For a fixed  $p$  setting, this “large  $n$ , diverging  $m$ ” asymptotic framework has been considered by Xie and Yang (2003). A future topic of interest is to consider large  $m$  together with large  $p$ .

Another interesting direction for future study is to consider a more flexible semi-parametric specification for the generalized estimating equations in the large- $p$  setting. In the classical “fixed  $p$ ” setting, GEE with partially linear model specification has been investigated by Lin and Carroll (2001a, 2001b), Lin and Ying (2001), He, Zhu and Fung (2002), Fan and Li (2004), Chiou and Müller (2005), Wang, Carroll and Lin (2005), Chen and Jin (2006), He, Fung and Zhu (2006) and Huang, Zhang and Zhou (2007), among others.

## APPENDIX

We use  $C$  to denote a generic positive constant that can vary from line to line.

PROOF OF (3.3). It suffices [Ortega and Rheinboldt (1970)] to show that  $\forall \varepsilon > 0$ , there exists a  $\Delta > 0$  such that for all  $n$  sufficiently large,

$P(\sup_{\|\beta_n - \beta_{n0}\| = \Delta\sqrt{p_n/n}} (\beta_n - \beta_{n0})^T \tilde{S}_n(\beta_n) < 0) \geq 1 - \varepsilon$ . We have

$$\begin{aligned} & (\beta_n - \beta_{n0})^T \tilde{S}_n(\beta_n) \\ &= (\beta_n - \beta_{n0})^T \tilde{S}_n(\beta_{n0}) + (\beta_n - \beta_{n0})^T \frac{\partial}{\partial \beta_n^T} \tilde{S}_n(\beta_n^*)(\beta_n - \beta_{n0}) \\ &\triangleq I_{n1} + I_{n2}, \end{aligned}$$

where  $\beta_n^*$  lies between  $\beta_{n0}$  and  $\beta_n$ . We first consider  $I_{n1}$ . For any  $\beta_n$  such that  $\|\beta_n - \beta_{n0}\| = \Delta\sqrt{\frac{p_n}{n}}$ , we have  $|I_{n1}| \leq \Delta\sqrt{\frac{p_n}{n}} \|\tilde{S}_n(\beta_{n0})\|$ . Note that

$$\begin{aligned} E[\|\tilde{S}_n(\beta_{n0})\|^2] &= E\left[\sum_{i=1}^n (\mathbf{Y}_i - \boldsymbol{\pi}_i(\beta_{n0}))^T \mathbf{X}_i \mathbf{X}_i^T (\mathbf{Y}_i - \boldsymbol{\pi}_i(\beta_{n0}))\right] \\ &\leq E\left[\sum_{i=1}^n \lambda_{\max}(\mathbf{X}_i \mathbf{X}_i^T) \|\mathbf{Y}_i - \boldsymbol{\pi}_i(\beta_{n0})\|^2\right] \\ &\leq C \operatorname{Tr}\left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T\right) = C \sum_{i=1}^n \sum_{j=1}^m \mathbf{X}_{ij}^T \mathbf{X}_{ij} = O(np_n), \end{aligned}$$

by assumption (A1). Thus,  $|I_{n1}| \leq \Delta O_p(p_n)$ . Next,

$$\begin{aligned} I_{n2} &= -(\beta_n - \beta_{n0})^T \left[ \sum_{i=1}^n \mathbf{X}_i^T A_i(\beta_{n0}) \mathbf{X}_i \right] (\beta_n - \beta_{n0}) \\ &\quad - (\beta_n - \beta_{n0})^T \left[ \sum_{i=1}^n \mathbf{X}_i^T (A_i(\beta_n^*) - A_i(\beta_{n0})) \mathbf{X}_i \right] (\beta_n - \beta_{n0}) \\ &\triangleq I_{n21} + I_{n22}. \end{aligned}$$

Note that  $I_{n21} \leq -\lambda_{\min}(\mathbf{A}_i(\beta_0)) \lambda_{\min}(\sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i) \|\beta_n - \beta_{n0}\|^2 \leq -C p_n \Delta^2$ , by (A3). Since  $\frac{\partial}{\partial \beta_n} \mathbf{A}_{ij}(\beta_n) = \pi_{ij}(\beta_n)(1 - \pi_{ij}(\beta_n))(1 - 2\pi_{ij}(\beta_n)) \mathbf{X}_{ij}$ , we have

$$\begin{aligned} |I_{n22}| &\leq (\beta_n - \beta_{n0})^T \left[ \sum_{i=1}^n \sum_{j=1}^m |\mathbf{A}_{ij}(\beta_n^*) - \mathbf{A}_{ij}(\beta_{n0})| \mathbf{X}_{ij} \mathbf{X}_{ij}^T \right] (\beta_n - \beta_{n0}) \\ &\leq \sup_{i,j} \|\mathbf{X}_{ij}\| \cdot \|\beta_n^* - \beta_0\| \cdot \|\beta_n - \beta_0\|^2 \cdot \lambda_{\max}\left(\sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i\right) \\ &\leq O(\sqrt{p_n}) O_p(\sqrt{p_n/n}) (\Delta^2 p_n/n) O(n) = \Delta^2 o_p(p_n), \end{aligned}$$

by (A1)–(A3). Thus, for sufficiently large  $\Delta$ ,  $(\beta_n - \beta_{n0})^T \tilde{S}_n(\beta_n)$  is dominated by  $I_{n21}$ , which is large and negative for all sufficiently large  $n$ .  $\square$

PROOF OF (3.4). The proof is given in the online supplement.  $\square$

PROOF OF LEMMA 3.1. Let  $\mathbf{Q} = \{q_{j_1, j_2}\}_{1 \leq j_1, j_2 \leq m}$  denote the matrix  $\widehat{\mathbf{R}}^{-1} - \overline{\mathbf{R}}^{-1}$ . Then,

$$\begin{aligned} \mathbf{S}_n(\boldsymbol{\beta}_{n0}) - \overline{\mathbf{S}}_n(\boldsymbol{\beta}_{n0}) &= \sum_{i=1}^n \sum_{j_1=1}^m \sum_{j_2=1}^m q_{j_1, j_2} \mathbf{A}_{ij_1}^{1/2}(\boldsymbol{\beta}_{n0}) \mathbf{A}_{ij_2}^{-1/2}(\boldsymbol{\beta}_{n0}) (Y_{ij_2} - \pi_{ij_2}(\boldsymbol{\beta}_{n0})) \mathbf{X}_{ij_1} \\ &= \sum_{j_1=1}^m \sum_{j_2=1}^m q_{j_1, j_2} \left[ \sum_{i=1}^n \mathbf{A}_{ij_1}^{1/2}(\boldsymbol{\beta}_{n0}) \varepsilon_{ij_2}(\boldsymbol{\beta}_{n0}) \mathbf{X}_{ij_1} \right], \end{aligned}$$

where  $\varepsilon_{ij_2}(\boldsymbol{\beta}_{n0}) = \mathbf{A}_{ij_2}^{-1/2}(\boldsymbol{\beta}_{n0}) (Y_{ij_2} - \pi_{ij_2}(\boldsymbol{\beta}_{n0}))$ . Note that

$$\begin{aligned} E \left[ \left\| \sum_{i=1}^n \mathbf{A}_{ij_1}^{1/2}(\boldsymbol{\beta}_{n0}) \varepsilon_{ij_2}(\boldsymbol{\beta}_{n0}) \mathbf{X}_{ij_1} \right\|^2 \right] &= \sum_{i=1}^n \mathbf{A}_{ij_1}(\boldsymbol{\beta}_{n0}) E[\varepsilon_{ij_2}^2(\boldsymbol{\beta}_{n0})] \mathbf{X}_{ij_1}^T \mathbf{X}_{ij_1} \\ &\leq \sum_{i=1}^n \mathbf{X}_{ij_1}^T \mathbf{X}_{ij_1} = O(np_n). \end{aligned}$$

Thus,  $\| \sum_{i=1}^n \mathbf{A}_{ij_1}^{1/2}(\boldsymbol{\beta}_{n0}) \varepsilon_{ij_2}(\boldsymbol{\beta}_{n0}) \mathbf{X}_{ij_1} \| = O_p(\sqrt{np_n}) \forall 1 \leq j_1, j_2 \leq m$ . Since, by (A4),  $q_{j_1, j_2} = O_p(\sqrt{p_n/n}) \forall 1 \leq j_1, j_2 \leq m$ , the proof is complete.  $\square$

PROOF OF LEMMA 3.2. The derivation can be found in Pan (2002).  $\square$

PROOF OF LEMMA 3.3. Let  $\mathbf{H}_n(\boldsymbol{\beta}_n)$ ,  $\mathbf{E}_n(\boldsymbol{\beta}_n)$  and  $\mathbf{G}_n(\boldsymbol{\beta}_n)$  be defined the same as  $\overline{\mathbf{H}}_n(\boldsymbol{\beta}_n)$ ,  $\overline{\mathbf{E}}_n(\boldsymbol{\beta}_n)$  and  $\overline{\mathbf{G}}_n(\boldsymbol{\beta}_n)$ , respectively, but with  $\widehat{\mathbf{R}}$  replaced by  $\widehat{\mathbf{R}}$ . By Lemma 3.2, it is sufficient to prove the following three results:

$$\begin{aligned} \sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} \sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T [\mathbf{H}_n(\boldsymbol{\beta}_n) - \overline{\mathbf{H}}_n(\boldsymbol{\beta}_n)] \mathbf{b}_n| \\ \text{(A.1)} \quad = O_p(\sqrt{np_n}), \end{aligned}$$

$$\begin{aligned} \sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} \sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T [\mathbf{E}_n(\boldsymbol{\beta}_n) - \overline{\mathbf{E}}_n(\boldsymbol{\beta}_n)] \mathbf{b}_n| \\ \text{(A.2)} \quad = O_p(\sqrt{np_n}), \end{aligned}$$

$$\begin{aligned} \sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} \sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T [\mathbf{G}_n(\boldsymbol{\beta}_n) - \overline{\mathbf{G}}_n(\boldsymbol{\beta}_n)] \mathbf{b}_n| \\ \text{(A.3)} \quad = O_p(\sqrt{np_n}). \end{aligned}$$

We have

$$\begin{aligned} & |\mathbf{b}_n^T [\mathbf{H}_n(\boldsymbol{\beta}_n) - \bar{\mathbf{H}}_n(\boldsymbol{\beta}_n)] \mathbf{b}_n| \\ &= \left| \sum_{i=1}^n \mathbf{b}_n^T \mathbf{X}_i^T \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_n) [\widehat{\mathbf{R}}^{-1} - \bar{\mathbf{R}}^{-1}] \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_n) \mathbf{X}_i \mathbf{b}_n \right| \\ &\leq \|\widehat{\mathbf{R}}^{-1} - \bar{\mathbf{R}}^{-1}\| \lambda_{\max}(\mathbf{A}_i(\boldsymbol{\beta}_n)) \lambda_{\max} \left( \sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i \right) \|\mathbf{b}_n\|^2. \end{aligned}$$

By assumptions (A2) and (A4), (A.1) is proved. Next, note that

$$\begin{aligned} & |\mathbf{b}_n^T [\mathbf{E}_n(\boldsymbol{\beta}_n) - \bar{\mathbf{E}}_n(\boldsymbol{\beta}_n)] \mathbf{b}_n| \\ &= \frac{1}{2} \left| \sum_{i=1}^n \sum_{j=1}^m (1 - 2\pi_{ij}(\boldsymbol{\beta}_n)) \varepsilon_{ij}(\boldsymbol{\beta}_n) \mathbf{b}_n^T \mathbf{X}_i^T \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_n) \right. \\ &\quad \left. \times [\widehat{\mathbf{R}}^{-1} - \bar{\mathbf{R}}^{-1}] \mathbf{e}_j \mathbf{e}_j^T \mathbf{X}_i \mathbf{b}_n \right| \\ &\leq \sum_{i=1}^n \sum_{j=1}^m \mathbf{A}_{ij}^{-1/2}(\boldsymbol{\beta}_n) |\mathbf{b}_n^T \mathbf{X}_i^T \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_n) [\widehat{\mathbf{R}}^{-1} - \bar{\mathbf{R}}^{-1}] \mathbf{e}_j| \cdot |\mathbf{e}_j^T \mathbf{X}_i \mathbf{b}_n| \\ &\leq \sum_{i=1}^n \sum_{j=1}^m \mathbf{A}_{ij}^{-1/2}(\boldsymbol{\beta}_n) \|\widehat{\mathbf{R}}^{-1} - \bar{\mathbf{R}}^{-1}\| \cdot \|\mathbf{A}_i^{1/2}(\boldsymbol{\beta}_n)\| \cdot \|\mathbf{X}_i \mathbf{b}_n\|^2. \end{aligned}$$

Thus,

$$\begin{aligned} & \sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} \sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T [\mathbf{E}_n(\boldsymbol{\beta}_n) - \bar{\mathbf{E}}_n(\boldsymbol{\beta}_n)] \mathbf{b}_n| \\ &\leq C \|\widehat{\mathbf{R}}^{-1} - \bar{\mathbf{R}}^{-1}\| \cdot \sum_{i=1}^n \sum_{j=1}^m \sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} \mathbf{A}_{ij}^{-1/2}(\boldsymbol{\beta}_n) \sup_{\|\mathbf{b}_n\|=1} \|\mathbf{X}_i \mathbf{b}_n\|^2 \\ &= O_p(\sqrt{p_n/n}) O(n) = O_p(\sqrt{np_n}), \end{aligned}$$

by assumption (A3). (A.3) is proved similarly.  $\square$

PROOF OF LEMMA 3.4. By (3.6), it is sufficient to verify that

$$(A.4) \quad \sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} \sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T \bar{\mathbf{E}}_n(\boldsymbol{\beta}_n) \mathbf{b}_n| = O_p(\sqrt{n} p_n),$$

$$(A.5) \quad \sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} \sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T \bar{\mathbf{G}}_n(\boldsymbol{\beta}_n) \mathbf{b}_n| = O_p(\sqrt{n} p_n).$$

First, note that we have the following decomposition of  $\bar{\mathbf{E}}_n(\boldsymbol{\beta}_n)$ :

$$\begin{aligned}
\bar{\mathbf{E}}_n(\boldsymbol{\beta}_n) &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m (1 - 2\pi_{ij}(\boldsymbol{\beta}_{n0})) \varepsilon_{ij}(\boldsymbol{\beta}_{n0}) \mathbf{X}_i^T \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0}) \bar{\mathbf{R}}^{-1} \mathbf{e}_j \mathbf{e}_j^T \mathbf{X}_i \\
&\quad + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m (1 - 2\pi_{ij}(\boldsymbol{\beta}_{n0})) \varepsilon_{ij}(\boldsymbol{\beta}_{n0}) \mathbf{X}_i^T [\mathbf{A}_i^{1/2}(\boldsymbol{\beta}_n) - \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0})] \\
&\quad \quad \quad \times \bar{\mathbf{R}}^{-1} \mathbf{e}_j \mathbf{e}_j^T \mathbf{X}_i \\
&\quad + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m [(1 - 2\pi_{ij}(\boldsymbol{\beta}_n)) \mathbf{A}_{ij}^{-1/2}(\boldsymbol{\beta}_n) - (1 - 2\pi_{ij}(\boldsymbol{\beta}_{n0})) \mathbf{A}_{ij}^{-1/2}(\boldsymbol{\beta}_{n0})] \\
&\quad \quad \quad \times (Y_{ij} - \pi_{ij}(\boldsymbol{\beta}_{n0})) \mathbf{X}_i^T \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_n) \bar{\mathbf{R}}^{-1} \mathbf{e}_j \mathbf{e}_j^T \mathbf{X}_i \\
&\quad + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m (1 - 2\pi_{ij}(\boldsymbol{\beta}_n)) \mathbf{A}_{ij}^{-1/2}(\boldsymbol{\beta}_n) (\pi_{ij}(\boldsymbol{\beta}_{n0}) - \pi_{ij}(\boldsymbol{\beta}_n)) \\
&\quad \quad \quad \times \mathbf{X}_i^T \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_n) \bar{\mathbf{R}}^{-1} \mathbf{e}_j \mathbf{e}_j^T \mathbf{X}_i \\
&\triangleq \bar{\mathbf{E}}_{1n}(\boldsymbol{\beta}_{n0}) + \sum_{k=2}^4 \bar{\mathbf{E}}_{kn}(\boldsymbol{\beta}_n).
\end{aligned}$$

Thus, to prove (A.4), it suffices to verify that  $\sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T \bar{\mathbf{E}}_{1n}(\boldsymbol{\beta}_{n0}) \mathbf{b}_n| = O_P(\sqrt{n} p_n)$  and  $\sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} \sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T \bar{\mathbf{E}}_{kn}(\boldsymbol{\beta}_n) \mathbf{b}_n| = O_P(\sqrt{n} p_n)$ . We first prove that  $\sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T \bar{\mathbf{E}}_{1n}(\boldsymbol{\beta}_{n0}) \mathbf{b}_n| = O_P(\sqrt{n} p_n)$ , by verifying that  $\|\bar{\mathbf{E}}_{1n}(\boldsymbol{\beta}_{n0})\| = O_P(\sqrt{n} p_n)$ , where  $\|\bar{\mathbf{E}}_{1n}(\boldsymbol{\beta}_{n0})\| = \sqrt{\text{trace}(\bar{\mathbf{E}}_{1n}(\boldsymbol{\beta}_{n0}) \bar{\mathbf{E}}_{1n}^T(\boldsymbol{\beta}_{n0}))}$ :

$$\begin{aligned}
&E[\|\bar{\mathbf{E}}_{1n}(\boldsymbol{\beta}_{n0})\|^2] \\
&= \frac{1}{4} \sum_{i=1}^n \sum_{j_1=1}^m \sum_{j_2=1}^m (1 - 2\pi_{ij_1}(\boldsymbol{\beta}_{n0})) (1 - 2\pi_{ij_2}(\boldsymbol{\beta}_{n0})) E[\varepsilon_{ij_1}(\boldsymbol{\beta}_{n0}) \varepsilon_{ij_2}(\boldsymbol{\beta}_{n0})] \\
&\quad \quad \quad \times \text{trace}[\mathbf{X}_i^T \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0}) \bar{\mathbf{R}}^{-1} \mathbf{e}_{j_1} \mathbf{e}_{j_1}^T \mathbf{X}_i \mathbf{X}_i^T \mathbf{e}_{j_2} \mathbf{e}_{j_2}^T \\
&\quad \quad \quad \quad \quad \quad \times \bar{\mathbf{R}}^{-1} \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0}) \mathbf{X}_i] \\
&\leq C \sum_{i=1}^n \sum_{j_1=1}^m \sum_{j_2=1}^m |\mathbf{e}_{j_1}^T \mathbf{X}_i \mathbf{X}_i^T \mathbf{e}_{j_2} \mathbf{e}_{j_2}^T \bar{\mathbf{R}}^{-1} \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0}) \\
&\quad \quad \quad \times \mathbf{X}_i \mathbf{X}_i^T \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0}) \bar{\mathbf{R}}^{-1} \mathbf{e}_{j_1}|
\end{aligned}$$

$$\leq C \sum_{i=1}^n \sum_{j_1=1}^m \sum_{j_2=1}^m \|\mathbf{e}_{j_1}^T \mathbf{X}_i\| \cdot \|\mathbf{X}_i^T \mathbf{e}_{j_2}\| \cdot \|\mathbf{e}_{j_2}^T \bar{\mathbf{R}}^{-1} \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0}) \mathbf{X}_i\| \\ \times \|\mathbf{X}_i^T \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0}) \bar{\mathbf{R}}^{-1} \mathbf{e}_{j_1}\|.$$

Note that  $\|\mathbf{e}_{j_1}^T \mathbf{X}_i\| = \|\mathbf{X}_{ij_1}\|$ ,  $\|\mathbf{X}_i^T \mathbf{e}_{j_2}\| = \|\mathbf{X}_{ij_2}\|$ ,  $\|\mathbf{e}_{j_2}^T \bar{\mathbf{R}}^{-1} \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0}) \mathbf{X}_i\| \leq C(\text{trace}(\mathbf{X}_i \mathbf{X}_i^T))^{1/2}$  and  $\|\mathbf{X}_i^T \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0}) \bar{\mathbf{R}}^{-1} \mathbf{e}_{j_1}\| \leq C(\text{trace}(\mathbf{X}_i \mathbf{X}_i^T))^{1/2}$ . Thus,

$$E[\|\bar{\mathbf{E}}_{1n}(\boldsymbol{\beta}_{n0})\|^2] \leq C \sum_{i=1}^n \sum_{j_1=1}^m \sum_{j_2=1}^m \|\mathbf{X}_{ij_1}\| \cdot \|\mathbf{X}_{ij_2}\| \text{trace}(\mathbf{X}_i \mathbf{X}_i^T) \\ \leq C \cdot \max_{i,j} \|\mathbf{X}_{ij}\|^2 \text{trace}\left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T\right) = O(np_n^2),$$

by assumptions (A1) and (A3). This implies that  $\sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T \bar{\mathbf{E}}_{1n}(\boldsymbol{\beta}_{n0}) \mathbf{b}_n| = O_p(\sqrt{n} p_n)$ . Next, we have

$$|\mathbf{b}_n^T \bar{\mathbf{E}}_{2n}(\boldsymbol{\beta}_n) \mathbf{b}_n| \\ = \left| \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m (1 - 2\pi_{ij}(\boldsymbol{\beta}_{n0})) \varepsilon_{ij}^{1/2}(\boldsymbol{\beta}_{n0}) \mathbf{b}_n^T \mathbf{X}_i^T [\mathbf{A}_i^{1/2}(\boldsymbol{\beta}_n) - \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0})] \right. \\ \left. \times \bar{\mathbf{R}}^{-1} \mathbf{e}_j \mathbf{e}_j^T \mathbf{X}_i \mathbf{b}_n \right| \\ \leq C \sum_{i=1}^n \sum_{j=1}^m |\mathbf{b}_n^T \mathbf{X}_i^T [\mathbf{A}_i^{1/2}(\boldsymbol{\beta}_n) - \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0})] \bar{\mathbf{R}}^{-1} \mathbf{e}_j| \cdot |\mathbf{e}_j^T \mathbf{X}_i \mathbf{b}_n| \\ \leq C \sum_{i=1}^n \sum_{j=1}^m \|\mathbf{X}_i \mathbf{b}_n\|^2 \lambda_{\max}(\bar{\mathbf{R}}^{-1}) \max_j |\mathbf{A}_{ij}^{1/2}(\boldsymbol{\beta}_n) - \mathbf{A}_{ij}^{1/2}(\boldsymbol{\beta}_{n0})|.$$

Note that there exists some  $\boldsymbol{\beta}_n^*$  between  $\boldsymbol{\beta}_n$  and  $\boldsymbol{\beta}_{n0}$  such that

$$\mathbf{A}_{ij}^{1/2}(\boldsymbol{\beta}_n) - \mathbf{A}_{ij}^{1/2}(\boldsymbol{\beta}_{n0}) = \frac{1}{2} \mathbf{A}_{ij}^{1/2}(\boldsymbol{\beta}_n^*) (1 - 2\pi_{ij}(\boldsymbol{\beta}_n^*)) \mathbf{X}_{ij}^T (\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}) \\ \leq C \|\mathbf{X}_{ij}\| \cdot \|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\|.$$

Therefore,

$$\sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} \sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T \bar{\mathbf{E}}_{2n}(\boldsymbol{\beta}_n) \mathbf{b}_n| \\ \leq C \max_{i,j} \|\mathbf{X}_{ij}\| \sup_{\|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \leq \Delta \sqrt{p_n/n}} \|\boldsymbol{\beta}_n - \boldsymbol{\beta}_{n0}\| \cdot \lambda_{\max}\left(\sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i\right) \\ = O(\sqrt{p_n}) O(\sqrt{p_n/n}) O(n) = O(\sqrt{n} p_n).$$

Similarly, we can show that  $\sup_{\|\beta_n - \beta_{n0}\| \leq \Delta \sqrt{p_n/n}} \sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T \bar{\mathbf{E}}_{kn}(\beta_n) \mathbf{b}_n| = O(\sqrt{n} p_n)$ ,  $k = 3, 4$ . This proves (A.4). Similarly, we can prove (A.5).  $\square$

PROOF OF LEMMA 3.5. The proof is given in the online supplementary material.  $\square$

PROOF OF LEMMA 3.7. We write  $\alpha_n^T \bar{\mathbf{M}}_n^{-1/2}(\beta_{n0}) \bar{\mathbf{S}}_n(\beta_{n0}) = \sum_{i=1}^n Z_{ni}$ , where  $Z_{ni} = \alpha_n^T \bar{\mathbf{M}}_n^{-1/2}(\beta_{n0}) \mathbf{X}_i^T \mathbf{A}_i^{1/2}(\beta_{n0}) \bar{\mathbf{R}}^{-1} \boldsymbol{\varepsilon}_i(\beta_{n0})$ . Since  $\bar{\mathbf{M}}_n(\beta_{n0}) = \text{Cov}(\bar{\mathbf{S}}_n(\beta_{n0}))$ , we have  $\text{Var}(\alpha_n^T \bar{\mathbf{M}}_n^{-1/2}(\beta_{n0}) \bar{\mathbf{S}}_n(\beta_{n0})) = 1$ . To establish the asymptotic normality, it suffices to check the Lindberg condition, that is,  $\forall \varepsilon > 0$ ,  $\sum_{i=1}^n E[Z_{ni}^2 I(|Z_{ni}| > \varepsilon)] \rightarrow 0$ . By the Cauchy–Schwarz inequality,

$$\begin{aligned} Z_{ni}^2 &\leq \|\alpha_n^T \bar{\mathbf{M}}_n^{-1/2}(\beta_{n0}) \mathbf{X}_i^T \mathbf{A}_i^{1/2}(\beta_{n0}) \bar{\mathbf{R}}^{-1}\|^2 \cdot \|\boldsymbol{\varepsilon}_i(\beta_{n0})\|^2 \\ &\leq \lambda_{\max}(\bar{\mathbf{R}}^{-2}) \lambda_{\max}(\mathbf{A}_i(\beta_{n0})) (\alpha_n^T \bar{\mathbf{M}}_n^{-1/2}(\beta_{n0}) \mathbf{X}_i^T \mathbf{X}_i \bar{\mathbf{M}}_n^{-1/2}(\beta_{n0}) \alpha_n) \\ &\quad \times \|\boldsymbol{\varepsilon}_i(\beta_{n0})\|^2 \\ &\leq C \gamma_{ni} \|\boldsymbol{\varepsilon}_i(\beta_{n0})\|^2, \end{aligned}$$

where  $\gamma_{ni} \triangleq \alpha_n^T \bar{\mathbf{M}}_n^{-1/2}(\beta_{n0}) \mathbf{X}_i^T \mathbf{X}_i \bar{\mathbf{M}}_n^{-1/2}(\beta_{n0}) \alpha_n$ . Next, we will show that  $\max_{1 \leq i \leq n} \gamma_{ni} \rightarrow 0$  as  $n \rightarrow \infty$ . Note that  $\gamma_{ni} \leq \lambda_{\max}(\mathbf{X}_i^T \mathbf{X}_i) \lambda_{\min}^{-1}(\bar{\mathbf{M}}_n(\beta_{n0}))$ . Since  $\bar{\mathbf{M}}_n(\beta_{n0})$  is symmetric, to evaluate  $\lambda_{\min}(\bar{\mathbf{M}}_n(\beta_{n0}))$ ,  $\forall \mathbf{b}_n \in R^{p_n}$ , we have

$$\begin{aligned} \mathbf{b}_n^T \bar{\mathbf{M}}_n(\beta_{n0}) \mathbf{b}_n &\geq \lambda_{\min}(\mathbf{R}_0) \lambda_{\min}(\bar{\mathbf{R}}^{-2}) \sum_{i=1}^n \lambda_{\min}(\mathbf{A}_i(\beta_{n0})) \mathbf{b}_n^T \mathbf{X}_i^T \mathbf{X}_i \mathbf{b}_n \\ &\geq C \mathbf{b}_n^T \left( \sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i \right) \mathbf{b}_n \geq C \lambda_{\min} \left( \sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i \right) \|\mathbf{b}_n\|^2. \end{aligned}$$

Thus,  $\inf_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T \bar{\mathbf{M}}_n(\beta_{n0}) \mathbf{b}_n| \geq C \lambda_{\min}(\sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i)$  and this implies that  $\lambda_{\min} \times (\bar{\mathbf{M}}_n(\beta_{n0})) \geq C \lambda_{\min}(\sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i)$ . Therefore, we have

$$\gamma_{ni} \leq \frac{\lambda_{\max}(\mathbf{X}_i^T \mathbf{X}_i)}{C \lambda_{\min}(\sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i)} \leq \frac{\text{Tr}(\mathbf{X}_i^T \mathbf{X}_i)}{C \lambda_{\min}(\sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i)} = \frac{\sum_{j=1}^m \mathbf{X}_{ij}^T \mathbf{X}_{ij}}{C \lambda_{\min}(\sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i)}.$$

It follows that  $\max_{1 \leq i \leq n} \gamma_{ni} \leq O(n^{-1} p_n) = o(1)$ . We have

$$\sum_{i=1}^n E[Z_{ni}^2 I(|Z_{ni}| > \varepsilon)] \leq \sum_{i=1}^n C \gamma_{ni} E \left[ \|\boldsymbol{\varepsilon}_i(\beta_{n0})\|^2 I \left\{ \|\boldsymbol{\varepsilon}_i(\beta_{n0})\|^2 > \frac{\varepsilon^2}{C \gamma_{ni}} \right\} \right].$$

Note that  $\|\boldsymbol{\varepsilon}_i(\beta_{n0})\|^2$  is uniformly bounded, by assumption (A2). Thus, for all  $\varepsilon > 0$  and  $\delta > 0$ , there exists a positive integer  $N$  such that (1)  $I\{\|\boldsymbol{\varepsilon}_i(\beta_{n0})\|^2 >$



$\frac{\varepsilon^2}{C\gamma_{ni}}\} = 0$  for all  $n > N$ ; (2)  $\sum_{i=1}^N C\gamma_{ni} \leq \delta$  for all  $n$  sufficiently large. This ensures that

$$\sum_{i=1}^n C\gamma_{ni} E \left[ \|\boldsymbol{\varepsilon}_i(\boldsymbol{\beta}_{n0})\|^2 I \left\{ \|\boldsymbol{\varepsilon}_i(\boldsymbol{\beta}_{n0})\|^2 > \frac{\varepsilon^2}{C\gamma_{ni}} \right\} \right] \rightarrow 0.$$

Therefore, the Lindberg condition is verified.  $\square$

**PROOF OF THEOREM 3.10.** It is sufficient to show that for  $\mathbf{b}_n \in R^{p_n}$ ,

$$(A.6) \quad \sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T (\widehat{\boldsymbol{\Sigma}}_n - \boldsymbol{\Sigma}_n) \mathbf{b}_n| = o_p(n^{-1}).$$

We use the conclusion of Theorem 3.6 throughout the proof. Note that we can write  $\widehat{\boldsymbol{\Sigma}}_n - \boldsymbol{\Sigma}_n = I_{n1} + I_{n2} + I_{n3}$ , where

$$\begin{aligned} I_{n1} &= \mathbf{H}_n^{-1}(\widehat{\boldsymbol{\beta}}_n) [\widehat{\mathbf{M}}_n(\widehat{\boldsymbol{\beta}}_n) - \overline{\mathbf{M}}_n(\boldsymbol{\beta}_{n0})] \mathbf{H}_n^{-1}(\widehat{\boldsymbol{\beta}}_n), \\ I_{n2} &= [\mathbf{H}_n^{-1}(\widehat{\boldsymbol{\beta}}_n) - \overline{\mathbf{H}}_n^{-1}(\boldsymbol{\beta}_{n0})] \overline{\mathbf{M}}_n(\boldsymbol{\beta}_{n0}) \mathbf{H}_n^{-1}(\widehat{\boldsymbol{\beta}}_n), \\ I_{n3} &= \overline{\mathbf{H}}_n^{-1}(\boldsymbol{\beta}_{n0}) \overline{\mathbf{M}}_n(\boldsymbol{\beta}_{n0}) [\mathbf{H}_n^{-1}(\widehat{\boldsymbol{\beta}}_n) - \overline{\mathbf{H}}_n^{-1}(\boldsymbol{\beta}_{n0})]. \end{aligned}$$

Thus, (A.6) is implied by  $\sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T I_{ni} \mathbf{b}_n| = o_p(1)$ . We have

$$\begin{aligned} &\sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T I_{n1} \mathbf{b}_n| \\ &\leq \frac{\max(|\lambda_{\max}(\widehat{\mathbf{M}}_n(\widehat{\boldsymbol{\beta}}_n) - \overline{\mathbf{M}}_n(\boldsymbol{\beta}_{n0}))|, |\lambda_{\min}(\widehat{\mathbf{M}}_n(\widehat{\boldsymbol{\beta}}_n) - \overline{\mathbf{M}}_n(\boldsymbol{\beta}_{n0}))|)}{\lambda_{\min}^2(\mathbf{H}_n(\widehat{\boldsymbol{\beta}}_n))}. \end{aligned}$$

To evaluate the eigenvalues of  $\widehat{\mathbf{M}}_n(\widehat{\boldsymbol{\beta}}_n) - \overline{\mathbf{M}}_n(\boldsymbol{\beta}_{n0})$ , we have

$$\begin{aligned} &|\mathbf{c}_n^T [\widehat{\mathbf{M}}_n(\widehat{\boldsymbol{\beta}}_n) - \overline{\mathbf{M}}_n(\boldsymbol{\beta}_{n0})] \mathbf{c}_n| \\ &\leq |\mathbf{c}_n^T [\widehat{\mathbf{M}}_n(\widehat{\boldsymbol{\beta}}_n) - \widehat{\mathbf{M}}_n(\boldsymbol{\beta}_{n0})] \mathbf{c}_n| + |\mathbf{c}_n^T [\widehat{\mathbf{M}}_n(\boldsymbol{\beta}_{n0}) - \overline{\mathbf{M}}_n(\boldsymbol{\beta}_{n0})] \mathbf{c}_n| \end{aligned}$$

for  $\mathbf{c}_n \in R^{p_n}$ . Note that

$$\begin{aligned} &\sup_{\|\mathbf{c}_n\|=1} |\mathbf{c}_n^T [\widehat{\mathbf{M}}_n(\widehat{\boldsymbol{\beta}}_n) - \widehat{\mathbf{M}}_n(\boldsymbol{\beta}_{n0})] \mathbf{c}_n| \\ &\leq \sup_{\|\mathbf{c}_n\|=1} \left| \sum_{i=1}^n \mathbf{c}_n^T \mathbf{X}_i^T [\mathbf{A}_i^{1/2}(\widehat{\boldsymbol{\beta}}_n) - \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0})] \widehat{\mathbf{R}}^{-1} \boldsymbol{\varepsilon}_i(\widehat{\boldsymbol{\beta}}_n) \boldsymbol{\varepsilon}_i^T(\widehat{\boldsymbol{\beta}}_n) \right. \\ &\quad \left. \times \widehat{\mathbf{R}}^{-1} \mathbf{A}_i^{1/2}(\widehat{\boldsymbol{\beta}}_n) \mathbf{X}_i \mathbf{c}_n \right| \\ &+ \sup_{\|\mathbf{c}_n\|=1} \left| \sum_{i=1}^n \mathbf{c}_n^T \mathbf{X}_i^T \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0}) \widehat{\mathbf{R}}^{-1} \boldsymbol{\varepsilon}_i(\widehat{\boldsymbol{\beta}}_n) \boldsymbol{\varepsilon}_i^T(\widehat{\boldsymbol{\beta}}_n) \widehat{\mathbf{R}}^{-1} \right. \end{aligned}$$

$$\begin{aligned} & \times [\mathbf{A}_i^{1/2}(\widehat{\boldsymbol{\beta}}_n) - \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0})] \mathbf{X}_i \mathbf{c}_n \Big| \\ & + \sup_{\|\mathbf{c}_n\|=1} \left| \sum_{i=1}^n \mathbf{c}_n^T \mathbf{X}_i^T \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0}) \widehat{\mathbf{R}}^{-1} [\boldsymbol{\varepsilon}_i(\widehat{\boldsymbol{\beta}}_n) \boldsymbol{\varepsilon}_i^T(\widehat{\boldsymbol{\beta}}_n) - \boldsymbol{\varepsilon}_i(\boldsymbol{\beta}_{n0}) \boldsymbol{\varepsilon}_i^T(\boldsymbol{\beta}_{n0})] \right. \\ & \qquad \qquad \qquad \left. \times \widehat{\mathbf{R}}^{-1} \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0}) \mathbf{X}_i \mathbf{c}_n \right| \end{aligned}$$

$$\triangleq \sup_{\|\mathbf{c}_n\|=1} J_{n1} + \sup_{\|\mathbf{c}_n\|=1} J_{n2} + \sup_{\|\mathbf{c}_n\|=1} J_{n3}.$$

Note that

$$J_{n1} \leq \sum_{i=1}^n \|\mathbf{c}_n^T \mathbf{X}_i^T [\mathbf{A}_i^{1/2}(\widehat{\boldsymbol{\beta}}_n) - \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0})]\| \cdot \|\widehat{\mathbf{R}}^{-1} \boldsymbol{\varepsilon}_i(\widehat{\boldsymbol{\beta}}_n)\|^2 \cdot \|\mathbf{A}_i^{1/2}(\widehat{\boldsymbol{\beta}}_n) \mathbf{X}_i \mathbf{c}_n\|.$$

We have  $\|\mathbf{A}_i^{1/2}(\widehat{\boldsymbol{\beta}}_n) \mathbf{X}_i \mathbf{c}_n\| \leq \|\mathbf{X}_i \mathbf{c}_n\|$  and

$$\begin{aligned} \|\mathbf{c}_n^T \mathbf{X}_i^T [\mathbf{A}_i^{1/2}(\widehat{\boldsymbol{\beta}}_n) - \mathbf{A}_i^{1/2}(\boldsymbol{\beta}_{n0})]\| & \leq \|\mathbf{X}_i \mathbf{c}_n\| \max_j |\mathbf{A}_{ij}^{1/2}(\widehat{\boldsymbol{\beta}}_n) - \mathbf{A}_{ij}^{1/2}(\boldsymbol{\beta}_{n0})| \\ & \leq C \|\mathbf{X}_i \mathbf{c}_n\| \cdot \|\mathbf{X}_{ij}\| \cdot \|\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_{n0}\|. \end{aligned}$$

Furthermore,

$$\begin{aligned} \|\widehat{\mathbf{R}}^{-1} \boldsymbol{\varepsilon}_i(\widehat{\boldsymbol{\beta}}_n)\|^2 & = (\mathbf{Y}_i - \boldsymbol{\pi}_i(\widehat{\boldsymbol{\beta}}_n))^T \mathbf{A}_i^{-1/2}(\widehat{\boldsymbol{\beta}}_n) \widehat{\mathbf{R}}^{-2} \mathbf{A}_i^{-1/2}(\widehat{\boldsymbol{\beta}}_n) (\mathbf{Y}_i - \boldsymbol{\pi}_i(\widehat{\boldsymbol{\beta}}_n)) \\ & \leq \lambda_{\max}(\widehat{\mathbf{R}}^{-2}) \lambda_{\max}(\mathbf{A}_i^{-1}(\widehat{\boldsymbol{\beta}}_n)) \|\mathbf{Y}_i - \boldsymbol{\pi}_i(\widehat{\boldsymbol{\beta}}_n)\|^2 \leq C O_p(1). \end{aligned}$$

Thus,

$$\sup_{\|\mathbf{c}_n\|=1} J_{n1} \leq O_p(1) \|\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_{n0}\| \max_{i,j} \|\mathbf{X}_{ij}\| \lambda_{\max} \left( \sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i \right) = o_p(n).$$

Similarly,  $\sup_{\|\mathbf{c}_n\|=1} J_{n2} = o_p(n)$  and  $\sup_{\|\mathbf{c}_n\|=1} J_{n3} = o_p(n)$ . Thus,

$$\sup_{\|\mathbf{c}_n\|=1} |\mathbf{c}_n^T [\widehat{\mathbf{M}}_n(\widehat{\boldsymbol{\beta}}_n) - \widehat{\mathbf{M}}_n(\boldsymbol{\beta}_{n0})] \mathbf{c}_n| = o_p(n).$$

Similarly,  $\sup_{\|\mathbf{c}_n\|=1} |\mathbf{c}_n^T [\widehat{\mathbf{M}}_n(\boldsymbol{\beta}_{n0}) - \overline{\mathbf{M}}_n(\boldsymbol{\beta}_{n0})] \mathbf{c}_n| = o_p(n)$ . Finally, note that

$$\begin{aligned} \lambda_{\min}(\mathbf{H}_n(\widehat{\boldsymbol{\beta}}_n)) & \geq \lambda_{\min}(\widehat{\mathbf{R}}) \min_{i,j} (\pi_{ij}(\widehat{\boldsymbol{\beta}}_n)(1 - \pi_{ij}(\widehat{\boldsymbol{\beta}}_n))) \lambda_{\min} \left( \sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i \right) \\ & = O_p(n). \end{aligned}$$

Thus,  $\sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T I_{n1} \mathbf{b}_n| = o_p(n^{-1})$ . We can also prove that  $\sup_{\|\mathbf{b}_n\|=1} |\mathbf{b}_n^T I_{ni} \times \mathbf{b}_n| = o_p(n^{-1})$ ,  $i = 2, 3$ , by first noting that

$$\mathbf{H}_n^{-1}(\widehat{\boldsymbol{\beta}}_n) - \overline{\mathbf{H}}_n^{-1}(\boldsymbol{\beta}_{n0}) = [\mathbf{H}_n^{-1}(\widehat{\boldsymbol{\beta}}_n) - \overline{\mathbf{H}}_n^{-1}(\widehat{\boldsymbol{\beta}}_n)] + [\overline{\mathbf{H}}_n^{-1}(\widehat{\boldsymbol{\beta}}_n) - \overline{\mathbf{H}}_n^{-1}(\boldsymbol{\beta}_{n0})]$$

and then using the expressions

$$\begin{aligned} \mathbf{H}_n^{-1}(\widehat{\boldsymbol{\beta}}_n) - \overline{\mathbf{H}}_n^{-1}(\widehat{\boldsymbol{\beta}}_n) &= \overline{\mathbf{H}}_n^{-1}(\widehat{\boldsymbol{\beta}}_n)[\overline{\mathbf{H}}_n(\widehat{\boldsymbol{\beta}}_n) - \mathbf{H}_n(\widehat{\boldsymbol{\beta}}_n)]\mathbf{H}_n^{-1}(\widehat{\boldsymbol{\beta}}_n), \\ \overline{\mathbf{H}}_n^{-1}(\widehat{\boldsymbol{\beta}}_n) - \overline{\mathbf{H}}_n^{-1}(\boldsymbol{\beta}_{n0}) &= \overline{\mathbf{H}}_n^{-1}(\boldsymbol{\beta}_{n0})[\overline{\mathbf{H}}_n(\boldsymbol{\beta}_{n0}) - \overline{\mathbf{H}}_n(\widehat{\boldsymbol{\beta}}_n)]\overline{\mathbf{H}}_n^{-1}(\widehat{\boldsymbol{\beta}}_n). \quad \square \end{aligned}$$

PROOF OF COROLLARY 3.11. It is sufficient to show that

$$(A.7) \quad [(\mathbf{L}_n \widehat{\boldsymbol{\Sigma}}_n \mathbf{L}_n^T)^{-1/2} - (\mathbf{L}_n \boldsymbol{\Sigma}_n \mathbf{L}_n^T)^{-1/2}] \mathbf{L}_n (\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_{n0}) \rightarrow 0$$

in probability. Note that the left-hand side can be written as

$$[(\mathbf{L}_n \widehat{\boldsymbol{\Sigma}}_n \mathbf{L}_n^T)^{-1/2} (\mathbf{L}_n \boldsymbol{\Sigma}_n \mathbf{L}_n^T)^{1/2} - \mathbf{I}_l] (\mathbf{L}_n \boldsymbol{\Sigma}_n \mathbf{L}_n^T)^{-1/2} \mathbf{L}_n (\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_{n0})$$

and thus (A.7) is implied by

$$(\mathbf{L}_n \widehat{\boldsymbol{\Sigma}}_n \mathbf{L}_n^T)^{-1} (\mathbf{L}_n \boldsymbol{\Sigma}_n \mathbf{L}_n^T) - \mathbf{I}_l = (\mathbf{L}_n \widehat{\boldsymbol{\Sigma}}_n \mathbf{L}_n^T)^{-1} \mathbf{L}_n (\boldsymbol{\Sigma}_n - \widehat{\boldsymbol{\Sigma}}_n) \mathbf{L}_n^T = o_p(1).$$

Let  $\mathbf{u}_i$  denote the  $l \times 1$  unit vector with the  $i$ th element being 1 and all of the other elements being 0. Then, for all  $1 \leq i, j \leq l$ , we have, by the Cauchy–Schwarz inequality,

$$\begin{aligned} &|\mathbf{u}_i^T (\mathbf{L}_n \widehat{\boldsymbol{\Sigma}}_n \mathbf{L}_n^T)^{-1} \mathbf{L}_n (\boldsymbol{\Sigma}_n - \widehat{\boldsymbol{\Sigma}}_n) \mathbf{L}_n^T \mathbf{u}_j| \\ &\leq |\mathbf{u}_i^T (\mathbf{L}_n \widehat{\boldsymbol{\Sigma}}_n \mathbf{L}_n^T)^{-2} \mathbf{u}_i|^{1/2} |\mathbf{u}_j^T [\mathbf{L}_n (\boldsymbol{\Sigma}_n - \widehat{\boldsymbol{\Sigma}}_n) \mathbf{L}_n^T]^2 \mathbf{u}_j|^{1/2} \\ &\leq \frac{\|\mathbf{L}_n (\boldsymbol{\Sigma}_n - \widehat{\boldsymbol{\Sigma}}_n) \mathbf{L}_n^T\|}{\lambda_{\min}(\mathbf{L}_n \widehat{\boldsymbol{\Sigma}}_n \mathbf{L}_n^T)}. \end{aligned}$$

Now, for any  $l$ -dimensional vector such that  $\|\mathbf{b}\| = 1$ , we have

$$\begin{aligned} |\mathbf{b}^T \mathbf{L}_n \widehat{\boldsymbol{\Sigma}}_n \mathbf{L}_n^T \mathbf{b}| &\geq |\mathbf{b}^T \mathbf{L}_n \boldsymbol{\Sigma}_n \mathbf{L}_n^T \mathbf{b}| - |\mathbf{b}^T \mathbf{L}_n (\widehat{\boldsymbol{\Sigma}}_n - \boldsymbol{\Sigma}_n) \mathbf{L}_n^T \mathbf{b}| \\ &\geq \lambda_{\min}(\boldsymbol{\Sigma}_n) + o_p(n^{-1}) \\ &\geq \frac{\lambda_{\min}(\overline{\mathbf{M}}_n(\boldsymbol{\beta}_{n0}))}{\lambda_{\max}^2(\overline{\mathbf{H}}_n(\boldsymbol{\beta}_{n0}))} + o_p(n^{-1}), \end{aligned}$$

where the second inequality uses Theorem 3.10. By (3.9),  $\lambda_{\min}(\overline{\mathbf{M}}_n(\boldsymbol{\beta}_{n0})) \geq c_1 \lambda_{\min}(\sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i)$  for some positive constant  $c_1$ . Similarly, we can show that  $\lambda_{\max}(\overline{\mathbf{H}}_n(\boldsymbol{\beta}_{n0})) \leq c_2 \lambda_{\max}(\sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i)$  for some positive constant  $c_2$ . Thus,  $\lambda_{\min}(\mathbf{L}_n \widehat{\boldsymbol{\Sigma}}_n \mathbf{L}_n^T) \geq O_p(n^{-1})$ . This proves that  $\frac{\|\mathbf{L}_n (\boldsymbol{\Sigma}_n - \widehat{\boldsymbol{\Sigma}}_n) \mathbf{L}_n^T\|}{\lambda_{\min}(\mathbf{L}_n \widehat{\boldsymbol{\Sigma}}_n \mathbf{L}_n^T)} = o_p(1)$ , by Theorem 3.10.  $\square$

**Acknowledgments.** The author would like to thank the Associate Editor and two referees for their constructive and insightful comments that significantly improved this paper.

## SUPPLEMENTARY MATERIAL

**Supplement to “GEE analysis of clustered binary data with diverging number of covariates”** (DOI: [10.1214/10-AOS846SUPP](https://doi.org/10.1214/10-AOS846SUPP); .pdf). The proofs of (3.3), Lemma 3.5, (3.11) and Theorem 5.1 are provided in this supplementary article [Wang (2010)].

## REFERENCES

- BALAN, R. M. and SCHIOPU-KRATINA, I. (2005). Asymptotic results with generalized estimating equations for longitudinal data. *Ann. Statist.* **32** 522–541. [MR2163150](#)
- BAI, Z. and WU, Y. (1994). Limiting behavior of M-estimators of regression coefficients in high dimensional linear models, I. Scale-dependent case. *J. Multivariate Anal.* **51** 211–239. [MR1321295](#)
- CHAGANTY, N. R. and JOE, H. (2004). Efficiency of generalized estimating equations for binary responses. *J. Roy. Statist. Soc. Ser. B* **66** 851–860. [MR2102468](#)
- CHEN, K. and JIN, Z. (2006). Partial linear regression models for clustered data. *J. Amer. Statist. Assoc.* **101** 195–204. [MR2268038](#)
- CHEN, S. X., PENG, L. and QIN, Y.-L. (2009). Effects of data dimension on empirical likelihood. *Biometrika* **96** 711–722. [MR2538767](#)
- CHIOU, J.-M. and MÜLLER, H.-G. (2005). Estimated estimating equations: Semiparametric inference for clustered and longitudinal data. *J. Roy. Statist. Soc. Ser. B* **67** 531–553. [MR2168203](#)
- DONOHO, D. L. (2000). High-dimensional data analysis: The curses and blessings of dimensionality. In *Math Challenges of 21st Century* 1–32. Amer. Math. Soc., Providence, RI.
- FAN, J. and LI, R. (2004). New estimation and model selection procedures for semiparametric modeling in longitudinal data analysis. *J. Amer. Statist. Assoc.* **99** 710–723. [MR2090905](#)
- FAN, J. and LI, R. (2006). Statistical challenges with high dimensionality: Feature selection in knowledge discovery. In *Proceedings of the International Congress of Mathematicians* (M. Sanz-Sole, J. Soria, J. L. Varona and J. Verdera eds.) **III** 595–622. Eur. Math. Soc., Zürich. [MR2275698](#)
- FAN, J. and LV, J. (2010). A selective overview of variable selection in high dimensional feature space. *Statist. Sinica* **20** 101–148. [MR2640659](#)
- FAN, J. and PENG, H. (2004). Nonconcave penalized likelihood with a diverging number of parameters. *Ann. Statist.* **32** 928–961. [MR2065194](#)
- FITZMAURICE, G. M. (1995). A caveat concerning independence estimating equations with multivariate binary data. *Biometrics* **51** 309–317.
- HE, X., FUNG, W. K. and ZHU, Z. Y. (2006). Robust estimation in generalized partial linear models for clustered data. *J. Amer. Statist. Assoc.* **100** 1176–1184. [MR2236433](#)
- HE, X. and SHAO, Q. M. (2000). On parameters of increasing dimensions. *J. Multivariate Anal.* **73** 120–135. [MR1766124](#)
- HE, X., ZHU, Z. Y. and FUNG, W. K. (2002). Estimation in a semiparametric model for longitudinal data with unspecified dependence structure. *Biometrika* **89** 579–590. [MR1929164](#)
- HJORT, H. L., MCKEAGUE, I. W. and VAN KEILEGOM, I. (2009). Extending the scope of empirical likelihood. *Ann. Statist.* **37** 1079–1111. [MR2509068](#)
- HUANG, J., HOROWITZ, J. L. and MA, S. J. (2008). Asymptotic properties of bridge estimators in sparse high-dimensional regression models. *Ann. Statist.* **36** 587–613. [MR2396808](#)
- HUANG, J. Z., ZHANG, L. and ZHOU, L. (2007). Efficient estimation in marginal partially linear models for longitudinal/clustered data using splines. *Scand. J. Statist.* **34** 451–477. [MR2368793](#)
- HUBER, P. J. (1973). Robust regression: Asymptotics, conjectures and Monte Carlo. *Ann. Statist.* **1** 799–821. [MR0356373](#)
- LAM, C. and FAN, J. (2008). Profile-kernel likelihood inference with diverging number of parameters. *Ann. Statist.* **36** 2232–2260. [MR2458186](#)

- LI, B. (1997). On the consistency of generalized estimating equations. In *Selected Proceedings of the Symposium on Estimating Functions* 115–136. *IMS Lecture Notes—Monograph Series* **32**. IMS, Hayward, CA. [MR1837801](#)
- LIANG, K. Y. and ZEGER, S. L. (1986). Longitudinal data analysis using generalised linear models. *Biometrika* **73** 12–22. [MR0836430](#)
- LIN, D. Y. and YING, Z. (2001). Semiparametric and nonparametric regression analysis of longitudinal data. *J. Amer. Statist. Assoc.* **96** 103–112. [MR1952726](#)
- LIN, X. and CARROLL, R. J. (2001a). Semiparametric regression for clustered data using generalized estimating equations. *J. Amer. Statist. Assoc.* **96** 1045–1056. [MR1947252](#)
- LIN, X. and CARROLL, R. J. (2001b). Semiparametric regression for clustered data. *Biometrika* **88** 1179–1185. [MR1872228](#)
- MAMMEN, E. (1989). Asymptotics with increasing dimension for robust regression with applications to the bootstrap. *Ann. Statist.* **17** 382–400. [MR0981457](#)
- MCCULLAGH, P. and NELDER, J. A. (1989). *Generalized Linear Models*, 2nd ed. Chapman and Hall, London. [MR0727836](#)
- ORTEGA, J. M. and RHEINBOLDT, W. C. (1970). *Iterative Solution of Nonlinear Equations in Several Variables*, Academic Press, San Diego. [MR0273810](#)
- PAN, W. (2002). Goodness-of-fit tests for GEE with correlated binary data. *Scand. J. Statist.* **29** 101–110. [MR1894384](#)
- PORTNOY, S. (1984). Asymptotic behavior of M estimators of  $p$  regression parameters when  $p^2/n$  is large. I. Consistency. *Ann. Statist.* **12** 1298–1309. [MR0760690](#)
- PORTNOY, S. (1985). Asymptotic behavior of M estimators of  $p$  regression parameters when  $p^2/n$  is large. II. Normal approximation. *Ann. Statist.* **13** 1403–1417. [MR0811499](#)
- PORTNOY, S. (1988). Asymptotic properties of likelihood methods for exponential families when the number of parameters tends to infinity. *Ann. Statist.* **16** 356–366. [MR0924876](#)
- XIE, M. and YANG, Y. (2003). Asymptotics for generalized estimating equations with large cluster sizes. *Ann. Statist.* **31** 310–347. [MR1962509](#)
- WANG, J. L., XUE, L. G., ZHU, L. X. and CHONG, Y. S. (2010). Estimation for a partial-linear single-index model. *Ann. Statist.* **38** 246–274. [MR2589322](#)
- WANG, L. (2010). Supplement to “GEE analysis of clustered binary data with diverging number of covariates.” DOI: [10.1214/10-AOS846SUPP](https://doi.org/10.1214/10-AOS846SUPP).
- WANG, N., CARROLL, R. J. and LIN, X. (2005). Efficient semiparametric marginal estimation for longitudinal/clustered data. *J. Amer. Statist. Assoc.* **100** 147–157. [MR2156825](#)
- WELSH, A. H. (1989). On  $M$ -processes and  $M$ -estimation. *Ann. Statist.* **17** 337–361. [MR0981455](#)
- ZHU, L. P. and ZHU, L. X. (2009). On distribution-weighted partial least squares with diverging number of highly correlated predictors. *J. Roy. Statist. Soc. Ser. B* **71** 525–548.
- ZOU, H. and ZHANG, H. H. (2009). On the adaptive elastic-net with a diverging number of parameters. *Ann. Statist.* **37** 1733–1751. [MR2533470](#)

SCHOOL OF STATISTICS  
UNIVERSITY OF MINNESOTA  
224 CHURCH STREET, SE  
MINNEAPOLIS, MINNESOTA 55455  
USA  
E-MAIL: [lan@stat.umn.edu](mailto:lan@stat.umn.edu)