

# Uniform Inference in High-Dimensional Dynamic Panel Data Models With Approximately Sparse Fixed Effects

ANDERS BREDAHL KOCK\*

HAIHAN TANG<sup>†</sup>

February 22, 2018

## Abstract

We establish oracle inequalities for a version of the Lasso in high-dimensional fixed effects dynamic panel data models. The inequalities are valid for the coefficients of the dynamic and exogenous regressors. Separate oracle inequalities are derived for the fixed effects. Next, we show how one can conduct uniformly valid inference on the parameters of the model and construct a uniformly valid estimator of the asymptotic covariance matrix which is robust to conditional heteroskedasticity in the error terms. Allowing for conditional heteroskedasticity is important in dynamic models as the conditional error variance may be non-constant over time and depend on the covariates. Furthermore, our procedure allows for inference on high-dimensional subsets of the parameter vector of an increasing cardinality. We show that the confidence bands resulting from our procedure are asymptotically honest and contract at the optimal rate. This rate is different for the fixed effects than for the remaining parts of the parameter vector.

*Keywords:* Dynamic Panel Data, High-dimensional Data, Uniform Inference, Oracle Inequality, Lasso.

*JEL codes:* C12, C13, C23.

## 1 Introduction

Dynamic panel data models are widely used in economics and social sciences. They are extremely popular as workers, firms, and countries often differ due to unobserved factors. Furthermore, these units are often sampled repeatedly over time in many modern applications thus

---

\*University of Oxford, Aarhus University and CREATES, Department of Economics, Manor Road, Oxford, OX1 3UQ, UK. Email: anders.kock@economics.ox.ac.uk. Financial support from the Center for Research in the Econometric Analysis of Time Series (grant DNRF78) is gratefully acknowledged.

<sup>†</sup>Fanhai International School of Finance and School of Economics, Fudan University. Email: hh-tang@fudan.edu.cn.

allowing one to model their dynamic development. However, so far no work has been done on how to conduct inference in the high-dimensional dynamic fixed effects model

$$y_{i,t} = \sum_{l=1}^L \alpha_l y_{i,t-l} + x'_{i,t} \beta + \eta_i + \varepsilon_{i,t}, \quad i = 1, \dots, N, \quad \text{and} \quad t = 1, \dots, T \quad (1.1)$$

with potentially more parameters than observation. Here the presence of  $L$  lags of  $y_{i,t}$  allows for autoregressive dependence of  $y_{i,t}$  on its own past.  $x_{i,t}$  is a  $p_x \times 1$  vector of exogenous variables and  $\eta_i, i = 1, \dots, N$  are the  $N$  individual effects while  $\varepsilon_{i,t}$  are idiosyncratic error terms.<sup>1</sup> Applications of panel data are widespread: ranging from wage regressions where one seeks to explain workers' salaries, to models of economic growth determining the factors that impact growth over time of a panel of countries as in Islam (1995).

Recent years have witnessed a surge in availability of big data sets including many explanatory variables. For example, De Neve et al. (2012) have considered the effect of genes on happiness/life satisfaction. Controlling for many genes simultaneously clearly results in a vast set of explanatory variables, hence calling for techniques which can handle such a setting. High-dimensionality may also arise out of a desire to control for flexible functional forms by including various transformations, such as cross products, of the available explanatory variables. In the specific context of panel data models Andersen et al. (2012) investigated the causal effect of lightning density on economic growth using a US panel data set. These authors had access to many control variables compared to the sample size. For this reason, they decided to investigate the effect of lightning using several subsets of control variables instead of including all control variables simultaneously as one would ideally do. In this paper we show how one can achieve this ideal by proposing an inferential procedure for high-dimensional dynamic panel data models.

Much progress has also been made on the methodological side of high-dimensional models in the last decade. Among the most popular procedures is the Lasso of Tibshirani (1996) which sparked a lot of research on its properties. However, until recently, not much work had been done on inference in high-dimensional models for Lasso-type estimators as these possess a rather complicated limiting distribution even in the low dimensional case, see Knight and Fu (2000). This problem has been cleverly approached by unpenalized estimation after double selection by Belloni et al. (2012, 2014) or by desparsification in Zhang and Zhang (2014); van de Geer et al. (2014); Javanmard and Montanari (2013); Caner and Kock (2014).

The focus in the above mentioned work has been almost exclusively on independent data and often on the plain linear regression model while high-dimensional panel data has not been

---

<sup>1</sup>Alternative names for  $\eta_i$ s are fixed effects and unobserved heterogeneities. In this paper we shall use these three names interchangeably.

treated. Exceptions are Kock (2013) and Belloni et al. (2015) who have established oracle inequalities and asymptotically valid inference for a low-dimensional parameter in *static* panel data models, respectively. Kock (2016) has studied high-dimensional panel data models with correlated random effects while Caner and Zhang (2014) have studied the properties of penalized GMM, which can be used to estimate dynamic panel data models, in the case of fewer parameters than observations. Lu and Su (2016) and Su et al. (2016) have considered shrinkage estimators in panel data models with special structures such as interactive fixed effects or latent structures. To the best of our knowledge, no research has been conducted on inference in high-dimensional dynamic panel data models with more parameters than observations. Note that high-dimensionality may arise from three sources in the dynamic panel data model (1.1). These sources are the coefficients pertaining to the lagged left hand side variables ( $\alpha_l$ ), the exogenous variables ( $\beta$ ), as well as the fixed effects ( $\eta_i$ ). In particular, we shall see that (joint) inference involving an  $\eta_i$  behaves in a markedly different way from inference only involving the  $\alpha_l$ s and  $\beta$ . Furthermore, panel data differ from the classic linear regression model in that one does not have independence across  $t = 1, \dots, T$  for any  $i$  as consecutive observations in time can be highly correlated for any given individual. Ignoring this dependence may lead to gravely misleading inference even in low-dimensional panel data models. For that reason we shall make *no* assumptions on this dependence structure across  $t = 1, \dots, T$  for the  $x_{i,t}$ . Static panel data models are a special case of (1.1) corresponding to  $\alpha_l = 0$ ,  $l = 1, \dots, L$ .

Traditional approaches to inference in low-dimensional static panel data models have considered the  $N$  fixed effects  $\eta_i$  as nuisance parameters which have been removed by either taking first differences or demeaning the data over time for each individual  $i$ , see e.g., Wooldridge (2010); Arellano (2003); Baltagi (2008). In this paper we take the stand that the fixed effects may be of intrinsic interest. Estimating the  $\eta_i$ s precisely is essential for obtaining precise cross-sectional forecasts or predictions for each individual, like in credit scoring or in the estimation of probabilities of tax fraud (Arellano (2003) p11). Thus, we do not remove them by first differencing or demeaning. This allows us to test hypotheses involving  $\alpha, \beta$  and the  $\eta_i$ s. An important recently developed approach focusing on how to jointly determine group membership and parameter estimation is the Classifier-Lasso of Su et al. (2016). The authors propose a novel penalty function to achieve these goals.

In an alternative framework, some treat the  $N$  individual effects as  $N$  random variables. We would like to remark that this is not the framework we are adopting in most of the paper.<sup>2</sup>

---

<sup>2</sup>The two exceptions are in Section 5 (Monte Carlo) and Appendix A, where we need to adopt the alternative framework to specify a data generating process for the individual effects and to justify the weak sparsity on the  $N$  realised values of individual effects, respectively.

By treating  $\eta_i$ s as  $N$  parameters, we are essentially considering the  $N$  realisations of the  $N$  individual effects in the alternative framework. We shall impose that  $(\eta_1, \dots, \eta_N)$  is weakly sparse in a sense to be made precise in Section 2.2.

In an interesting recent paper dealing with the low-dimensional case, Bonhomme and Manresa (2015) have assumed a different type of structure, namely grouping, on the fixed effects. However, in the high-dimensional setting we are considering, weak sparsity works well as just explained.

Our inferential procedure is closest in spirit to the one in van de Geer et al. (2014), which in turn builds on Zhang and Zhang (2014), who cleverly used nodewise regressions to *desparsify* the Lasso and to construct an approximate inverse of the non-invertible sample Gram matrix in the context of the linear regression model. In particular, we show how nodewise regressions can be used to construct one of the blocks of the approximate inverse of the empirical Gram matrix in dynamic panel data models. As opposed to van de Geer et al. (2014), we do not require the inverse covariance matrix of the covariates to be exactly sparse. It suffices that the rows of the inverse covariance matrix are weakly sparse. Thus, none of its entries needs to be zero.

We contribute by first establishing an oracle inequalities for a version of the Lasso in dynamic panel data models for all groups of parameters. As can be expected, the fixed effects turn out to behave differently from the remaining parameters. Next, we show how joint asymptotically gaussian inference may be conducted on the three types of parameters in (1.1). In particular, we show that hypotheses involving an increasing number of parameters can be tested and provide a uniformly consistent estimator of the asymptotic covariance matrix which is robust to conditional heteroskedasticity in the error terms. Thus, we introduce a feasible procedure for inference in high-dimensional heteroskedastic dynamic panel data models. Allowing for conditional heteroskedasticity is important in dynamic models like the one considered here as the conditional variance is known to often depend on the current state of the process, see e.g. Engle (1982). Thus, assuming the error terms to be independent of the covariates and possessing a constant variance is not reasonable. Next, we show that confidence bands constructed by our procedure are asymptotically honest (uniform) in the sense of Li (1989) over a certain subset of the parameter space. Finally, we show that the confidence bands have uniformly the optimal rate of contraction for all types of parameters. Thus, the honesty is not bought at the price of wide confidence bands as is the case for sparse estimators, c.f. Pötscher (2009). Simulations reveal that our procedure performs well in terms of size, power, and coverage rate of the constructed intervals.

The rest of the paper is organized as follows. Section 2 introduces the estimator and provides

oracle inequalities for all types of parameters. Next, Section 3 shows how limiting Gaussian inference may be conducted and provides a feasible estimator of the covariance matrix which is robust to heteroskedasticity even in the case where the number of parameter estimates for which we seek the limiting distribution diverges with the sample size. Section 4 shows that confidence intervals constructed by our procedure are honest and contract at the optimal rate for all types of parameters. Section 5 studies our estimator in Monte Carlo experiments while Section 6 concludes. Appendix A contains sufficient conditions for some of our assumptions. Appendices B - E contain the proofs for our oracle inequality, desparsification, inference, and honest confidence intervals, respectively. Appendix F contains further auxiliary lemmas needed in Appendices A - E. Appendix G contains some technical expositions omitted in the main text.

## 2 The Model

### 2.1 Notation

For  $x \in \mathbb{R}^n$ , let  $\|x\|_0 = \sum_{i=1}^n 1(x_i \neq 0)$ ,  $\|x\| = \sqrt{\sum_{i=1}^n x_i^2}$ ,  $\|x\|_1 = \sum_{i=1}^n |x_i|$  and  $\|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$  denote the  $\ell_0$ ,  $\ell_2$ ,  $\ell_1$  and  $\ell_\infty$  norms, respectively. Let  $e_m$  denote the unit column vector with  $m$ th entry being 1 in some Euclidean space whose dimension depends on the context. If the argument of  $\|\cdot\|_\infty$  is a matrix, then  $\|\cdot\|_\infty$  denotes the maximal absolute element of the matrix. For some generic set  $R \subseteq \{1, \dots, n\}$ , let  $x_R \in \mathbb{R}^{|R|}$  denote the vector obtained by extracting the elements of  $x \in \mathbb{R}^n$  whose indices are in  $R$ , where  $|R|$  denotes the cardinality of  $R$ ;  $R^c = \{1, \dots, n\} \setminus R$ . For an  $n \times n$  matrix  $A$ ,  $A_R$  denotes the submatrix consisting of the rows and columns indexed by  $R$ .  $\otimes$  is the Kronecker product. For real numbers  $a, b$  let  $a \vee b$  and  $a \wedge b$  denote  $\max(a, b)$  and  $\min(a, b)$ , respectively. For two real sequences  $(a_n)$  and  $(b_n)$ ,  $a_n \lesssim b_n$  means that  $a_n \leq Cb_n$  for some fixed, finite and positive constant  $C$  for all  $n \geq 1$ . Similarly, we write  $a_n \asymp b_n$  if there exist constants  $0 < a_1 \leq a_2$  such that  $a_1 b_n \leq a_n \leq a_2 b_n$  for all  $n \geq 1$ .  $\text{sgn}(\cdot)$  is the sign function.  $\text{maxeval}(\cdot)$  and  $\text{mineval}(\cdot)$  are the maximal and minimal eigenvalues of the argument, respectively. For some vector  $x \in \mathbb{R}^n$ ,  $\text{diag}(x)$  returns an  $n \times n$  diagonal matrix with  $x$  supplying the diagonal entries.

The model in (1.1) can be rewritten as

$$y_{i,t} = z'_{i,t} \alpha + \eta_i + \varepsilon_{i,t}, \quad i = 1, \dots, N, \quad t = 1, \dots, T \quad (2.1)$$

where  $z_{i,t} := (y_{i,t-1}, \dots, y_{i,t-L}, x'_{i,t})'$  and  $\alpha := (\alpha_1, \dots, \alpha_L, \beta')'$  are  $p \times 1$  vectors ( $p = p_x + L$ ). We shall adopt the joint asymptotics approach in the sense that  $L$ ,  $p_x$ ,  $p$ ,  $T$  and  $N$  diverge to infinity jointly. We assume that initial observations  $y_{i,0}, y_{i,-1}, \dots, y_{i,1-L}$  are available for  $i = 1, \dots, N$ . We conjecture that (2.1) could also be extended with time effects. However, to

keep expressions and assumptions reasonably simple, we do not pursue this possibility in this work.

The three sources of high-dimensionality in (2.1) are  $p_x$ ,  $L$  and  $N$  as all of these can be increasing sequences. Sometimes one thinks of the number of lags,  $L$ , as being fixed and in that case only two sources remain. Next, (2.1) may be written more compactly as

$$y_i = Z_i' \alpha + \eta_i \iota + \varepsilon_i,$$

where  $Z_i := (z_{i,1}, \dots, z_{i,T})$  is a  $p \times T$  matrix,  $y_i := (y_{i,1}, \dots, y_{i,T})'$ ,  $\varepsilon_i := (\varepsilon_{i,1}, \dots, \varepsilon_{i,T})'$ , and  $\iota$  is a  $T \times 1$  vector of ones. Then, one can write

$$y = (Z \quad D) \begin{pmatrix} \alpha \\ \eta \end{pmatrix} + \varepsilon = \Pi \gamma + \varepsilon,$$

where  $Z := (Z_1, \dots, Z_N)'$ ,  $y := (y_1', \dots, y_N')'$  and  $\varepsilon := (\varepsilon_1', \dots, \varepsilon_N')'$ .  $\eta := (\eta_1, \dots, \eta_N)'$  contains the fixed effects,  $D := I_N \otimes \iota$ , and  $\Pi := (Z, D)$ . Finally,  $\gamma := (\alpha', \eta')'$  contains all  $p + N$  parameters of the model. Thus the dynamic panel data model (1.1) can be written more compactly as something resembling a linear regression model. There are several differences, however. First, blocks of rows in the data matrix  $\Pi$  may be highly dependent. Second, we shall see that the estimators of  $\alpha$  and  $\eta$  have markedly different properties as a result of the fact that the probabilistic properties of the blocks of a properly scaled version of the Gram matrix pertaining to  $\Pi$  are very different. Third, imposing weak sparsity on  $\alpha$  and  $\eta$  implies that the oracle inequalities which we use as a stepping stone towards inference do not follow directly from the technique in, e.g., Bickel et al. (2009). In fact, we do not get explicit expressions for the upper bounds but instead characterize them as solutions to certain quadratic equations in two variables.

## 2.2 Weak Sparsity and the Panel Lasso

As explained in the introduction, we treat  $\eta$  as a parameter to be estimated. However,  $\eta$  is not entirely unrestricted but assumed to be weakly sparse<sup>3</sup> in the sense

$$\sum_{i=1}^N |\eta_i|^\nu \leq E$$

for some  $0 < \nu < 1$  and  $E > 0$ . Weak sparsity does not require any of the fixed effects to be zero but instead restricts the “sum” of all the fixed effects.  $E$  can be large in the sense that it tends to infinity but the smaller it is, the sharper will our results be. It is appropriate to stress

---

<sup>3</sup>The term weakly sparse is borrowed from Negahban et al. (2012).

that the fixed effects cannot be entirely unrestricted. In particular we shall see that the oracle inequalities (Theorem 1) require  $E = O(\sqrt{N(\log(p \vee N))^{3\nu}/T^\nu})$  whereas the CLT and uniform inference (Theorems 2 and 3) require  $E = o(\sqrt{N(\log(p \vee N))^{3\nu}/T^\nu})$ . Notice that this is more restrictive than in classic low dimensional fixed effects models where  $E = O(N)$ . Thus, our framework also excludes many models of interest. We believe, however, that our results provide a useful first step towards uniform inference in high-dimensional dynamic panel data models. Lemma 1 in Appendix A provides sufficient conditions and discussion for the rate assumptions on  $E$  mentioned above to be satisfied, i.e. sufficient conditions for the fixed effects to be weakly sparse. Our strengthened assumption allows us to use the Lasso in the high-dimensional setting.

Note that the presence of many control variables in a high-dimensional model leaves less variation to be explained by the unobserved heterogeneities and these are therefore likely to be small in magnitude making the weak sparsity assumption reasonable. Thus, weak sparsity actually becomes more reasonable the larger the number of control variables is.

We also assume that  $\alpha$  is weakly sparse, i.e.,

$$\sum_{j=1}^p |\alpha_j|^\nu \leq E_1$$

for some  $E_1 > 0$ . Weak sparsity is a strict generalization of exact sparsity in the sense that if only  $s_\alpha$  elements of  $\alpha$  are non-zero and none of these exceeds a constant  $K$ , then  $\sum_{j=1}^p |\alpha_j|^\nu \leq s_\alpha K^\nu$ . Thus,  $E_1 = s_\alpha K^\nu$  works. Alternatively, exact sparsity of  $\alpha$  can be handled as the boundary case as  $\nu \rightarrow 0$  upon defining  $0^0 = 0$  such that  $E_1$  will equal the number of non-zero entries of  $\alpha$ . We shall see that the oracle inequalities (Theorem 1) require  $E_1 = O(\sqrt{N^{1-\nu}(\log(p \vee N))^{3\nu}/T^\nu})$  whereas the CLT and uniform inference (Theorems 2 and 3) require  $E_1 = o(\sqrt{N^{1-\nu}(\log(p \vee N))^{3\nu}/T^\nu})$ .

Note that we use the same  $\nu$  for both  $\alpha$  and  $\eta$  simply for convenience.

### 2.3 The Objective Function and Assumptions

Our starting point for inference is the minimiser  $\hat{\gamma} = (\hat{\alpha}', \hat{\eta}')'$  of the following panel Lasso objective function

$$L(\gamma) = \|y - \Pi\gamma\|^2 + 2\lambda\|\alpha\|_1 + 2\frac{\lambda}{\sqrt{N}}\|\eta\|_1. \quad (2.2)$$

As usual  $\lambda$  is a positive regularization sequence. Note that we penalize  $\alpha$  and  $\eta$  differently to reflect the fact that we have  $NT$  observations to estimate  $\alpha_j$  for  $j = 1, \dots, p$  while only  $T$  observations are available to estimate each  $\eta_i$ . Penalizing the fixed effects is not new and was already done in Koenker (2004) and Galvao and Montes-Rojas (2010) in a low dimensional

panel-quantile model. Furthermore, the penalization fits well with the weak sparsity assumption on the fixed effects and may increase efficiency of  $\hat{\alpha}$  as found in Galvao and Montes-Rojas (2010).

For practical implementation it is very convenient that we only have one penalty parameter  $\lambda$  instead of having separate penalty parameters for  $\alpha$  and  $\eta$ . The minimization problem can be solved easily as it simply corresponds to a weighted Lasso with known weights. However, the probabilistic analysis of the properly scaled Gram matrix is different from the one for the standard Lasso as it must be broken into several steps. We now turn to the assumptions imposed for our inferential procedure.

**Assumption 1.** For each  $T \in \mathbb{N}$ ,  $\{(x'_{i,1}, \dots, x'_{i,T}, \varepsilon'_i)\}_{i=1}^N$  is an independent sequence and

$$\mathbb{E}[\varepsilon_{i,t} | y_{i,t-1}, \dots, y_{i,1-L}, x_{i,t}, \dots, x_{i,1}] = 0 \quad \text{for} \quad i = 1, \dots, N, \quad t = 1, \dots, T.$$

Assumption 1 imposes independence across  $i = 1, \dots, N$  which is standard in the panel data literature, see e.g. Wooldridge (2010) or Arellano (2003). Note however, that we do not assume the data to be identically distributed across  $i = 1, \dots, N$ . Assumption 1 also implies, by iterated expectations, that the error terms form a martingale difference sequence with respect to the filtration generated by the variables in the above conditioning set and thus restricts the degree of dependence in the error terms across  $t$  (in particular, they are uncorrelated).<sup>4</sup> However, it still allows for considerable dependence over time, as higher moments than the first are not restricted. Furthermore, the error terms need not be identically distributed over time for any individual. Note that the increasing number of lags of  $y_{i,t}$  also whiten the error terms in practice. We also note that Assumption 1 does not rule out that the error terms are conditionally heteroskedastic. In panel data terminology, both lags of  $y_{i,t}$  and  $x_{i,t}$  are called *predetermined* or *weakly exogenous*. Finally, one can of course also include lags of the  $x_{i,t}$  in (1.1) as these are also weakly exogenous.

In order to introduce the next assumption define the scaled empirical Gram matrix

$$\Psi_N = S^{-1} \Pi' \Pi S^{-1} = \begin{pmatrix} \frac{1}{NT} Z' Z & \frac{1}{T\sqrt{N}} Z' D \\ \frac{1}{T\sqrt{N}} D' Z & I_N \end{pmatrix} \quad \text{where} \quad S = \begin{pmatrix} \sqrt{NT} I_p & 0 \\ 0 & \sqrt{T} I_N \end{pmatrix}$$

When  $p+N > NT$ ,  $\Psi_N$  is singular. However, to conduct inference it suffices that a compatibility type condition tailored to the panel data structure is satisfied. We shall see that it actually suffices that

$$\Psi = \begin{pmatrix} \Psi_Z & 0 \\ 0 & I_N \end{pmatrix} := \begin{pmatrix} \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \mathbb{E}[z_{i,t} z'_{i,t}] & 0 \\ 0 & I_N \end{pmatrix}.$$

---

<sup>4</sup>It can also be verified that  $\{\varepsilon_{i,t}\}_{t=1}^T$  forms a martingale difference sequence with respect to the natural filtration for all  $i = 1, \dots, N$ . This is because the  $\varepsilon_{i,t}$  are (linear) functions of the variables in the conditioning set in Assumption 1.



satisfies such a compatibility condition since  $\Psi_N$  will be shown to be close to  $\Psi$  in an appropriate sense. To be precise, define for integers  $r_1 \in \{1, \dots, p\}$  and  $r_2 \in \{1, \dots, N\}$

$$\kappa^2(A, r_1, r_2) := \min_{\substack{R_1 \subseteq \{1, \dots, p\}, |R_1| \leq r_1 \\ R_2 \subseteq \{1, \dots, N\}, |R_2| \leq r_2 \\ R := R_1 \cup (R_2 + p)}} \min_{\substack{\delta \in \mathbb{R}^{p+N} \setminus \{0\} \\ \|\delta_{R^c}\|_1 \leq 4\|\delta_R\|_1}} \frac{\delta' A \delta}{\frac{1}{r_1 + r_2} \|\delta_R\|_1^2}$$

which is reminiscent of the restricted eigenvalue condition in Bickel et al. (2009). Define the auxiliary parameters

$$\begin{aligned} \alpha_j^* &:= \alpha_j 1\{|\alpha_j| \geq \Xi_1\}, & J_1 &= \{j : \alpha_j^* \neq 0, j = 1, \dots, p\}, & s_1 &:= |J_1|, \\ \eta_i^* &:= \eta_i 1\{|\eta_i| \geq \Xi_2\}, & J_2 &= \{i : \eta_i^* \neq 0, i = 1, \dots, N\}, & s_2 &:= |J_2|. \end{aligned}$$

for  $\Xi_1, \Xi_2 \geq 0$  the details of which will be made precise in Appendix B.

**Assumption 2.**  $\kappa^2 := \kappa^2(\Psi, s_1, s_2)$  is bounded away from zero.

Assumption 2 is rather innocent as it is trivially satisfied when the  $\Psi_Z$  is positive definite. Since  $\Psi_Z$  is the *population* second moment matrix of  $z_{i,t}$  this is a standard assumption. Compatibility type conditions are typical in the literature and various versions and their inter-relationship have been investigated in van de Geer et al. (2009).

**Assumption 3.** *There exist absolute positive constants  $C$  and  $K$  such that*

- (a)  $\varepsilon_{i,t}$  are uniformly subgaussian; that is,  $\mathbb{P}(|\varepsilon_{i,t}| \geq \epsilon) \leq \frac{1}{2} K e^{-C\epsilon^2}$  for every  $\epsilon \geq 0$ ,  $i = 1, \dots, N$  and  $t = 1, \dots, T$ .
- (b)  $z_{i,t}$  are uniformly subgaussian; that is,  $\sup_{\|v\| \leq 1} \mathbb{P}(|v' z_{i,t}| \geq \epsilon) \leq \frac{1}{2} K e^{-C\epsilon^2}$  for every  $\epsilon \geq 0$ ,  $i = 1, \dots, N$  and  $t = 1, \dots, T$ .

In the context of the high-dimensional plain static regression model it is common practice to assume the error terms as well as the covariates to be subgaussian. However, this assumption is not as innocent in the context of the dynamic panel data model (1.1) as  $y_{i,t}$  is generated by the model and its properties are thus completely determined by those of  $x_{i,t}, \varepsilon_{i,t}$  as well as the parameters of the model. Lemma 2 in Appendix A shows that  $y_{i,t}$  is subgaussian if  $x_{i,t}$  and  $\varepsilon_{i,t}$  satisfy this property and the parameters are well-behaved. In particular, a wide class of (causal) stationary processes are included, though GARCH processes are excluded. Note also, that Assumption 3 imposes subgaussianity of the initial values  $y_{i,0}, \dots, y_{i,1-L}$  for all  $i = 1, \dots, N$ . Belloni et al. (2012) have analysed high-dimensional IV models without imposing subgaussianity by utilising moderate deviation inequalities for self-normalised sums of independent random variables. However, the dynamic panel data setting considered here induces a complicated

dependence structure such that these tools are not readily applicable. Further examples of papers which do not rely on subgaussianity are Belloni et al. (2014) and Caner and Kock (2014) both of which work in a setting of independent observations.

## 2.4 The Oracle Inequalities

With the above assumptions in place we are ready to state our first result. Set  $E_1 = N^{-\frac{\nu}{2}}E$ . The reason for such choice of  $E_1$  is to balance the terms of the upper bounds of the oracle inequalities. Define  $\mathcal{F}(\nu, E) := \{\alpha \in \mathbb{R}^p : \sum_{j=1}^p |\alpha_j|^\nu \leq N^{-\frac{\nu}{2}}E\} \times \{\eta \in \mathbb{R}^N : \sum_{i=1}^N |\eta_i|^\nu \leq E\}$ .

**Theorem 1 (Oracle inequalities).** *Let Assumptions 1 - 3 hold. Choose  $\lambda = \sqrt{4MNT(\log(p \vee N))^3}$  for some  $M > 0$  ( $M$  does not depend on any other constant). Then the following inequalities are valid with probability at least*

$$1 - Ap^{1-AM^{1/3}} - AN^{1-AM^{1/3}} - A(p^2 + pN) \exp\left(-B \left\{ \frac{N}{E^2 [(\log(p \vee N))^3/T]^{-\nu}} \right\}^{1/3}\right)$$

provided  $E = O(\sqrt{N(\log(p \vee N))^{3\nu}/T^\nu})$ , where  $A$  and  $B$  are generic positive constants.

$$\begin{aligned} \frac{1}{NT} \|\Pi(\hat{\gamma} - \gamma)\|^2 &\leq \left(\frac{240}{\kappa^2} + 40\right) \frac{\lambda}{\sqrt{N}NT} E \left(\frac{\lambda}{\sqrt{NT}}\right)^{1-\nu} \\ \|\hat{\alpha} - \alpha\|_1 &\leq \left(\frac{240}{\kappa^2} + 40\right) \frac{1}{\sqrt{N}} E \left(\frac{\lambda}{\sqrt{NT}}\right)^{1-\nu} \\ \|\hat{\eta} - \eta\|_1 &\leq \left(\frac{240}{\kappa^2} + 40\right) E \left(\frac{\lambda}{\sqrt{NT}}\right)^{1-\nu}. \end{aligned}$$

Theorem 1 provides oracle inequalities for the prediction error as well as the estimation error of the parameter vectors. Moreover, the above bounds are valid uniformly over  $\mathcal{F}(\nu, E)$ . While these bounds are of independent interest, we primarily use them as means towards our ultimate end of conducting (joint) inference on  $\alpha$  and  $\eta$ . We stress that the bounds in Theorem 1 are finite sample bounds; they hold for any fixed values of  $N$  and  $T$ . The novel feature of our oracle inequalities is that  $E$ , the "size" of  $\eta$  or  $\alpha$ , is allowed to grow even when we want the upper bound of  $\|\hat{\eta} - \eta\|_1$  or  $\|\hat{\alpha} - \alpha\|_1$  go to zero. The special case of exact sparsity of  $\alpha$  corresponds to  $\nu = 0$  upon defining  $0^0 = 0$  and  $N^{-\frac{\nu}{2}}E = E_1$  being the sparsity index of  $\alpha$ .

Theorem 1 or its corollary (Corollary 1 in Appendix B) imply that for fixed  $T$ , it is possible for  $\|\hat{\alpha} - \alpha\|_1 \xrightarrow{p} 0$  by allowing  $N \rightarrow \infty$  only, provided that  $E$  does not grow too fast. Hence our panel Lasso estimator  $\hat{\alpha}$  does not run into the problem of Nickell bias (i.e., inconsistency) (see Nickell (1981)). This fact is due to the weak sparsity assumption we imposed on the model and use of an estimator consistent with this assumption. This framework does, however, exclude many models of interest.

We also note that the oracle inequalities are not obtained in an entirely standard manner as the weak sparsity in dynamic panel data models calls for a different proof technique which yields the upper bounds as solutions to certain quadratic equations. Furthermore, we remark that in analogy to oracle inequalities in the plain linear regression model the number of covariates in  $x_{i,t}$  ( $p_x$ ) may increase at an exponential rate in  $NT$  without hindering the right hand sides of the oracle inequalities in converging to zero. Finally, we do not assume independence across  $t = 1, \dots, T$  for any individual thus extending the standard probabilistic analysis as well. Instead we use concentration inequalities for martingales to obtain bounds almost as sharp as in the completely independent case. However, there is a small cost of not assuming independence:  $\lambda$  has a factor  $\sqrt{(\log(p \vee N))^3}$  instead of  $\sqrt{\log(p \vee N)}$ . If one restricts the dependence structure of  $\{x_{i,t}\}_{t=1}^T$  for every  $i = 1, \dots, N$  to be, e.g., strongly mixing then one can use concentration inequalities for mixing processes such as in Merlevède et al. (2011). Restricting the dependence structure this way will allow  $E$  to increase faster. The focus on the  $\ell_1$ -norm in the oracle inequalities for  $\alpha$  and  $\eta$  is due to the fact that an upper bound in this norm will be particularly useful when developing our uniformly valid inference procedure in the following sections.

### 3 Inference

In this section we show how to conduct inference on  $\gamma$  and first discuss how desparsification as proposed in van de Geer et al. (2014) works in our context.

#### 3.1 The De-sparsified Lasso Estimator $\tilde{\gamma}$

First, observe that  $L(\gamma)$  in (2.2) is convex in  $\gamma$  and in order for  $\hat{\gamma}$  to be a minimiser of  $L$ , 0 must belong to the subdifferential of  $L(\gamma)$  at  $\hat{\gamma}$ , i.e.

$$0 \in \partial L(\hat{\gamma}) = \begin{pmatrix} -2Z'(y - \Pi\hat{\gamma}) + 2\lambda\hat{\kappa}_1 \\ -2D'(y - \Pi\hat{\gamma}) + 2\frac{\lambda}{\sqrt{N}}\hat{\kappa}_2 \end{pmatrix}$$

where  $\hat{\kappa}_1$  and  $\hat{\kappa}_2$  are  $p \times 1$  and  $N \times 1$  vectors, respectively, such that  $\hat{\kappa}_{1j} \in [-1, 1]$  with  $\hat{\kappa}_{1j} = \text{sgn}(\hat{\alpha}_j)$  if  $\hat{\alpha}_j \neq 0$  for  $j = 1, \dots, p$ . Similarly,  $\hat{\kappa}_{2i} \in [-1, 1]$  with  $\hat{\kappa}_{2i} = \text{sgn}(\hat{\eta}_i)$  if  $\hat{\eta}_i \neq 0$  for  $i = 1, \dots, N$ . Hence,

$$-\Pi'(y - \Pi\hat{\gamma}) + \begin{pmatrix} \lambda\hat{\kappa}_1 \\ \frac{\lambda}{\sqrt{N}}\hat{\kappa}_2 \end{pmatrix} = 0. \quad (3.1)$$

Using that  $y = \Pi\gamma + \varepsilon$  and multiplying by  $S^{-1}$  from the left yields

$$\Psi_N S(\hat{\gamma} - \gamma) + S^{-1} \begin{pmatrix} \lambda\hat{\kappa}_1 \\ \frac{\lambda}{\sqrt{N}}\hat{\kappa}_2 \end{pmatrix} = S^{-1}\Pi'\varepsilon.$$

In order to derive the limiting distribution of  $S(\hat{\gamma} - \gamma)$  one would usually proceed by isolating  $S(\hat{\gamma} - \gamma)$  which implies inverting  $\Psi_N$ . However, when  $p + N > NT$ ,  $\Psi_N$  is not invertible. The idea of van de Geer et al. (2014) and Javanmard and Montanari (2013) is to circumvent this problem by using an approximate inverse of  $\Psi_N$  and controlling the asymptotic approximation error. Suppose that a matrix  $\hat{\Theta}$  is a reasonable approximation to the inverse of  $\Psi_N$ . We shall explicitly construct  $\hat{\Theta}$  in the next section. Then we may write

$$\hat{\gamma} = \gamma - S^{-1}\hat{\Theta}S^{-1} \begin{pmatrix} \lambda\hat{\kappa}_1 \\ \frac{\lambda}{\sqrt{N}}\hat{\kappa}_2 \end{pmatrix} + S^{-1}\hat{\Theta}S^{-1}\Pi'\varepsilon - S^{-1}\Delta$$

where  $\Delta := (\hat{\Theta}\Psi_N - \mathbf{I})S(\hat{\gamma} - \gamma)$  is the error resulting from using an approximate inverse  $\hat{\Theta}$  of  $\Psi_N$  as opposed to an exact inverse. The term  $S^{-1}\hat{\Theta}S^{-1} \begin{pmatrix} \lambda\hat{\kappa}_1 \\ \frac{\lambda}{\sqrt{N}}\hat{\kappa}_2 \end{pmatrix}$  in the above display is the bias incurred by  $\hat{\gamma}$  due to shrinkage of the parameters in (2.2). As this bias term is known one may add it back to  $\hat{\gamma}$  in order to define the debiased estimator

$$\tilde{\gamma} = \hat{\gamma} + S^{-1}\hat{\Theta}S^{-1} \begin{pmatrix} \lambda\hat{\kappa}_1 \\ \frac{\lambda}{\sqrt{N}}\hat{\kappa}_2 \end{pmatrix} = \gamma + S^{-1}\hat{\Theta}S^{-1}\Pi'\varepsilon - S^{-1}\Delta$$

The new estimator  $\tilde{\gamma}$  is no longer sparse as it includes a bias correction term to the sparse Lasso estimator  $\hat{\gamma}$ . Therefore, we will also refer to it as the *de-sparsified* Lasso estimator in the dynamic panel context. For any  $(p+N) \times 1$  vector  $\rho$  with  $\|\rho\| = 1$  we shall study the asymptotic behaviour of

$$\rho'S(\tilde{\gamma} - \gamma) = \rho'\hat{\Theta}S^{-1}\Pi'\varepsilon - \rho'\Delta. \quad (3.2)$$

A central limit theorem for  $\rho'\hat{\Theta}S^{-1}\Pi'\varepsilon$  as well as asymptotic negligibility of  $\rho'\Delta$  will yield asymptotically gaussian inference. Furthermore, we shall provide a uniformly consistent estimator of the asymptotic variance of  $\rho'\hat{\Theta}S^{-1}\Pi'\varepsilon$  even in the presence of conditional heteroskedasticity. A leading special case of (3.2) is when one is only interested in the asymptotic distribution of  $\tilde{\gamma}_j$  corresponding to  $\rho = e_j$  being the  $j$ th basis vector of  $\mathbb{R}^{p+N}$ . In general, we will be interested in the asymptotic distribution of a subset  $H \subseteq \{1, \dots, p+N\}$  of the indices of  $\gamma$  with cardinality  $h$  and shall show that asymptotically honest (uniformly valid) gaussian inference is possible in the presence of heteroskedasticity even for  $h \rightarrow \infty$  and  $H$  involving elements of  $\alpha$  and  $\eta$ .

### 3.2 Construction of $\hat{\Theta}$

As is clear from the discussion above we need a good choice for  $\hat{\Theta}$ . In particular we shall show that

$$\hat{\Theta} = \begin{pmatrix} \hat{\Theta}_Z & 0 \\ 0 & I_N \end{pmatrix}$$

works well. Here  $\hat{\Theta}_Z$  will be constructed using nodewise regressions as in van de Geer et al. (2014) and we show that this is possible even when the rows of  $Z$  are not independent and identically distributed. The construction of  $\hat{\Theta}_Z$  parallels the one in van de Geer et al. (2014) to a high extent but importantly for our context we do not need the rows of  $\Psi_Z^{-1}$  to be sparse for the nodewise regressions to work well. The details of the construction of  $\hat{\Theta}_Z$  are given in Appendix G.

### 3.3 Asymptotic Properties of the Approximate Inverse

In order to show that  $\rho' \hat{\Theta} S^{-1} \Pi' \varepsilon$  is asymptotically Gaussian one needs to understand the limiting behaviour of  $\hat{\Theta}$  constructed above. We show that  $\hat{\Theta}$  is close to

$$\Theta = \begin{pmatrix} \Theta_Z & 0 \\ 0 & I_N \end{pmatrix} := \begin{pmatrix} \Psi_Z^{-1} & 0 \\ 0 & I_N \end{pmatrix}$$

in an appropriate sense. To this end, note that by Yuan (2010)

$$\Theta_{Z,j,j} = \left[ \Psi_{Z,j,j} - \Psi_{Z,j,-j} \Psi_{Z,-j,-j}^{-1} \Psi_{Z,-j,j} \right]^{-1} \quad \text{and} \quad \Theta_{Z,j,-j} = -\Theta_{Z,j,j} \Psi_{Z,j,-j} \Psi_{Z,-j,-j}^{-1}, \quad (3.3)$$

where  $\Theta_{Z,j,j}$  is the  $j$ th diagonal entry of  $\Theta_Z$ ,  $\Theta_{Z,j,-j}$  is the  $1 \times (p-1)$  vector obtained by removing the  $j$ th entry of the  $j$ th row of  $\Theta_Z$ ,  $\Psi_{Z,-j,-j}$  is the submatrix of  $\Psi_Z$  with the  $j$ th row and column removed,  $\Psi_{Z,j,-j}$  is the  $j$ th row of  $\Psi_Z$  with its  $j$ th entry removed,  $\Psi_{Z,-j,j}$  is the  $j$ th column of  $\Psi_Z$  with its  $j$ th entry removed. Next, let  $z_{i,t,j}$  be the  $j$ th element of  $z_{i,t}$  and  $z_{i,t,-j}$  be all elements except the  $j$ th. Define the  $(p-1) \times 1$  vector

$$\phi_j := \underset{\delta \in \mathbb{R}^{p-1}}{\operatorname{argmin}} \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \mathbb{E}[z_{i,t,j} - z'_{i,t,-j} \delta]^2$$

such that

$$\phi_j = \left( \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \mathbb{E}[z_{i,t,-j} z'_{i,t,-j}] \right)^{-1} \left( \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \mathbb{E}[z_{i,t,-j} z_{i,t,j}] \right) = \Psi_{Z,-j,-j}^{-1} \Psi_{Z,-j,j}. \quad (3.4)$$

Therefore,  $\Theta_{Z,j,-j} = -\Theta_{Z,j,j} \phi'_j$  showing that  $\Theta_{Z,j,-j}$  and  $\phi'_j$  only differ by a multiplicative constant. In particular,  $j$ th row of  $\Theta_Z$  is exactly sparse if and only if  $\phi_j$  is exactly sparse. More

generally, we shall exploit below that weak sparsity of one implies weak sparsity of the other.

Furthermore, defining  $\zeta_{j,i,t} := z_{i,t,j} - z'_{i,t,-j}\phi_j$  we may write

$$z_{i,t,j} = z'_{i,t,-j}\phi_j + \zeta_{j,i,t}, \quad \text{for } i = 1, \dots, N, \quad t = 1, \dots, T.$$

where by the definition of  $\phi_j$

$$\frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \mathbb{E}[z_{i,t,-j}\zeta_{j,i,t}] = 0. \quad (3.5)$$

Thus, in light of Theorem 1, it is sensible that the Lasso estimator  $\hat{\phi}_j$  defined in (13.1) is close to the population regression coefficients  $\phi_j$  (we shall make this more formal in Appendix C). Next, defining

$$\tau_j^2 := \mathbb{E} \left[ \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T (z_{i,t,j} - z'_{i,t,-j}\phi_j)^2 \right] = \Psi_{Z,j,j} - \Psi_{Z,j,-j} \Psi_{Z,-j,-j}^{-1} \Psi_{Z,-j,j} = \frac{1}{\Theta_{Z,j,j}}$$

observe  $\Theta_{Z,j,-j} = -\phi'_j/\tau_j^2$ . Thus, we can write  $\Theta_Z = T^{-1}C$  where  $T = \text{diag}(\tau_1^2, \dots, \tau_p^2)$  and  $C$  is defined similarly to  $\hat{C}$  but with  $\phi_j$  replacing  $\hat{\phi}_j$  for  $j = 1, \dots, p$ . Finally, let  $\Theta_{Z,j}$  denote the  $j$ th row of  $\Theta_Z$  written as a column vector. In Lemma 9 in Appendix C we will see that  $\hat{\phi}_j$  and  $\hat{\tau}_j^2$  are close to  $\phi_j$  and  $\tau_j^2$ , respectively such that  $\hat{\Theta}_{Z,j}$  is close to  $\Theta_{Z,j}$  which is the desired control of  $\hat{\Theta}_{Z,j}$ . Write  $\rho = (\rho'_1, \rho'_2)'$  with  $\|\rho\| = 1$ , where  $\rho_1 \in \mathbb{R}^p$  and  $\rho_2 \in \mathbb{R}^N$ . Hence define

$$H = H_1 \cup (H_2 + p) := \{j : \rho_{1j} \neq 0\} \cup (\{i : \rho_{2i} \neq 0\} + p),$$

with  $|H_1| = h_{1,N} = h_1$ ,  $|H_2| = h_{2,N} = h_2$  and  $|H| = h = h_1 + h_2$ . In dynamic panel data models it may not be reasonable to assume that the rows of the inverse second moment matrix  $\Psi_Z^{-1} = \Theta_Z$ , i.e.  $\Theta_{Z,j}$  are sparse. Paralleling Section 2.2 we shall instead assume that the  $\Theta_{Z,j}$  are weakly sparse and assume that

$$\sum_{k=1}^{p-1} |\phi_{j,k}|^\vartheta = \tau_j^{2\vartheta} \sum_{l \neq j}^p |\Theta_{Z,j,l}|^\vartheta \leq G_j \quad (3.6)$$

for some  $0 < \vartheta < 1$  and  $G_j > 0$ . Define  $\bar{G} := \max_{j \in H_1} G_j$ .

**Assumption 4.**

(a) *mineval*( $\Psi_Z$ ) is uniformly bounded away from zero and *maxeval*( $\Psi_Z$ ) is uniformly bounded from above.

(b)  $\bar{G}\lambda_{node}^{1-\vartheta} = O(1)$ .

(c) There exist positive constants  $C$  and  $K$  such that  $\zeta_{j,i,t}$  are uniformly subgaussian; that is,  $\mathbb{P}(|\zeta_{j,i,t}| \geq \epsilon) \leq \frac{1}{2}Ke^{-C\epsilon^2}$  for every  $\epsilon > 0$ ,  $i = 1, \dots, N$ ,  $t = 1, \dots, T$  and  $j = 1, \dots, p$ .

Assumption 4(a) is standard and strengthens Assumption 2 slightly. Recall that the population matrix  $\Psi_Z$  can have full rank even when the empirical counterpart  $\Psi_N$  has rank zero – which it has when  $p+N > NT$ . Note that Assumption 4(a) implies that  $\tau_j^2$  is uniformly bounded away from zero as  $\tau_j^2 = 1/\Theta_{Z,j,j} \geq 1/\text{maxeval}(\Theta_Z) = \text{mineval}(\Psi_Z)$ . Similarly,  $\tau_j^2 \leq \text{maxeval}(\Psi_Z)$  implying that  $\tau_j^2$  is bounded in (3.6). Therefore, weak sparsity of  $\phi_{j,k}$  translates into weak sparsity of the rows of  $\Theta$ . Notice that we generalize the cross sectional results of van de Geer et al. (2014) by not imposing the inverse covariance (second moment matrix) of  $z_{i,t}$  to have sparse rows. When  $z_{i,t}$  is gaussian exact sparsity of  $\Psi_Z^{-1}$  is related to the notion of conditional independence: the  $(j,k)$ th entry of  $\Psi_Z^{-1}$  being zero is equivalent to  $z_{i,t,j}$  being independent of  $z_{i,t,k}$  conditional on the remaining variables in  $z_{i,t}$ . This is hard to justify in dynamic panel data models. First, it does not sound reasonable for  $x_{i,t}$ s to be mostly conditionally independent given the lagged variables. Second, adjacent lagged variables  $y_{i,t-l}$  and  $y_{i,t-l-1}$  ( $l = 1, \dots, L+1$ ) are not independent even after conditioning on all the other variables in  $z_{i,t}$ . Thus, it is important to relax the exact sparsity assumption on the rows of  $\Theta_Z$  in the context of dynamic panel data models.

Part (b) restricts the rate of growth of  $\tilde{G}$ . As we shall choose  $\lambda_{node} \asymp \sqrt{\frac{\log^3(p)}{N}}$  it implies in particular that  $\tilde{G} = O((N/\log^3(p))^{\frac{1-\vartheta}{2}})$ . Part (c) imposes subgaussianity on the error terms from the nodewise regressions.

### 3.4 The Asymptotic Distribution of $\tilde{\gamma}$

In this section we formalise the discussion in Section 3.1 as Theorem 2. To this end, define

$$\Sigma_{\Pi\varepsilon} = \mathbb{E}(S^{-1}\Pi'\varepsilon\varepsilon'\Pi S^{-1}) = \begin{pmatrix} \mathbb{E}[Z'\varepsilon\varepsilon'Z/(NT)] & \mathbb{E}[Z'\varepsilon\varepsilon'D/(\sqrt{NT})] \\ \mathbb{E}[D'\varepsilon\varepsilon'Z/(\sqrt{NT})] & \mathbb{E}[D'\varepsilon\varepsilon'D/T] \end{pmatrix} = \begin{pmatrix} \Sigma_{1,N} & \Sigma_{2,N} \\ \Sigma'_{2,N} & \Sigma_{3,N} \end{pmatrix}.$$

and note that

$$\Sigma_{1,N} = \mathbb{E} \left[ \frac{1}{NT} \sum_{i=1}^N \sum_{j=1}^N Z_i \varepsilon_i \varepsilon'_j Z'_j \right] = \frac{1}{NT} \sum_{i=1}^N \mathbb{E} [Z_i \varepsilon_i \varepsilon'_i Z'_i] = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \mathbb{E} [\varepsilon_{i,t}^2 z_{i,t} z'_{i,t}],$$

where the second and third equality both follow from Assumption 1. Likewise,  $\Sigma_{3,N} = \frac{1}{T} \sum_{i=1}^N \mathbb{E} [d_i \varepsilon_i \varepsilon'_i d'_i] = \frac{1}{T} \sum_{i=1}^N \sum_{t=1}^T \mathbb{E} [\varepsilon_{i,t}^2 d_{i,t} d'_{i,t}] = \text{diag}(\frac{1}{T} \sum_{t=1}^T \mathbb{E}[\varepsilon_{1,t}^2], \dots, \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\varepsilon_{N,t}^2])$ , where  $d'_i$  is the  $i$ th  $T \times N$  block of  $D$ , and  $d_{i,t}$  is a  $N \times 1$  zero vector with the  $i$ th entry replaced by 1. In the same manner,  $\Sigma_{2,N} = \frac{1}{\sqrt{NT}} \sum_{i=1}^N \mathbb{E} [z_i \varepsilon_i \varepsilon'_i d'_i] = \frac{1}{\sqrt{NT}} \sum_{i=1}^N \sum_{t=1}^T \mathbb{E} [\varepsilon_{i,t}^2 z_{i,t} d'_{i,t}]$ . In words,  $\Sigma_{2,N}$  is a  $p \times N$  matrix with its  $i$ th column being  $\frac{1}{\sqrt{NT}} \sum_{t=1}^T \mathbb{E} [z_{i,t} \varepsilon_{i,t}^2]$ . Finally, motivated by the above, define the feasible sample counterpart of  $\Sigma_{\Pi\varepsilon}$  as

$$\hat{\Sigma}_{\Pi\varepsilon} = \begin{pmatrix} \hat{\Sigma}_{1,N} & \hat{\Sigma}_{2,N} \\ \hat{\Sigma}'_{2,N} & \hat{\Sigma}_{3,N} \end{pmatrix} := \begin{pmatrix} \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \hat{\varepsilon}_{i,t}^2 z_{i,t} z'_{i,t} & \frac{1}{\sqrt{NT}} \sum_{i=1}^N \sum_{t=1}^T \hat{\varepsilon}_{i,t}^2 z_{i,t} d'_{i,t} \\ \frac{1}{\sqrt{NT}} \sum_{i=1}^N \sum_{t=1}^T \hat{\varepsilon}_{i,t}^2 d_{i,t} z'_{i,t} & \frac{1}{T} \sum_{i=1}^N \sum_{t=1}^T \hat{\varepsilon}_{i,t}^2 d_{i,t} d'_{i,t} \end{pmatrix},$$

where  $\hat{\varepsilon}_{i,t} := y_{i,t} - z'_{i,t}\hat{\alpha} - \hat{\eta}_i$ . One could also consider constructing  $\hat{\varepsilon}_{i,t}$  based on the de-sparsified estimates. However, this would require running the nodewise regressions for all variables and not only those pertaining to the coefficients in the hypothesis being tested resulting in a much more computationally demanding procedure. The following assumptions are imposed to establish the validity of asymptotically gaussian inference of our procedure.

**Assumption 5.** *Let  $\tilde{p} := p \vee N \vee T$  and assume*

(a)

$$\frac{(h_1 \vee h_2 1\{h_1 \neq 0\})^2 \bar{G}^2 \left[ \frac{(\log p)^3}{N} \right]^{-\vartheta} (\log \tilde{p})^7}{N} = o(1), \quad \frac{(\log(N \vee T))^3 1\{h_2 \neq 0\}}{T} = o(1).$$

(b)

$$\frac{\left( h_1^2 \bar{G}^2 \left[ \frac{(\log p)^3}{N} \right]^{-\vartheta} \vee N h_2^2 \right) \left[ E \left( \frac{(\log(p \vee N))^3}{T} \right)^{-\nu/2} \right] (\log \tilde{p})^5}{NT} = o(1).$$

(c)

$$\frac{(h_1 \vee h_2) \left[ \left( \bar{G} \left[ \frac{(\log p)^3}{N} \right]^{-\vartheta/2} \vee (\log \tilde{p})^2 \right) 1\{h_1 \neq 0\} \vee 1\{h_2 \neq 0\} \right] \left[ E^2 \left( \frac{(\log(p \vee N))^3}{T} \right)^{-\nu} \right] (\log \tilde{p})^4}{N} = o(1).$$

(d) *mineval( $\Sigma_{\Pi_E}$ ) is uniformly bounded away from zero and maxeval( $\Sigma_{1,N}$ ) is uniformly bounded from above.*

Assumption 5 is slightly stronger than what we actually need in order to prove Theorem 2 but it is less cluttered in terms of notation. Assumption 5 restricts the rate at which  $p$ ,  $T$ ,  $E$ ,  $\bar{G}$ ,  $h_1$  and  $h_2$  are allowed to increase as none of these are assumed to be bounded. First, note that  $p = L + p_x$  only enters through its logarithm. Thus, we can allow for very high-dimensional models. Furthermore,  $h_1$  as well as  $h_2$  are allowed to increase with the sample size such that hypotheses of an increasing dimension involving  $\alpha$  and  $\eta$  can be tested. In the classic setting where one is only interested in testing hypotheses on  $\alpha$  one has that  $h_2 = 0$  and Assumption 5 simplifies. The case of hypotheses only involving the fixed effects  $\eta$  corresponds to  $h_1 = 0$  and again the assumptions simplify. We also note that Assumption 5 requires  $\bar{G}$  and  $h_1$  necessarily to be  $o(N^{\frac{1-\vartheta}{2}})$ ,  $E$  necessarily to be  $o(\sqrt{N(\log(p \vee N))^{3\nu}/T^\nu})$ , and  $h_2$  necessarily to be  $o(T^{1/2})$ . The restrictions on  $h_1$  and  $h_2$ , i.e. the number of common coefficients and fixed effects involved in the hypothesis, thus clearly encompass the classical setting where one tests only a fixed number of parameters ( $h_1$  and  $h_2$  fixed). Assumption 5 is satisfied if, for example,  $p = N, T = N^{1/2}, \nu = \vartheta = 0.5, E = N^{1/6}$  and  $\bar{G} = N^{1/7}$ . Thus, while we allow these quantities to diverge, the rate at which they do so can not be too fast



**Theorem 2.** *Let Assumptions 1, 3, 4, and 5 be satisfied. If, furthermore,  $\{\varepsilon_{i,t}\}_{t=1}^T$  is an independent sequence for all  $i = 1, \dots, N$ , then*

$$\frac{\rho' S(\tilde{\gamma} - \gamma)}{\sqrt{\rho' \hat{\Theta} \hat{\Sigma}_{\Pi\varepsilon} \hat{\Theta}' \rho}} \xrightarrow{d} N(0, 1), \quad (3.7)$$

where  $\rho = (\rho'_1, \rho'_2)'$  is a  $(p + N) \times 1$  vector, with  $\|\rho\| = 1$ ,  $\rho_1 \in \mathbb{R}^p$  and  $\rho_2 \in \mathbb{R}^N$ . Moreover,

$$\sup_{\gamma \in \mathcal{F}(\nu, E)} |\rho' \hat{\Theta} \hat{\Sigma}_{\Pi\varepsilon} \hat{\Theta}' \rho - \rho' \Theta \Sigma_{\Pi\varepsilon} \Theta' \rho| = o_p(1). \quad (3.8)$$

Finally, for every fixed set  $H \subseteq \{1, \dots, N + p\}$  with cardinality  $h$ , we have

$$[S_H(\tilde{\gamma}_H - \gamma_H)]' (\hat{\Theta} \hat{\Sigma}_{\Pi\varepsilon} \hat{\Theta}')^{-1}_H [S_H(\tilde{\gamma}_H - \gamma_H)] \xrightarrow{d} \chi^2_h. \quad (3.9)$$

Theorem 2 provides sufficient conditions under which our procedure allows for asymptotically gaussian inference. We stress again that hypotheses involving an increasing number of parameters can be tested and that the total number of parameters in the model may be much larger than the sample size. Furthermore, the error terms are allowed to be conditionally heteroskedastic and we provide a consistent estimator of the asymptotic covariance matrix even for the case of hypotheses involving an increasing number of parameters. Indeed, this estimator converges uniformly over  $\mathcal{F}(\nu, E)$  even for high-dimensional covariance matrices – a property we use in Theorem 3 to establish the honesty (uniform validity) over  $\mathcal{F}(\nu, E)$  of confidence intervals based on (3.7). van de Geer et al. (2014) have derived similar results in the setting of the homoskedastic linear cross sectional model for the case of inference on a low-dimensional parameter. Thus, our results can be seen as an extension to dynamic panel data models as well as to inference involving many parameters. We stress again that we relax their assumption of the inverse covariance matrix  $\Theta_Z$  being exactly sparse which is important in dynamic models like ours. Furthermore, relaxing the homoskedasticity assumption is important as volatility is known to vary over time in dynamic models, see e.g. Engle (1982), and the conditional volatility often depends on the state of the process. Theorem 2 is also related to Belloni et al. (2015) who consider inference in static panel data models for a low-dimensional parameter of interest.

The classic setup where one is only interested in inference on  $\alpha$  corresponds to  $\rho_2 = 0$  such that  $\sqrt{NT} \rho'_1 (\tilde{\alpha} - \alpha)$  is asymptotically gaussian with variance equal to the limit of  $\rho'_1 \Theta_Z \Sigma_{1,N} \Theta'_Z \rho_1$  (assumed to exist for illustration). If, furthermore,  $\varepsilon_{i,t}$  is homoskedastic with variance  $\sigma^2$  and independent of  $z_{i,t}$  for all  $i = 1, \dots, N$  and  $t = 1, \dots, T$ , it follows from the definition of  $\Sigma_{1,N}$  that this variance equals the limit of  $\sigma^2 \rho'_1 \Theta_Z \rho_1 = \sigma^2 \rho'_1 \Psi_Z^{-1} \rho_1$ . The leading special case where one is interested in testing a hypothesis on the  $j$ 'th entry of  $\alpha$  corresponds to  $\rho_1 = e_j$ . Similar reasoning shows that in the case where one is testing hypotheses involving

fixed effects only, corresponding to  $\rho_1 = 0$ , one has that  $\rho'_2 \sqrt{T} (\tilde{\eta} - \eta)$  is asymptotically gaussian with variance  $\sigma^2$ . This simple form of the variance follows from the asymptotic independence of the components of  $\tilde{\eta}$ . Note that the different rates of convergence for  $\tilde{\alpha}$  and  $\tilde{\eta}$  are in accordance with Theorem 1.

(3.9) is a straightforward consequence of (3.7) and reveals that classical  $\chi^2$  inference can be carried out in the usual manner. Thus, asymptotically valid  $\chi^2$ -inference can be performed in order to test a hypothesis involving  $h$  parameters. Wald tests of general restrictions of the type  $H_0 : g(\gamma) = 0$  (where  $g : \mathbb{R}^{p+N} \rightarrow \mathbb{R}^h$  is differentiable in an open neighborhood around  $\gamma$  and has derivative matrix of rank  $h$ ) can now also be constructed in the usual manner, see e.g. Davidson (2000) Chapter 12, even when  $p + N > NT$  which has hitherto been impossible.

Finally, the independence assumption on  $\varepsilon_{i,t}$  across  $t$  is needed only if one tests hypotheses involving  $\{\eta_i\}_{i=1}^N$  ( $h_2 \neq 0$ ). Weaker assumptions on the error terms, such as strong mixing, are possible at the expense of more involved expressions but will not be pursued here.<sup>5</sup>

## 4 Honest Confidence Intervals

In this section we show that the confidence bands based on (3.7) are honest (uniformly valid) and contract at the optimal rate. The precise result is contained in the following theorem.

**Theorem 3.** *Let Assumptions 1, 3, 4, and 5 be satisfied. If, furthermore,  $\{\varepsilon_{i,t}\}_{t=1}^T$  is an independent sequence for all  $i = 1, \dots, N$ , then, for all  $\rho \in \mathbb{R}^{p+N}$  with  $\|\rho\| = 1$ ,*

$$\sup_{t \in \mathbb{R}} \sup_{\gamma \in \mathcal{F}(\nu, E)} \left| \mathbb{P} \left( \frac{\rho' S(\tilde{\gamma} - \gamma)}{\sqrt{\rho' \hat{\Theta} \hat{\Sigma}_{\Pi \varepsilon} \hat{\Theta}' \rho}} \leq t \right) - \Phi(t) \right| = o(1), \quad (4.1)$$

where  $\Phi(\cdot)$  is the CDF of the standard normal distribution. Furthermore, define  $\tilde{\sigma}_{\alpha,j} := \sqrt{[\hat{\Theta}_Z \hat{\Sigma}_{1,N} \hat{\Theta}_Z]_{jj}}$  and  $\tilde{\sigma}_{\eta,i} := \sqrt{[\hat{\Sigma}_{3,N}]_{ii}}$  for  $j = 1, \dots, p$  and  $i = 1, \dots, N$ , respectively. Then,

$$\liminf_{N \rightarrow \infty} \inf_{\gamma \in \mathcal{F}(\nu, E)} \mathbb{P} \left( \alpha_j \in \left[ \tilde{\alpha}_j - z_{1-\delta/2} \frac{\tilde{\sigma}_{\alpha,j}}{\sqrt{NT}}, \tilde{\alpha}_j + z_{1-\delta/2} \frac{\tilde{\sigma}_{\alpha,j}}{\sqrt{NT}} \right] \right) \geq 1 - \delta, \quad (4.2)$$

$$\liminf_{N \rightarrow \infty} \inf_{\gamma \in \mathcal{F}(\nu, E)} \mathbb{P} \left( \eta_i \in \left[ \tilde{\eta}_i - z_{1-\delta/2} \frac{\tilde{\sigma}_{\eta,i}}{\sqrt{T}}, \tilde{\eta}_i + z_{1-\delta/2} \frac{\tilde{\sigma}_{\eta,i}}{\sqrt{T}} \right] \right) \geq 1 - \delta, \quad (4.3)$$

for  $j = 1, \dots, p$  and  $i = 1, \dots, N$ , respectively, where  $z_{1-\delta/2}$  is the  $1-\delta/2$  percentile of the standard normal distribution. Finally, letting  $\text{diam}([a, b]) = b - a$  be the length of an interval  $[a, b]$  in the

---

<sup>5</sup>The reason that the martingale difference assumption on  $\varepsilon_{i,t}$  across  $t$  (i.e., Assumption 1) is not enough for Theorem 2 is that when one is estimating the asymptotic variance  $\rho' \Theta \Sigma_{\Pi \varepsilon} \Theta' \rho$  in Theorem 2, one needs to invoke a concentration inequality for  $T^{-1} \sum_{t=1}^T (\varepsilon_{i,t}^2 - \mathbb{E}[\varepsilon_{i,t}^2])$  (see the second display above (10.25) in Appendix D) and squares of martingale differences need not be martingale differences.

real line, we have

$$\sup_{\gamma \in \mathcal{F}(\nu, E)} \text{diam} \left( \left[ \tilde{\alpha}_j - z_{1-\delta/2} \frac{\tilde{\sigma}_{\alpha,j}}{\sqrt{NT}}, \tilde{\alpha}_j + z_{1-\delta/2} \frac{\tilde{\sigma}_{\alpha,j}}{\sqrt{NT}} \right] \right) = O_p \left( \frac{1}{\sqrt{NT}} \right), \quad (4.4)$$

$$\sup_{\gamma \in \mathcal{F}(\nu, E)} \text{diam} \left( \left[ \tilde{\eta}_i - z_{1-\delta/2} \frac{\tilde{\sigma}_{\eta,i}}{\sqrt{T}}, \tilde{\eta}_i + z_{1-\delta/2} \frac{\tilde{\sigma}_{\eta,i}}{\sqrt{T}} \right] \right) = O_p \left( \frac{1}{\sqrt{T}} \right), \quad (4.5)$$

for  $j = 1, \dots, p$  and  $i = 1, \dots, N$ , respectively.

(4.1) reveals that the convergence to the normal distribution in Theorem 2 is uniform over  $\mathcal{F}(\nu, E)$ . Since the de-sparsified Lasso is not a sparse estimator this uniform convergence does not contradict the work of Leeb and Pötscher (2005). Next, (4.2) is a direct consequence of (4.1) and reveals that the de-sparsified Lasso produces confidence bands which are *honest* (uniform) over  $\mathcal{F}(\nu, E)$ . Honest confidence bands are important in practical applications of dynamic panel data models as they guarantee the existence of an  $N_0$ , not depending on  $\gamma \in \mathcal{F}(\nu, E)$ , such that  $\left[ \tilde{\alpha}_j - z_{1-\delta/2} \frac{\tilde{\sigma}_{\alpha,j}}{\sqrt{NT}}, \tilde{\alpha}_j + z_{1-\delta/2} \frac{\tilde{\sigma}_{\alpha,j}}{\sqrt{NT}} \right]$  covers  $\alpha_j$  with probability not much smaller than  $1 - \delta$ . Here the important point is that one and the same  $N_0$  guarantees this coverage, irrespective of the true value of  $\gamma \in \mathcal{F}(\nu, E)$ . On the other hand, pointwise consistent confidence bands only guarantee that

$$\inf_{\gamma \in \mathcal{F}(\nu, E)} \liminf_{N \rightarrow \infty} \mathbb{P} \left( \alpha_j \in \left[ \tilde{\alpha}_j - z_{1-\delta/2} \frac{\tilde{\sigma}_{\alpha,j}}{\sqrt{NT}}, \tilde{\alpha}_j + z_{1-\delta/2} \frac{\tilde{\sigma}_{\alpha,j}}{\sqrt{NT}} \right] \right) \geq 1 - \delta,$$

implying that the value of  $N$  needed in order to guarantee a coverage of close to  $1 - \delta$  may depend on the *unknown* true parameter. Thus, for some parameter values one may have to sample more data points to achieve the desired coverage than for others which is unfortunate as one does not know for which parameters this is the case. An honest confidence set  $S_N$  for  $\alpha_j$  can of course trivially be obtained by setting  $S_N = \mathbb{R}$ . However, this is clearly not very informative and therefore (4.4) is reassuring as it guarantees that the length of the honest confidence interval contracts at the optimal rate. In particular, the confidence bands are uniformly narrow over  $\mathcal{F}(\nu, E)$  in the sense that for any  $\epsilon > 0$  there exists an  $M > 0$  such that  $\text{diam} \left( \left[ \tilde{\alpha}_j - z_{1-\delta/2} \frac{\tilde{\sigma}_{\alpha,j}}{\sqrt{NT}}, \tilde{\alpha}_j + z_{1-\delta/2} \frac{\tilde{\sigma}_{\alpha,j}}{\sqrt{NT}} \right] \right) \leq \frac{M}{\sqrt{NT}}$  for all  $\gamma \in \mathcal{F}(\nu, E)$  with probability at least  $1 - \epsilon$ . Therefore, our confidence bands are not only honest, they are also very *informative* as they contract as fast as possible and this contraction is uniform over  $\mathcal{F}(\nu, E)$ . Since the de-sparsified Lasso is not a sparse estimator, this fast contraction does not contradict inequality 6 in Theorem 2 of Pötscher (2009) who shows that honest confidence bands based on sparse estimators must be large.

Similarly to the confidence bands pertaining to  $\alpha$ , the ones for the fixed effects are also honest and contract at the optimal rate. Note that this rate is again slower than the one for  $\alpha$ . It is also worth remarking that the above inference results are valid without any sort of lower bound on the non-zero coefficients.

## 5 Monte Carlo

In this section we investigate the finite sample properties of our estimator by means of simulations. In the panel Lasso regression, because the regularization parameter for  $\alpha$  is  $2\lambda$  while that for  $\eta$  is  $2\lambda/\sqrt{N}$ , the option `penalty.factor` in the command `glmnet` is used to adjust this. The options `standardize` and `intercept` in the command `glmnet` are set to `TRUE` and `FALSE`, respectively.<sup>6</sup> The results only changed marginally when we tried other combinations. In the nodewise regression, the option `penalty.factor` in the command `glmnet` is not needed. The options `standardize` and `intercept` in the command `glmnet` are, again, set to `TRUE` and `FALSE`, respectively. The results only changed marginally when we tried `standardize=FALSE`. We suppressed the intercept because it is not needed in the nodewise regression.

Both  $\lambda$  and  $\lambda_{node}$  are chosen via BIC by the formula given in (9.4.9) in Davidson (2000). For example, in the panel Lasso regression, the  $\lambda$  chosen by BIC minimises

$$\log \|y - \Pi\hat{\gamma}(\lambda^*)\|^2 + \frac{\log(NT)}{NT} \|\hat{\gamma}(\lambda^*)\|_0,$$

among a grid of 100 values of  $\lambda^*$  chosen by `glmnet`. Ten-fold cross validation (`cv.glmnet`) was also considered, but this did not alter the results much while being considerably slower. One could also consider adapting the data dependent choice proposed by Belloni et al. (2012) in the context of IV models to the setting of dynamic panel data. We leave it for future work to establish theoretical performance guarantees on these procedures in the setting of high-dimensional dynamic panel data models.

The data generating process is (1.1) and in all experiments  $(\alpha_1, \alpha_2, \alpha_3, \alpha_4) = (0.9, 0, 0, -0.3)$  such that the roots of the corresponding lag polynomial lie outside the unit disk implying stationarity of  $\{y_t\}$ . In practice, one might not know the true lag length and it is usual to specify a reasonably large number of lags (to test downwards). To reflect this in our simulations, we always included 5 lags but also experimented with more than 5 lags. The results were not sensitive to this.

For each  $i = 1, \dots, N$ , the  $x_{i,t}$  are generated according to the autoregressive structure

$$x_{i,t} = a_x x_{i,t-1} + e_{distur,i,t},$$

where the  $e_{distur,i,t}$  are  $p_x \times 1$  random disturbance vectors independent across  $i$  and  $t$ .  $a_x$  is an autoregressive scalar which controls the temporal dependence of  $x_{i,t}$ . For simplicity, we restrict  $a_x$  to be the same across  $i$ . When  $a_x = 0$ , we have temporal independence across  $t$  for  $x_{i,t}$ . Since Assumption 1 does not restrict any temporal dependence of  $x_{i,t}$ , we set  $a_x = 0.5$ . Our

---

<sup>6</sup>We use the publically available `glmnet` package, Friedman et al. (2010), for R, Team (2000).

simulation results are reasonably robust to the choice of  $a_x$ . The covariance matrix of  $e_{distur,i,t}$  is chosen to have a Toeplitz structure with the  $(i, j)$ th entry equal to  $\rho^{|i-j|}$  with  $\rho = 0.75$ . We also experimented with other choices of  $\rho$  which did not change the results dramatically. Furthermore, we also tried to let the covariance matrix of  $e_{distur,i,t}$  be block-diagonal. Again, this did not alter our results.

We allow the fixed effect  $\eta_i$  to depend on the initial observation of  $x_i$ :

$$\eta_i = x'_{i,1} b_\eta / \sqrt{\log p_x} \quad i = 1, \dots, N,$$

where  $b_\eta$  is a  $p_x \times 1$  vector whose entries are drawn from standard normal and normalized to have unit  $\ell_1$ -norm. Note that  $|\eta_i| \leq \|x_{i,1}/\sqrt{\log p_x}\|_\infty \|b_\eta\|_1 = \|x_{i,1}/\sqrt{\log p_x}\|_\infty$ . If  $x_{i,1}$  is multivariate normal, then  $\|x_{i,1}\|_\infty = O_p(\sqrt{\log p_x})$ . In this sense,  $\eta_i$  is bounded. However,  $\eta$  is not necessarily weakly sparse and thus we actually also investigate how robust our results are to violations of this assumption. Of course our estimator performed much better in the truly weakly sparse setting than the setting we present here (results available upon request).

As our theory allows for heteroskedasticity, we also investigate the effect of this. To be precise, we consider error terms of the form  $\varepsilon_{i,t} = u_{i,t} (x_{i,t,1}/\sqrt{2} + b_x x_{i,t,2})$  where  $u_{i,t}$  is independent of  $y_{i,t-1}, \dots, y_{i,1-L}$  and  $x_{i,t}, \dots, x_{i,1}$ .  $b_x$  is chosen such that the unconditional variance of  $\varepsilon_{i,t}$  is the same as the one of  $u_{i,t}$  which in turn equals the one from the homoskedastic case. A simple calculation reveals that  $b_x = (-\sqrt{2}\rho + \sqrt{2\rho^2 + 2 - 4a_x^2})/2$ . Note that  $\varepsilon_{i,t}$  constructed this way satisfies Assumption 1. The reason we ensure that the unconditional variance is the same as in the homoskedastic case is that we do not want any findings in the heteroskedastic case to be driven by a plain change in the unconditional variance.

Our estimator is compared to the least squares oracle which only includes variables with non-zero coefficients in addition to those variables we wish to test hypotheses about. Thus, it is an oracle which knows the relevant control variables. When sample size allows it, that is when  $p + N \leq NT$ , we also implement naive least squares including all variables. This estimator is numerically equivalent to the often used within estimator. Finally, we implemented the de-sparsified conservative Lasso of Caner and Kock (2014). However, this only improved the results slightly and so we do not report these results here. The number of Monte Carlo replications is 1,000 for all setups and we consider the performance of our estimator along the following dimensions:

1. Estimation error: We compute the root mean square errors (RMSE) of all procedures averaged over the Monte Carlo replications.
2. Coverage rate: We calculate the coverage rate of a gaussian confidence interval constructed

as in Theorem 3. This is done for three coefficients of regressors in  $x_{i,t}$ .

3. Length of confidence interval: We calculate the length of the three confidence intervals considered in point 2 above.
4. Size: We evaluate the size of the  $\chi^2$ -test in Theorem 2 for a hypothesis involving the same three parameters for which we construct confidence intervals in point 2 above.
5. Power: We evaluate the power of the  $\chi^2$ -test in point 4 above.

All tests are carried out at the 5% level of significance and all confidence intervals have a nominal coverage of 95%. Furthermore, as the oracle inequalities in Theorem 1 are for the plain Lasso, the root mean square errors are reported for this instead of the de-sparsified Lasso. Moreover, to compute the estimation error of the de-sparsified Lasso, one needs to run nodewise regressions for all  $p$  columns of  $Z$ . Note that, to conduct joint inference on  $h$  coefficients, one only needs to run nodewise regressions for the  $h$  variables whose coefficients are involved in the hypothesis being tested. As our models are dynamic, we allow for a burn-in period of 1,000 observations when generating the data.

The following experiments were carried out

- Experiment 1: (moderate-dimensional setting):  $N = 20$  and  $T = 10$ .  $\beta$  is  $100 \times 1$  with five equidistant non-zero entries equaling one. Thus,  $p = 105$  and the number of non-zero entries in  $\alpha$  is seven. In total,  $\gamma = (\alpha', \eta')'$  is  $125 \times 1$ . The disturbances of  $x_{i,t}$ ,  $e_{distur,i,t}$ , are gaussian and  $\varepsilon_{i,t}$  are standard gaussian. We test the true hypothesis

$$H_0 : (\gamma_7, \gamma_{27}, \gamma_{47}) = (0, 0, 0)$$

by the  $\chi^2_3$  test described in Theorem 2 in order to gauge the size of the test. The power is investigated by the hypothesis

$$H_0 : (\gamma_7, \gamma_{27}, \gamma_{47}) = (0.4, 0, 0).$$

The following variations of this setting are considered

- (a) The baseline case described so far.
- (b) Same as (a) but with heteroskedastic errors.
- (c) Same as (b) but  $e_{distur,i,t}$  and  $\varepsilon_{i,t}$  are  $t$ -distributed with 3 degrees of freedom. In this case,  $\eta_i$  may not even be  $O_p(1)$ .

- Experiment 2: (high-dimensional setting).  $N = 20$  and  $T = 10$ .  $\beta$  is  $400 \times 1$  with five equidistant non-zero entries equaling one. Thus,  $p = 405$  and the number of non-zero entries in  $\alpha$  is seven. In total,  $\gamma = (\alpha', \eta')'$  is  $425 \times 1$ . The disturbances of  $x_{i,t}$ ,  $e_{distur,i,t}$ , are gaussian and  $\varepsilon_{i,t}$  are standard gaussian. We test the true hypothesis

$$H_0 : (\gamma_7, \gamma_{87}, \gamma_{167}) = (0, 0, 0)$$

by the  $\chi_3^2$  test described in Theorem 2 in order to gauge the size of the test. The power is investigated by the hypothesis

$$H_0 : (\gamma_7, \gamma_{87}, \gamma_{167}) = (0.4, 0, 0).$$

The following variations of this setting are considered

- (a) The baseline case described so far.
  - (b) Same as (a) but with heteroskedastic errors.
  - (c) Same as (b) but  $e_{distur,i,t}$  and  $\varepsilon_{i,t}$  are  $t$ -distributed with 3 degrees of freedom. In this case,  $\eta_i$  may not even be  $O_p(1)$ .
- Experiment 3: (increase  $T$ ): As Experiment 2 but with  $T = 40$ .
  - Experiment 4: (increase  $N$ ): As Experiment 2 but with  $N = 40$ .
  - Experiment 5: (high-dimensional setting 2).  $N = 20$  and  $T = 40$ .  $\beta$  is  $1005 \times 1$  with 15 equidistant non-zero entries equaling one. Thus,  $p = 1010$  and the number of non-zero entries in  $\alpha$  is seventeen. In total,  $\gamma = (\alpha', \eta')'$  is  $1030 \times 1$ . The disturbances of  $x_{i,t}$ ,  $e_{distur,i,t}$ , are gaussian and  $\varepsilon_{i,t}$  are standard gaussian. We test the true hypothesis

$$H_0 : (\gamma_7, \gamma_{74}, \gamma_{141}) = (0, 0, 0)$$

by the  $\chi_3^2$  test described in Theorem 2 in order to gauge the size of the test. The power is investigated by the hypothesis

$$H_0 : (\gamma_7, \gamma_{74}, \gamma_{141}) = (0.4, 0, 0).$$

The following variations of this setting are considered

- (a) The baseline case described so far.
- (b) Same as (a) but with heteroskedastic errors.
- (c) Same as (b) but  $e_{distur,i,t}$  and  $\varepsilon_{i,t}$  are  $t$ -distributed with 3 degrees of freedom. In this case,  $\eta_i$  may not even be  $O_p(1)$ .

		RMSE		Coverage			Length			Size	Power
		$\alpha$	$\eta$	$\gamma_7$	$\gamma_{27}$	$\gamma_{47}$	$\gamma_7$	$\gamma_{27}$	$\gamma_{47}$		
1(a)	LS	23.744	16.382	0.762	0.786	0.766	0.552	0.549	0.553	0.412	0.741
	DL	3.061	8.528	0.892	0.918	0.884	0.395	0.396	0.396	0.150	0.852
	Ora	0.523	7.226	0.933	0.939	0.919	0.402	0.403	0.404	0.074	0.907
1(b)	LS	23.796	16.444	0.732	0.760	0.747	0.548	0.543	0.547	0.453	0.747
	DL	3.011	8.298	0.920	0.914	0.903	0.408	0.389	0.391	0.135	0.846
	Ora	0.524	7.079	0.920	0.931	0.937	0.401	0.393	0.398	0.092	0.904
1(c)	LS	46.140	51.979	0.753	0.772	0.743	0.910	0.883	0.889	0.432	0.605
	DL	4.747	23.939	0.912	0.901	0.883	0.619	0.544	0.567	0.159	0.632
	Ora	1.200	23.596	0.907	0.937	0.913	0.662	0.617	0.617	0.091	0.651

Table 1: Experiment 1. LS, DL and Ora: least squares including all variables, de-sparsified Lasso and least squares oracle. RMSE: root mean square error. Coverage: the coverage rate of the asymptotic 95% confidence intervals. Length: the average length of the asymptotic 95% confidence intervals. Size: size of the correct hypothesis  $H_0 : (\gamma_7, \gamma_{27}, \gamma_{47}) = (0, 0, 0)$ . Power: the probability to reject the false  $H_0 : (\gamma_7, \gamma_{27}, \gamma_{47}) = (0.4, 0, 0)$ .

Table 1 contains the results of experiment 1. Setting 1(a) reveals that the RMSE of the Lasso are lower than those for least squares including all variables but higher than those of least squares only including the relevant variables. This is the case for  $\alpha$  as well as the fixed effects. Next, it is very encouraging that the coverage probabilities for the de-sparsified Lasso are close to the ones based on the oracle. The lengths of the confidence intervals are also comparable for those two procedures while the ones based on the within estimator are considerably wider while still having a lower coverage. The oracle and the de-sparsified Lasso both produce tests which are a bit oversized but they are still much better than the within estimator. The same is true when it comes to power.

Experiment 1(b) adds heteroskedasticity to the error terms and none of the procedures is affected by this.

In Panel 1(c) the random variables have heavy tails. Overall, and as expected, all procedures suffer from this. However, it is worth mentioning that the coverage rate of the confidence intervals does not decrease. Instead, the length of these intervals increases to reflect the larger uncertainty. The size of the significance test is not affected either while the power suffers.

Next, we turn to experiment 2(a) which is high-dimensional. The results can be found in



		RMSE		Coverage			Length			Size	Power
		$\alpha$	$\eta$	$\gamma_7$	$\gamma_{87}$	$\gamma_{167}$	$\gamma_7$	$\gamma_{87}$	$\gamma_{167}$		
LS											
2(a)	DL	4.209	8.333	0.875	0.893	0.881	0.386	0.385	0.386	0.189	0.841
	Ora	0.513	7.103	0.919	0.918	0.924	0.402	0.403	0.403	0.110	0.922
LS											
2(b)	DL	4.165	8.322	0.896	0.872	0.861	0.407	0.379	0.381	0.189	0.825
	Ora	0.535	7.074	0.906	0.913	0.929	0.401	0.396	0.397	0.101	0.899
LS											
2(c)	DL	7.602	22.895	0.916	0.868	0.882	0.622	0.543	0.551	0.193	0.602
	Ora	1.074	21.724	0.922	0.944	0.944	0.657	0.619	0.632	0.076	0.674

Table 2: Experiment 2. LS, DL and Ora: least squares including all variables, de-sparsified Lasso and least squares oracle. RMSE: root mean square error. Coverage: the coverage rate of the asymptotic 95% confidence intervals. Length: the average length of the asymptotic 95% confidence intervals. Size: size of the correct hypothesis  $H_0 : (\gamma_7, \gamma_{87}, \gamma_{167}) = (0, 0, 0)$ . Power: the probability to reject the false  $H_0 : (\gamma_7, \gamma_{87}, \gamma_{167}) = (0.4, 0, 0)$ .

Table 2. As expected, the estimation error is higher for the Lasso than for the oracle. However, it is encouraging that the confidence intervals produced by the de-sparsified Lasso have coverage which is as almost as good as the one for the oracle. In fact, the coverage rate is close to identical to the one in the above moderate-dimensional simulation. The confidence bands based on the de-sparsified Lasso are actually slightly shorter than the ones based on the oracle which may explain their slightly lower coverage. The significance test is again a bit oversized for the oracle as well as the de-sparsified Lasso but the size is not far from the one in Table 1. Power is also virtually unaffected by the increase in dimension.

Experiment 2(b) adds heteroskedasticity and the results are not affected by this. Finally, the addition of heavy tails in Experiment 2(c) makes the estimators less precise. However, and as in the moderate-dimensional setting above, the coverage remains high since the confidence bands get wider. The size of the significance test is unaffected while the power goes down for the oracle as well as the de-sparsified Lasso.

In Table 3,  $T$  has been increased to 40 compared to Table 2. This results in lower estimation errors for the Lasso as well as oracle assisted least squares. The coverage rates of the confidence bands also improve and get closer to the nominal rate. At the same time, the bands also become

		RMSE		Coverage			Length			Size	Power
		$\alpha$	$\eta$	$\gamma_7$	$\gamma_{87}$	$\gamma_{167}$	$\gamma_7$	$\gamma_{87}$	$\gamma_{167}$		
3(a)	LS	40.341	7.367	0.815	0.827	0.796	0.266	0.266	0.268	0.325	0.993
	DL	1.208	2.121	0.931	0.918	0.925	0.190	0.190	0.190	0.098	1.000
	Ora	0.223	3.208	0.943	0.945	0.933	0.186	0.187	0.187	0.052	1.000
3(b)	LS	40.351	7.376	0.800	0.823	0.813	0.267	0.265	0.267	0.318	0.993
	DL	1.169	2.139	0.923	0.915	0.924	0.200	0.189	0.189	0.104	1.000
	Ora	0.233	3.235	0.955	0.940	0.953	0.189	0.185	0.186	0.060	1.000
3(c)	LS	88.649	29.738	0.766	0.813	0.837	0.460	0.452	0.448	0.315	0.821
	DL	2.630	7.689	0.931	0.922	0.913	0.359	0.309	0.319	0.090	0.926
	Ora	0.610	10.759	0.930	0.963	0.951	0.329	0.299	0.302	0.056	0.962

Table 3: Experiment 3. LS, DL and Ora: least squares including all variables, de-sparsified Lasso and least squares oracle. RMSE: root mean square error. Coverage: the coverage rate of the asymptotic 95% confidence intervals. Length: the average length of the asymptotic 95% confidence intervals. Size: size of the correct hypothesis  $H_0 : (\gamma_7, \gamma_{87}, \gamma_{167}) = (0, 0, 0)$ . Power: the probability to reject the false  $H_0 : (\gamma_7, \gamma_{87}, \gamma_{167}) = (0.4, 0, 0)$ .

more narrow. The size of the significance test also improves and the power of the oracle and the de-sparsified Lasso is 1. As above, adding heteroskedasticity does not alter the results. The consequences of heavy tails are also the same: higher estimation error, no change in coverage of the confidence bands, wider bands, unchanged size, but lower power.

Table 4 increases  $N$  to 40 compared to Table 2. This results in more fixed effects to be estimated. Thus, it is not surprising that the estimation error for  $\alpha$  goes down while the one for  $\eta$  increases. The coverage rates of the oracle as well as the de-sparsified Lasso improve compared to Table 2. However, the length of the confidence bands does not decrease as much as when  $T$  was increased in the previous experiment. The size of the significance test decreases when increasing  $N$  while power is close to 1. Adding heteroskedasticity has no consequences while the presence of heavy tails has the usual effect.

Table 5 contains a setting with more than 1000 variables. The main message of the previous tables prevails even in this setting: the coverages of the Lasso-based confidence intervals are almost as high as the ones based on the oracle. On the other hand, the bands of the former are now slightly wider than the ones of the latter. Both procedures have power close to one while the Lasso-based test is a bit oversized compared to the oracle based test. Heteroskedasticity

		RMSE		Coverage			Length			Size	Power
		$\alpha$	$\eta$	$\gamma_7$	$\gamma_{87}$	$\gamma_{167}$	$\gamma_7$	$\gamma_{87}$	$\gamma_{167}$		
LS											
4(a)	DL	2.145	14.002	0.924	0.891	0.920	0.261	0.262	0.263	0.123	0.988
	Ora	0.351	13.570	0.936	0.934	0.930	0.282	0.283	0.285	0.065	0.999
LS											
4(b)	DL	2.145	14.073	0.933	0.904	0.911	0.275	0.260	0.263	0.114	0.979
	Ora	0.368	13.469	0.926	0.931	0.940	0.283	0.281	0.282	0.076	1.000
LS											
4(c)	DL	4.305	42.303	0.917	0.899	0.903	0.456	0.386	0.390	0.139	0.825
	Ora	0.782	41.269	0.926	0.934	0.926	0.465	0.428	0.435	0.087	0.870

Table 4: Experiment 4. LS, DL and Ora: least squares including all variables, de-sparsified Lasso and least squares oracle. RMSE: root mean square error. Coverage: the coverage rate of the asymptotic 95% confidence intervals. Length: the average length of the asymptotic 95% confidence intervals. Size: size of the correct hypothesis  $H_0 : (\gamma_7, \gamma_{87}, \gamma_{167}) = (0, 0, 0)$ . Power: the probability to reject the false  $H_0 : (\gamma_7, \gamma_{87}, \gamma_{167}) = (0.4, 0, 0)$ .

does not affect the results. The consequences of heavy tails are the same as in the previous experiments.

		RMSE		Coverage			Length			Size	Power
		$\alpha$	$\eta$	$\gamma_7$	$\gamma_{74}$	$\gamma_{141}$	$\gamma_7$	$\gamma_{74}$	$\gamma_{141}$		
LS											
5(a)	DL	3.342	2.463	0.922	0.903	0.912	0.209	0.209	0.208	0.124	0.997
	Ora	0.546	3.305	0.956	0.942	0.949	0.187	0.187	0.187	0.062	1.000
LS											
5(b)	DL	3.327	2.432	0.913	0.916	0.896	0.218	0.207	0.208	0.127	0.998
	Ora	0.556	3.261	0.942	0.951	0.921	0.190	0.186	0.186	0.065	1.000
LS											
5(c)	DL	7.294	7.273	0.936	0.910	0.916	0.363	0.307	0.304	0.101	0.920
	Ora	1.072	9.693	0.952	0.936	0.951	0.326	0.297	0.293	0.061	0.955

Table 5: Experiment 5. LS, DL and Ora: least squares including all variables, de-sparsified Lasso and least squares oracle. RMSE: root mean square error. Coverage: the coverage rate of the asymptotic 95% confidence intervals. Length: the average length of the asymptotic 95% confidence intervals. Size: size of the correct hypothesis  $H_0 : (\gamma_7, \gamma_{74}, \gamma_{141}) = (0, 0, 0)$ . Power: the probability to reject the false  $H_0 : (\gamma_7, \gamma_{74}, \gamma_{141}) = (0.4, 0, 0)$ .

## 6 Conclusion

This paper has considered inference in high-dimensional dynamic panel data models with fixed effects. In particular we have shown how hypotheses involving an increasing number of parameters can be tested. These hypotheses can involve parameters from all groups of variables in the model. As a stepping stone towards this inference we constructed a uniformly valid estimator of the covariance matrix of the parameter estimates which is robust towards conditional heteroskedasticity. We also stress that our theory does not require the inverse covariance matrix of the covariates to be exactly sparse.

Next, we showed that confidence bands based on our procedure are asymptotically honest and contract at the optimal rate. This rate of contraction depends on which type of parameter is under consideration. Simulations revealed that our procedure works well in finite samples. Future work may include relaxing the weak sparsity assumption on the inverse covariance matrix  $\Theta_Z$  as well as extending our results to non-linear panel data models.

## References

- Andersen, T. B., J. Bentzen, C.-J. Dalgaard, and P. Selaya (2012). Lightning, IT diffusion, and economic growth across U.S. states. *The Review of Economics and Statistics* 94(4), 903–924.
- Arellano, M. (2003). Panel data econometrics. *OUP Catalogue*.
- Baltagi, B. (2008). *Econometric analysis of panel data*, Volume 1. John Wiley & Sons.
- Belloni, A., D. Chen, V. Chernozhukov, and C. Hansen (2012). Sparse models and methods for optimal instruments with an application to eminent domain. *Econometrica* 80(6), 2369–2429.
- Belloni, A., V. Chernozhukov, and C. Hansen (2014). Inference on treatment effects after selection among high-dimensional controls. *The Review of Economic Studies* 81(2), 608–650.
- Belloni, A., V. Chernozhukov, C. Hansen, and D. Kozbur (2015). Inference in high-dimensional panel models with an application to gun control. *Journal of Business & Economic Statistics* 34(4), 590–605.
- Bickel, P. J., Y. Ritov, and A. B. Tsybakov (2009). Simultaneous analysis of Lasso and Dantzig selector. *The Annals of Statistics* 37(4), 1705–1732.
- Bonhomme, S. and E. Manresa (2015). Grouped patterns of heterogeneity in panel data. *Econometrica* 83(3), 1147–1184.
- Bühlmann, P. and S. van de Geer (2011). *Statistics for High-Dimensional Data*. Springer.
- Caner, M. and A. B. Kock (2014). Asymptotically honest confidence regions for high dimensional parameters by the de-sparsified conservative Lasso. *arXiv preprint arXiv:1410.4208*.
- Caner, M. and H. H. Zhang (2014). Adaptive elastic net for generalized methods of moments. *Journal of Business & Economic Statistics* 32(1), 30–47.
- Davidson, J. (2000). *Econometric Theory*. Blackwell Publishers.
- De Neve, J.-E., N. A. Christakis, J. H. Fowler, and B. S. Frey (2012). Genes, economics, and happiness. *Journal of Neuroscience, Psychology, and Economics* 5(4), 193.
- Engle, R. F. (1982). Autoregressive conditional heteroscedasticity with estimates of the variance of united kingdom inflation. *Econometrica: Journal of the Econometric Society*, 987–1007.
- Fan, X., I. Grama, and Q. Liu (2012). Large deviation exponential inequalities for supermartingales. *Electronic Communications in Probability* 17(59), 1–8.

- Friedman, J., T. Hastie, and R. Tibshirani (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software* 33(1), 1–22.
- Galvao, A. F. and G. V. Montes-Rojas (2010). Penalized quantile regression for dynamic panel data. *Journal of Statistical Planning and Inference* 140(11), 3476–3497.
- Horn, R. A. and C. R. Johnson (1990). *Matrix analysis*. Cambridge University Press.
- Islam, N. (1995). Growth empirics: a panel data approach. *The Quarterly Journal of Economics*, 1127–1170.
- Javanmard, A. and A. Montanari (2013). Confidence intervals and hypothesis testing for high-dimensional regression. *arXiv preprint arXiv:1306.3171*.
- Knight, K. and W. Fu (2000). Asymptotics for Lasso-type estimators. *The Annals of Statistics* 28(5), 1356–1378.
- Kock, A. B. (2013). Oracle efficient variable selection in random and fixed effects panel data models. *Econometric Theory* 29, 115–152.
- Kock, A. B. (2016). Oracle inequalities, variable selection and uniform inference in high-dimensional correlated random effects panel data models. *Journal of Econometrics* 195(1), 71–85.
- Koenker, R. (2004). Quantile regression for longitudinal data. *Journal of Multivariate Analysis* 91(1), 74–89.
- Leeb, H. and B. M. Pötscher (2005). Model selection and inference: Facts and fiction. *Econometric Theory* 21, 21–59.
- Lesigne, E. and D. Volny (2001). Large deviations for martingales. *Stochastic Processes and Their Applications* 96(1), 143–159.
- Li, K.-C. (1989). Honest confidence regions for nonparametric regression. *The Annals of Statistics*, 1001–1008.
- Lu, X. and L. Su (2016). Shrinkage estimation of dynamic panel data models with interactive fixed effects. *Journal of Econometrics* 190(1), 148–175.
- McLeish, D. L. (1974). Dependent central limit theorems and invariance principles. *The Annals of Probability* 2(4), 620–628.

- Merlevède, F., M. Peligrad, and E. Rio (2011). A Bernstein type inequality and moderate deviations for weakly dependent sequences. *Probability Theory and Related Fields* 151(3-4), 435–474.
- Negahban, S. N., P. Ravikumar, M. J. Wainwright, and B. Yu (2012). A unified framework for high-dimensional analysis of m-estimators with decomposable regularizers. *Statistical Science* 27(4), 538–557.
- Nickell, S. (1981). Biases in dynamic models with fixed effects. *Econometrica* 49(6), 1417–1426.
- Pötscher, B. M. (2009). Confidence sets based on sparse estimators are necessarily large. *Sankhyā: The Indian Journal of Statistics, Series A (2008-)*, 1–18.
- Su, L., Z. Shi, and P. C. Phillips (2016). Identifying latent structures in panel data. *Econometrica* 84(6), 2215–2264.
- Team, R. C. (2000). R language definition. *Vienna, Austria: R foundation for statistical computing*.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B* 58, 267–288.
- van de Geer, S., P. Bühlmann, Y. Ritov, and R. Dezeure (2014). On asymptotically optimal confidence regions and tests for high-dimensional models. *The Annals of Statistics* 42(3), 1166–1202.
- van de Geer, S. A., P. Bühlmann, et al. (2009). On the conditions used to prove oracle results for the lasso. *Electronic Journal of Statistics* 3, 1360–1392.
- van der Vaart, A. and J. A. Wellner (1996). *Weak Convergence and Empirical Processes*. Springer.
- Wooldridge, J. M. (2010). *Econometric analysis of cross section and panel data*. MIT press.
- Yuan, M. (2010). High dimensional inverse covariance matrix estimation via linear programming. *Journal of Machine Learning Research*.
- Zhang, C.-H. and S. S. Zhang (2014). Confidence intervals for low dimensional parameters in high dimensional linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 76(1), 217–242.

## 7 Appendix A

### 7.1 Examples of Weakly Sparse Fixed Effects

The following lemma provides sufficient conditions for the required weak sparsity of the fixed effects  $\eta_i$ .

**Lemma 1.** *Let  $0 < \nu < 1$  be given and assume that  $\max_{1 \leq i \leq N} \mathbb{E}|\eta_i|^\nu = O(\sqrt{(\log(p \vee N))^{3\nu}/(NT^\nu)})$  and  $\max_{1 \leq i \leq N} \mathbb{E}|\eta_i|^{r\nu} = O(\sqrt{(\log(p \vee N))^{3r\nu}/T^{r\nu}})$  for some  $r \geq 2$ . Then,  $\sum_{i=1}^N |\eta_i|^\nu = O_p(\sqrt{N(\log(p \vee N))^{3\nu}/T^\nu})$ <sup>7</sup>. The sufficient conditions are met if, for example,*

1.  $\eta_i \sim N(0, \sigma_i^2)$  with  $\max_{1 \leq i \leq N} \sigma_i^2 = O((\log(p \vee N))^3/(N^{1/\nu}T))$
2.  $\eta_i \sim \text{Uniform}[-a, a]$  with  $a = O(\sqrt{(\log(p \vee N))^3/(N^{1/\nu}T)})$

*Proof.* It suffices to show that  $\sum_{i=1}^N \mathbb{E}|\eta_i|^\nu = O(\sqrt{N(\log(p \vee N))^{3\nu}/T^\nu})$  and  $|\sum_{i=1}^N (|\eta_i|^\nu - \mathbb{E}|\eta_i|^\nu)| = O_p(\sqrt{N(\log(p \vee N))^{3\nu}/T^\nu})$ . First,

$$\sum_{i=1}^N \mathbb{E}|\eta_i|^\nu \leq N \max_{1 \leq i \leq N} \mathbb{E}|\eta_i|^\nu = O(\sqrt{N(\log(p \vee N))^{3\nu}/T^\nu})$$

since  $\max_{1 \leq i \leq N} \mathbb{E}|\eta_i|^\nu = O(\sqrt{(\log(p \vee N))^{3\nu}/(NT^\nu)})$ .

Second, for any  $t > 0$  it follows by Rosenthal's inequality that for any  $r \geq 2$  (the constant  $C$  may change from line to line)

$$\begin{aligned} \mathbb{P}\left(\left|\sum_{i=1}^N (|\eta_i|^\nu - \mathbb{E}|\eta_i|^\nu)\right| > t\right) &\leq C \frac{\left[\sum_{i=1}^N \mathbb{E}(|\eta_i|^\nu - \mathbb{E}|\eta_i|^\nu)^2\right]^{r/2} + \sum_{i=1}^N \mathbb{E}|\eta_i|^\nu - \mathbb{E}|\eta_i|^\nu}{t^r} \\ &\leq C \frac{\left[\sum_{i=1}^N \mathbb{E}|\eta_i|^{2\nu}\right]^{r/2} + \sum_{i=1}^N \mathbb{E}|\eta_i|^{r\nu}}{t^r} \\ &\leq C \frac{N^{r/2} \max_{1 \leq i \leq N} \mathbb{E}|\eta_i|^{r\nu} + N \max_{1 \leq i \leq N} \mathbb{E}|\eta_i|^{r\nu}}{t^r} \\ &\leq C \frac{N^{r/2} \max_{1 \leq i \leq N} \mathbb{E}|\eta_i|^{r\nu}}{t^r} \end{aligned}$$

where the second and third inequality both use the convexity of  $x \mapsto x^a$  for  $a \geq 1$  to invoke Jensen's inequality repeatedly. Thus, setting  $t = M \cdot \sqrt{N(\log(p \vee N))^{3\nu}/T^\nu}$  for some constant  $M > 0$ , it is seen that it suffices that  $\max_{1 \leq i \leq N} \mathbb{E}|\eta_i|^{r\nu} = O((\log(p \vee N))^{\frac{3}{2}r\nu}/T^{\frac{1}{2}r\nu})$ .

The primitive conditions for the normal and uniform distributions follow from the fact that  $\mathbb{E}|\eta_i|^s \sim \sigma_i^s$ ,  $a^s$ , respectively. □

---

<sup>7</sup>If the big  $O$ s are replaced by small  $o$ s in the hypothesis of the lemma the big  $O$  in the conclusion of the lemma can also be replaced by a small  $o$ .



Note that the sufficient conditions for weak sparsity of the fixed effects are given entirely in terms of their moments. The lemma does *not* restrict the correlation between the fixed effects and the common covariates. By choosing  $r = 2$  one sees that the  $\eta_i$  need not even have two moments. The specific example of the normal distribution shows that the maximal variance of the fixed effects must tend to zero. While this is not an innocent assumption and restricts the size of the  $\eta_i$  recall that the dimension  $p$  of  $z_{i,t}$  in

$$y_{i,t} = z'_{i,t}\alpha + \eta_i + \varepsilon_{i,t}, \quad i = 1, \dots, N, \quad t = 1, \dots, T$$

can grow almost exponentially in the sample size leaving very little of the variation in  $y_{i,t}$  to be explained by  $\eta_i$ . Put differently, the observed covariates in  $z_{i,t}$  will be so numerous that very little variation is left for the fixed effects to explain.

## 7.2 Sufficient Conditions for $y_{i,t}$ to be Subgaussian

The following lemma provides sufficient conditions for  $y_{i,t}$  to inherit the subgaussianity from the covariates and the error terms. It allows for a wide range of models but rules out dynamic panel data models which are explosive or contain unit roots.

**Lemma 2.** *Let  $x_{i,t}$  and  $\varepsilon_{i,t}$  be uniformly subgaussian for  $i = 1, \dots, N$ ,  $t = 1, \dots, T$  (as defined in Assumption 3(b) and 3(a), respectively) and assume that  $\|\beta\| \leq C$  for some  $C > 0$ . Furthermore,  $\max_{1 \leq i \leq N} |\eta_i|$  is bounded uniformly in  $N$ . Then, if all roots of  $1 - \sum_{j=1}^L \alpha_j z^j$  ( $\alpha_1, \dots, \alpha_L$  fixed) are outside the unit disc,  $y_{i,t}$  is uniformly subgaussian for  $i = 1, \dots, N$  and  $t = 1, \dots, T$ .*

*Proof of Lemma 2.* Let  $y_t = \sum_{j=1}^L \alpha_j y_{t-j} + u_t$  be an AR( $L$ ) process with roots outside the unit disc. Write the companion form as  $\xi_t = F\xi_{t-1} + v_t$ . Then, by the monotone convergence theorem for Orlicz norms, see van der Vaart and Wellner (1996) exercise 6, page 105,  $\|\|\xi_t\|\|_{\psi_2} \leq \|\sum_{j=1}^{\infty} \|F^j\|_{\ell_2} \|v_{t-j}\|\|_{\psi_2} = \sum_{j=1}^{\infty} \|F^j\|_{\ell_2} \|v_{t-j}\|_{\psi_2} = \sum_{j=1}^{\infty} \|F^j\|_{\ell_2} \|u_{t-j}\|_{\psi_2}$ , where  $\|\cdot\|_{\ell_2}$  is the  $\ell_2$  induced norm, and the last equality used that  $v_t$  is  $L \times 1$  with only one non-zero entry equaling  $u_t$ . By Corollary 5.6.14 in Horn and Johnson (1990) there exists a  $1 > \delta > 0$  such that  $\|F^j\|_{\ell_2} \leq (1 - \delta)^j$  for  $j$  sufficiently large. Thus, if  $\|u_t\|_{\psi_2}$  is uniformly bounded we conclude  $\|y_t\|_{\psi_2} \leq \|\|\xi_t\|\|_{\psi_2} \leq K$  for some  $K > 0$ . Thus, in our context it suffices to show that  $\|x'_{i,t}\beta + \eta_i + \varepsilon_{i,t}\|_{\psi_2}$  is uniformly bounded as  $y_{i,t} = \sum_{j=1}^L \alpha_j y_{i,t-j} + x'_{i,t}\beta + \eta_i + \varepsilon_{i,t} = \sum_{j=1}^L \alpha_j y_{i,t-j} + u_{i,t}$  with  $u_{i,t} = x'_{i,t}\beta + \eta_i + \varepsilon_{i,t}$ . But  $\|x'_{i,t}\beta + \eta_i + \varepsilon_{i,t}\|_{\psi_2} \leq \|\beta' x_{i,t}\|_{\psi_2} + \|\eta_i\|_{\psi_2} + \|\varepsilon_{i,t}\|_{\psi_2} \leq \|\beta\| \sup_{\|v\| \leq 1} \|v' x_{i,t}\|_{\psi_2} + \|\eta_i\|_{\psi_2} + \|\varepsilon_{i,t}\|_{\psi_2}$  which is bounded by the assumptions made.  $\square$

## 8 Appendix B

### 8.1 Proof of Theorem 1

This appendix proves Theorem 1. To this end, we introduce some auxiliary lemmas. Define the events

$$\mathcal{A} = \left\{ \|Z'\varepsilon\|_\infty \leq \frac{\lambda}{2}, \quad \|D'\varepsilon\|_\infty \leq \frac{\lambda}{2\sqrt{N}} \right\}, \quad \mathcal{B} = \left\{ \kappa^2(\Psi_N, s_1, s_2) \geq \frac{\kappa^2}{2} \right\}.$$

**Lemma 3.** *On the event  $\mathcal{A}$ , the following inequalities are valid*

$$\|\Pi(\hat{\gamma} - \gamma)\|^2 + \lambda \|\hat{\alpha} - \alpha\|_1 + \frac{\lambda}{\sqrt{N}} \|\hat{\eta} - \eta\|_1 \leq 4\lambda \|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1 + \frac{4\lambda}{\sqrt{N}} \|\hat{\eta}_{J_2} - \eta_{J_2}\|_1 + 4\lambda E_1 \Xi_1^{1-\nu} + \frac{4\lambda}{\sqrt{N}} E \Xi_2^{1-\nu}; \quad (8.1)$$

$$\|\hat{\alpha}_{J_1^c} - \alpha_{J_1^c}\|_1 + \frac{1}{\sqrt{N}} \|\hat{\eta}_{J_2^c} - \eta_{J_2^c}\|_1 \leq 3\|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1 + \frac{3}{\sqrt{N}} \|\hat{\eta}_{J_2} - \eta_{J_2}\|_1 + 4E_1 \Xi_1^{1-\nu} + \frac{4}{\sqrt{N}} E \Xi_2^{1-\nu}. \quad (8.2)$$

*Proof.* By the minimizing property of the Lasso,

$$\|y - \Pi\hat{\gamma}\|^2 + 2\lambda \|\hat{\alpha}\|_1 + 2\frac{\lambda}{\sqrt{N}} \|\hat{\eta}\|_1 \leq \|y - \Pi\gamma\|^2 + 2\lambda \|\alpha\|_1 + 2\frac{\lambda}{\sqrt{N}} \|\eta\|_1$$

such that inserting  $y = \Pi\gamma + \varepsilon$  yields

$$\|\Pi(\hat{\gamma} - \gamma)\|^2 \leq 2\varepsilon' \Pi(\hat{\gamma} - \gamma) + 2\lambda(\|\alpha\|_1 - \|\hat{\alpha}\|_1) + 2\frac{\lambda}{\sqrt{N}}(\|\eta\|_1 - \|\hat{\eta}\|_1). \quad (8.3)$$

Note that on  $\mathcal{A}$

$$2\varepsilon' \Pi(\hat{\gamma} - \gamma) \leq 2\|\varepsilon' Z\|_\infty \|\hat{\alpha} - \alpha\|_1 + 2\|\varepsilon' D\|_\infty \|\hat{\eta} - \eta\|_1 \leq \lambda \|\hat{\alpha} - \alpha\|_1 + \frac{\lambda}{\sqrt{N}} \|\hat{\eta} - \eta\|_1.$$

Using this and adding  $\lambda \|\hat{\alpha} - \alpha\|_1 + \frac{\lambda}{\sqrt{N}} \|\hat{\eta} - \eta\|_1$  to both sides of (8.3) gives

$$\begin{aligned} & \|\Pi(\hat{\gamma} - \gamma)\|^2 + \lambda \|\hat{\alpha} - \alpha\|_1 + \frac{\lambda}{\sqrt{N}} \|\hat{\eta} - \eta\|_1 \\ & \leq 2\lambda(\|\alpha\|_1 - \|\hat{\alpha}\|_1 + \|\hat{\alpha} - \alpha\|_1) + 2\frac{\lambda}{\sqrt{N}}(\|\eta\|_1 - \|\hat{\eta}\|_1 + \|\hat{\eta} - \eta\|_1) \\ & \leq 2\lambda(\|\alpha_{J_1}\|_1 - \|\hat{\alpha}_{J_1}\|_1 + \|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1 + 2E_1 \Xi_1^{1-\nu}) + 2\frac{\lambda}{\sqrt{N}}(\|\eta_{J_2}\|_1 - \|\hat{\eta}_{J_2}\|_1 + \|\hat{\eta}_{J_2} - \eta_{J_2}\|_1 + 2E \Xi_2^{1-\nu}) \\ & \leq 4\lambda \|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1 + 4\frac{\lambda}{\sqrt{N}} \|\hat{\eta}_{J_2} - \eta_{J_2}\|_1 + 4\lambda E_1 \Xi_1^{1-\nu} + \frac{4\lambda}{\sqrt{N}} E \Xi_2^{1-\nu}, \end{aligned}$$

where the second inequality is due to, taking  $\eta$  as an illustration,

$$\begin{aligned} \|\eta_{J_2^c}\|_1 - \|\hat{\eta}_{J_2^c}\|_1 + \|\hat{\eta}_{J_2^c} - \eta_{J_2^c}\|_1 & \leq 2\|\eta_{J_2^c}\|_1 = 2 \sum_{i=1}^N |\eta_i| 1\{|\eta_i| < \Xi_2\} \\ & < 2\Xi_2^{1-\nu} \sum_{i=1}^N |\eta_i|^\nu 1\{|\eta_i| < \Xi_2\} \leq 2E \Xi_2^{1-\nu}. \end{aligned}$$

We proved (8.1). (8.2) follows trivially from this.  $\square$

## 8.2 Deterministic Oracle Inequalities

Here we provide oracle inequalities which are valid on  $\mathcal{A} \cap \mathcal{B}$ .

**Lemma 4.** *Let  $E_1 = N^{-\frac{\nu}{2}}E$ . Choose  $\Xi_1 = \lambda_N/(NT)$  and  $\Xi_2 = \lambda/(\sqrt{NT})$ . Let Assumption 2 hold. Then on the event  $\mathcal{A} \cap \mathcal{B}$  one has for any positive constant  $\lambda$ ,*

$$\begin{aligned}\|\Pi(\hat{\gamma} - \gamma)\|^2 &\leq \left(\frac{240}{\kappa^2} + 40\right) \frac{\lambda}{\sqrt{N}} E \left(\frac{\lambda}{\sqrt{NT}}\right)^{1-\nu} \\ \|\hat{\alpha} - \alpha\|_1 &\leq \left(\frac{240}{\kappa^2} + 40\right) \frac{1}{\sqrt{N}} E \left(\frac{\lambda}{\sqrt{NT}}\right)^{1-\nu} \\ \|\hat{\eta} - \eta\|_1 &\leq \left(\frac{240}{\kappa^2} + 40\right) E \left(\frac{\lambda}{\sqrt{NT}}\right)^{1-\nu}.\end{aligned}$$

*Proof.* By (8.1) of Lemma 3, which is valid on  $\mathcal{A}$ ,

$$\|\Pi(\hat{\gamma} - \gamma)\|^2 \leq 4\lambda \|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1 + 4\frac{\lambda}{\sqrt{N}} \|\hat{\eta}_{J_2} - \eta_{J_2}\|_1 + 4\lambda E_1 \Xi_1^{1-\nu} + 4\frac{\lambda}{\sqrt{N}} E \Xi_2^{1-\nu}. \quad (8.4)$$

Consider the auxiliary event

$$\mathcal{C} := \left\{ E_1 \Xi_1^{1-\nu} + \frac{1}{\sqrt{N}} E \Xi_2^{1-\nu} \leq \frac{1}{4} \|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1 + \frac{1}{4\sqrt{N}} \|\hat{\eta}_{J_2} - \eta_{J_2}\|_1 \right\}.$$

On the event  $\mathcal{A} \cap \mathcal{C}$ , from (8.2) of Lemma 3, we have

$$\|\hat{\alpha}_{J_1^c} - \alpha_{J_1^c}\|_1 + \frac{1}{\sqrt{N}} \|\hat{\eta}_{J_2^c} - \eta_{J_2^c}\|_1 \leq 4\|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1 + 4\frac{1}{\sqrt{N}} \|\hat{\eta}_{J_2} - \eta_{J_2}\|_1. \quad (8.5)$$

In order to apply the compatibility condition, re-parametrise the vector  $\delta$  in the definition of the compatibility condition as follows. Let  $b^1$  and  $b^2$  be  $p \times 1$  and  $N \times 1$  vectors, respectively, with  $b = (b^1, b^2)'$  defined as

$$\begin{pmatrix} b^1 \\ b^2 \end{pmatrix} := \begin{pmatrix} I_p & 0 \\ 0 & \sqrt{N}I_N \end{pmatrix} \begin{pmatrix} \delta_1 \\ \delta_2 \end{pmatrix}.$$

Hence, that  $\kappa^2(\Psi_N, r_1, r_2)$  is bounded away from zero for integers  $r_1 \in \{1, \dots, p\}$  and  $r_2 \in \{1, \dots, N\}$  is equivalent to

$$\kappa^2(\Psi_N, r_1, r_2) := \min_{\substack{R_1 \subseteq \{1, \dots, p\}, |R_1| \leq r_1 \\ R_2 \subseteq \{1, \dots, N\}, |R_2| \leq r_2 \\ R := R_1 \cup (R_2 + p)}} \min_{\substack{b \in \mathbb{R}^{p+N} \setminus \{0\}, \\ \|b_{R_1^c}^1\|_1 + \frac{1}{\sqrt{N}} \|b_{R_2^c}^2\|_1 \\ \leq 4\|b_{R_1}^1\|_1 + \frac{4}{\sqrt{N}} \|b_{R_2}^2\|_1}} \frac{\|\Pi b\|^2}{2} > 0. \quad (8.6)$$

By (8.5), our estimator satisfies the constraint of the just introduced version of the compatibility condition and so

$$\begin{aligned}\|\Pi(\hat{\gamma} - \gamma)\|^2 &\geq \frac{\kappa^2(\Psi_N, s_1, s_2)NT}{s_1 + s_2} \left\| \begin{pmatrix} \hat{\alpha}_{J_1} - \alpha_{J_1} \\ (\hat{\eta}_{J_2} - \eta_{J_2})/\sqrt{N} \end{pmatrix} \right\|_1^2 \\ &\geq \frac{\kappa^2(\Psi_N, s_1, s_2)NT}{s_1 + s_2} \left( \|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1^2 + \frac{1}{N} \|\hat{\eta}_{J_2} - \eta_{J_2}\|_1^2 \right) \\ &\geq \frac{\kappa^2 NT}{2(s_1 + s_2)} \left( \|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1^2 + \frac{1}{N} \|\hat{\eta}_{J_2} - \eta_{J_2}\|_1^2 \right),\end{aligned}$$

where the last inequality is valid on  $\mathcal{B}$ . Hence, on  $\mathcal{A} \cap \mathcal{B} \cap \mathcal{C}$  upon combining with (8.4) one has,

$$\begin{aligned} & \frac{\kappa^2 NT}{2(s_1 + s_2)} \left( \|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1^2 + \frac{1}{N} \|\hat{\eta}_{J_2} - \eta_{J_2}\|_1^2 \right) \\ & \leq 4\lambda \|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1 + \frac{4\lambda}{\sqrt{N}} \|\hat{\eta}_{J_2} - \eta_{J_2}\|_1 + 4\lambda E_1 \Xi_1^{1-\nu} + \frac{4\lambda}{\sqrt{N}} E \Xi_2^{1-\nu} \\ & \leq 5\lambda \|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1 + \frac{5\lambda}{\sqrt{N}} \|\hat{\eta}_{J_2} - \eta_{J_2}\|_1, \end{aligned}$$

which, since  $\kappa^2 > 0$  by Assumption 2, is equivalent to

$$\|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1^2 - \frac{10\lambda(s_1 + s_2)}{\kappa^2 NT} \|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1 + \frac{1}{N} \|\hat{\eta}_{J_2} - \eta_{J_2}\|_1^2 - \frac{10\lambda(s_1 + s_2)}{\kappa^2 N^{3/2} T} \|\hat{\eta}_{J_2} - \eta_{J_2}\|_1 \leq 0.$$

Let  $x = \|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1$ ,  $y = \|\hat{\eta}_{J_2} - \eta_{J_2}\|_1$ ,  $a = \frac{10\lambda(s_1 + s_2)}{\kappa^2 NT}$ ,  $b = \frac{1}{N}$  and  $c = \frac{10\lambda(s_1 + s_2)}{\kappa^2 N^{3/2} T}$ . Thus one has

$$x^2 - ax + by^2 - cy \leq 0. \quad (8.7)$$

First bound  $x = \|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1$ . For every  $y$  the values of  $x$  that satisfy the above quadratic inequality form an interval in  $\mathbb{R}_+$ . The right end point of this interval is the desired upper bound on  $x$ . Clearly, by the solution formula for the roots of a second degree polynomial, this right end point is a decreasing function in  $by^2 - cy$ . Hence, we first minimize the polynomial  $by^2 - cy$  to find the largest possible value of  $x$  which satisfies (8.7). This yields  $y = c/2b$  and the corresponding value of  $by^2 - cy$  is  $-c^2/(4b)$ . Hence, our desired upper bound on  $x$  is the largest solution of  $x^2 - ax - \frac{c^2}{4b} \leq 0$ . By the standard solution formula for the roots of a quadratic polynomial this yields

$$\|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1 = x \leq \frac{a + \sqrt{a^2 + c^2/b}}{2} \leq a + \frac{c}{2\sqrt{b}}. \quad (8.8)$$

Switching the roles of  $x$  and  $y$ , one gets a similar bound on  $y = \|\hat{\eta}_{J_2} - \eta_{J_2}\|_1$ , namely

$$\|\hat{\eta}_{J_2} - \eta_{J_2}\|_1 = y \leq \frac{c + \sqrt{c^2 + ba^2}}{2b} \leq \frac{c}{b} + \frac{a}{2\sqrt{b}}. \quad (8.9)$$

Inserting the definitions of  $a, b$  and  $c$  into (8.8) and (8.9), we get

$$\|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1 \leq \frac{15\lambda(s_1 + s_2)}{\kappa^2 NT} \quad (8.10)$$

$$\|\hat{\eta}_{J_2} - \eta_{J_2}\|_1 \leq \frac{15\lambda(s_1 + s_2)}{\kappa^2 N^{1/2} T}. \quad (8.11)$$

Therefore, on  $\mathcal{A} \cap \mathcal{B} \cap \mathcal{C}$ , it follows from (8.1) that

$$\begin{aligned}
\|\Pi(\hat{\gamma} - \gamma)\|^2 &\leq 4\lambda\|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1 + \frac{4\lambda}{\sqrt{N}}\|\hat{\eta}_{J_2} - \eta_{J_2}\|_1 + 4\lambda E_1 \Xi_1^{1-\nu} + \frac{4\lambda}{\sqrt{N}} E \Xi_2^{1-\nu} \\
&\leq \frac{120\lambda^2(s_1 + s_2)}{\kappa^2 NT} + 4\lambda E_1 \Xi_1^{1-\nu} + \frac{4\lambda}{\sqrt{N}} E \Xi_2^{1-\nu} \\
\|\hat{\alpha} - \alpha\|_1 &\leq 4\|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1 + \frac{4}{\sqrt{N}}\|\hat{\eta}_{J_2} - \eta_{J_2}\|_1 + 4E_1 \Xi_1^{1-\nu} + \frac{4}{\sqrt{N}} E \Xi_2^{1-\nu} \\
&\leq \frac{120\lambda(s_1 + s_2)}{\kappa^2 NT} + 4E_1 \Xi_1^{1-\nu} + \frac{4}{\sqrt{N}} E \Xi_2^{1-\nu} \\
\|\hat{\eta} - \eta\|_1 &\leq 4\sqrt{N}\|\hat{\alpha}_{J_1} - \alpha_{J_1}\|_1 + 4\|\hat{\eta}_{J_2} - \eta_{J_2}\|_1 + 4\sqrt{N} E_1 \Xi_1^{1-\nu} + 4E \Xi_2^{1-\nu} \\
&\leq \frac{120\lambda(s_1 + s_2)}{\kappa^2 \sqrt{N} T} + 4\sqrt{N} E_1 \Xi_1^{1-\nu} + 4E \Xi_2^{1-\nu}.
\end{aligned}$$

On  $\mathcal{A} \cap \mathcal{C}^c$  one has trivial oracle inequalities via (8.1) of Lemma 3. To be precise,

$$\begin{aligned}
\|\Pi(\hat{\gamma} - \gamma)\|^2 &\leq 20\lambda E_1 \Xi_1^{1-\nu} + 20\lambda \frac{E \Xi_2^{1-\nu}}{\sqrt{N}} \\
\|\hat{\alpha} - \alpha\|_1 &\leq 20E_1 \Xi_1^{1-\nu} + 20 \frac{E \Xi_2^{1-\nu}}{\sqrt{N}} \\
\|\hat{\eta} - \eta\|_1 &\leq 20\sqrt{N} E_1 \Xi_1^{1-\nu} + 20E \Xi_2^{1-\nu}.
\end{aligned}$$

These inequalities are valid on event  $\mathcal{A} \cap \mathcal{B} \cap \mathcal{C}^c$  too. Synchronising constants, using that  $(\mathcal{A} \cap \mathcal{B} \cap \mathcal{C}^c) \cup (\mathcal{A} \cap \mathcal{B} \cap \mathcal{C}) = \mathcal{A} \cap \mathcal{B}$ , and recognising that

$$\begin{aligned}
s_1 &:= \sum_{j=1}^p 1\{|\alpha_j| \geq \Xi_1\} = \sum_{j=1}^p 1\{|\alpha_j|^\nu \geq \Xi_1^\nu\} \leq E_1 \Xi_1^{-\nu} \\
s_2 &:= \sum_{i=1}^N 1\{|\eta_i| \geq \Xi_2\} = \sum_{i=1}^N 1\{|\eta_i|^\nu \geq \Xi_2^\nu\} \leq E \Xi_2^{-\nu},
\end{aligned}$$

we arrive at

$$\begin{aligned}
\|\Pi(\hat{\gamma} - \gamma)\|^2 &\leq \frac{120\lambda^2(E_1 \Xi_1^{-\nu} + E \Xi_2^{-\nu})}{\kappa^2 NT} + 20\lambda E_1 \Xi_1^{1-\nu} + 20\lambda \frac{E \Xi_2^{1-\nu}}{\sqrt{N}} \\
\|\hat{\alpha} - \alpha\|_1 &\leq \frac{120\lambda(E_1 \Xi_1^{-\nu} + E \Xi_2^{-\nu})}{\kappa^2 NT} + 20E_1 \Xi_1^{1-\nu} + 20 \frac{E \Xi_2^{1-\nu}}{\sqrt{N}} \\
\|\hat{\eta} - \eta\|_1 &\leq \frac{120\lambda(E_1 \Xi_1^{-\nu} + E \Xi_2^{-\nu})}{\kappa^2 \sqrt{N} T} + 20\sqrt{N} E_1 \Xi_1^{1-\nu} + 20E \Xi_2^{1-\nu}.
\end{aligned}$$

The deterministic oracle inequalities follow upon choosing  $\Xi_1 = \frac{\lambda}{NT}$ ,  $\Xi_2 = \frac{\lambda}{\sqrt{N} T}$  and  $E_1 = N^{-\nu/2} E$ .

□

### 8.3 A Lower Bound on $\mathbb{P}(\mathcal{A})$

For the proof of Lemma 5 below, we shall use Orlicz norms as defined in van der Vaart and Wellner (1996): Let  $\psi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  be a non-decreasing, convex function with  $\psi(0) = 0$  and

$\lim_{x \rightarrow \infty} \psi(x) = \infty$ . Then, the Orlicz norm of a random variable  $X$  is given by

$$\|X\|_\psi = \inf \left\{ C > 0 : \mathbb{E} \psi(|X|/C) \leq 1 \right\},$$

where, as usual,  $\inf \emptyset = \infty$ . We will use Orlicz norms for  $\psi(x) = \psi_b(x) = e^{x^b} - 1$  for various values of  $b$ . The following lemma provides a lower bound on the probability of  $\mathcal{A}$ .

**Lemma 5.** *Let  $\lambda = \sqrt{4MNT(\log(p \vee N))^3}$  for some  $M > 0$ . By Assumptions 1 and 3, we have*

$$\mathbb{P}(\mathcal{A}) \geq 1 - Ap^{1-BM^{1/3}} - AN^{1-BM^{1/3}},$$

for positive constants  $A$  and  $B$ .

*Proof.* Consider the event  $\{\|Z'\varepsilon\|_\infty > \lambda/2\}$  first. To this end, let  $z_{j,l}$  denote the  $j$ th entry of the  $l$ th column of  $Z$ , i.e. the  $j$ th entry of  $(z_{1,1,l}, z_{1,2,l}, \dots, z_{1,T,l}, z_{2,1,l}, \dots, z_{N,T,l})'$ . Similarly, we write  $\varepsilon_j$  for the  $j$ th entry of  $\varepsilon$ . Now note that  $j \mapsto (\lceil \frac{j}{T} \rceil, j - \lfloor \frac{j}{T} \rfloor T)$  is a bijection from  $\{1, \dots, NT\}$  to  $\{1, \dots, N\} \times \{1, \dots, T\}$  where  $\lfloor x \rfloor$  denotes the greatest integer strictly less than  $x$  and  $\lceil x \rceil$  the smallest integer greater than or equal to  $x \in \mathbb{R}$ . In case the  $l$ th column of  $Z$  corresponds to one of the lags of the left hand side variable, assume for concreteness the  $k$ th lag, define  $\mathcal{F}_n = \sigma(y_{\lceil \frac{j}{T} \rceil, j - \lfloor \frac{j}{T} \rfloor T}, \dots, y_{\lceil \frac{j}{T} \rceil, j - \lfloor \frac{j}{T} \rfloor T - L}, \varepsilon_{\lceil \frac{j}{T} \rceil, j - \lfloor \frac{j}{T} \rfloor T}, 1 \leq j \leq n, y_{\lceil \frac{n+1}{T} \rceil, 0}, \dots, y_{\lceil \frac{n+1}{T} \rceil, 1-L})$  and  $S_{n,l} = \sum_{j=1}^n z_{j,l} \varepsilon_j = \sum_{j=1}^n y_{\lceil \frac{j}{T} \rceil, j - \lfloor \frac{j}{T} \rfloor T - k} \varepsilon_{\lceil \frac{j}{T} \rceil, j - \lfloor \frac{j}{T} \rfloor T}$ . Thus,

$$\begin{aligned} \mathbb{E}[S_{n,l} | \mathcal{F}_{n-1}] &= \sum_{j=1}^{n-1} y_{\lceil \frac{j}{T} \rceil, j - \lfloor \frac{j}{T} \rfloor T - k} \varepsilon_{\lceil \frac{j}{T} \rceil, j - \lfloor \frac{j}{T} \rfloor T} + \mathbb{E}[y_{\lceil \frac{n}{T} \rceil, j - \lfloor \frac{n}{T} \rfloor T - k} \varepsilon_{\lceil \frac{n}{T} \rceil, n - \lfloor \frac{n}{T} \rfloor T} | \mathcal{F}_{n-1}] \\ &= S_{n-1,l} + y_{\lceil \frac{n}{T} \rceil, j - \lfloor \frac{n}{T} \rfloor T - k} \mathbb{E}[\varepsilon_{\lceil \frac{n}{T} \rceil, n - \lfloor \frac{n}{T} \rfloor T} | \mathcal{F}_{n-1}]. \end{aligned}$$

Using that  $(\lceil \frac{n}{T} \rceil, n - \lfloor \frac{n}{T} \rfloor T)$  is a unique pair  $(i, t) \in \{1, \dots, N\} \times \{1, \dots, T\}$  we have that

$$\mathbb{E}[\varepsilon_{\lceil \frac{n}{T} \rceil, n - \lfloor \frac{n}{T} \rfloor T} | \mathcal{F}_{n-1}] = \mathbb{E}[\varepsilon_{i,t} | \mathcal{F}_{n-1}] = \mathbb{E}[\varepsilon_{i,t} | \sigma(y_{i,t-1}, \dots, y_{i,1-L}, \varepsilon_{i,t-1}, \dots, \varepsilon_{i,1})]$$

<sup>8</sup>where the last equality follows from the assumption of independence across  $1 \leq i \leq N$  (Assumption 1). By Assumption 1, this conditional expectation equals zero as the  $\varepsilon_{i,s}$  is a linear function of  $y_{i,s}, \dots, y_{i,s-L}$  and  $x_{i,s}$  for  $1 \leq s \leq t-1$ . Thus,  $S_{n,l}$  is a martingale with mean zero (i.e., the increments are martingale differences by the above argument). A similar argument applies when the  $l$ th column of  $Z$  equals  $(x_{1,1,k}, \dots, x_{1,T,k}, x_{2,1,k}, \dots, x_{N,T,k})'$  for some  $1 \leq k \leq p_x$  such that every row of  $Z'\varepsilon$  is a zero mean martingale.

Next, note that by Assumption 3, for all  $1 \leq j \leq NT$ ,  $1 \leq l \leq p$  and  $\epsilon > 0$ , one has

$$\mathbb{P}(|z_{j,l} \varepsilon_j| \geq \epsilon) \leq \mathbb{P}(|z_{j,l}| \geq \sqrt{\epsilon}) + \mathbb{P}(|\varepsilon_j| \geq \sqrt{\epsilon}) \leq Ke^{-C\epsilon}.$$

---

<sup>8</sup>For  $t = 1$ , the last expression in the above display is to be read as absence of conditioning on the error terms.

It follows from Lemma 2.2.1 in van der Vaart and Wellner (1996) that  $\|z_{j,l}\varepsilon_j\|_{\psi_1} \leq (1+K)/C$ . Then, by the definition of the Orlicz norm,  $\mathbb{E}[e^{C/(1+K)|z_{j,l}\varepsilon_j|}] \leq 2$ . Now use Proposition 2 in Appendix F with  $D = C/(1+K)$ ,  $\alpha = 1/3$  and  $C_1 = 2$  to conclude

$$\mathbb{P}\left(\|Z'\varepsilon\|_\infty > \frac{\lambda}{2}\right) \leq \sum_{l=1}^p \mathbb{P}\left(\left|\sum_{j=1}^{NT} z_{j,l}\varepsilon_j\right| > \frac{\lambda}{2NT}NT\right) = pAe^{-B\log(p\vee N)M^{1/3}} \leq Ap^{1-BM^{1/3}}.$$

Note also that the upper bound of the preceding probability becomes arbitrarily small for sufficiently large  $N$  and  $M$  such that we also conclude

$$\|Z'\varepsilon\|_\infty = O_p(\lambda). \quad (8.12)$$

Next, consider the event  $\{\|D'\varepsilon\|_\infty > \lambda/(2\sqrt{N})\}$ . Using Assumption 1 a small calculation shows that all entries of  $D'\varepsilon$  are zero mean martingales with respect to the natural filtration. As above, Assumption 3 and Lemma 2.2.1 in van der Vaart and Wellner (1996) yield  $\|\varepsilon_{i,t}\|_{\psi_2} \leq \left(\frac{1+K/2}{C}\right)^{1/2}$  such that by the second to last inequality on page 95 in van der Vaart and Wellner (1996) one has  $\|\varepsilon_{i,t}\|_{\psi_1} \leq \|\varepsilon_{i,t}\|_{\psi_2}(\log 2)^{-1/2} \leq \left(\frac{1+K/2}{C}\right)^{1/2}(\log 2)^{-1/2}$  for all  $i$  and  $t$ . Then using the definition of the Orlicz norm,  $\mathbb{E}[\exp((\frac{C}{1+K/2})^{1/2}(\log 2)^{1/2}|\varepsilon_{i,t}|)] \leq 2$  and Proposition 2 in Appendix F with  $D = \left(\frac{C}{1+K/2}\right)^{1/2}(\log 2)^{1/2}$ ,  $\alpha = 1/3$  and  $C_1 = 2$  implies

$$\mathbb{P}\left(\|D'\varepsilon\|_\infty > \frac{\lambda}{2\sqrt{N}}\right) \leq \sum_{i=1}^N \mathbb{P}\left(\left|\sum_{t=1}^T \varepsilon_{i,t}\right| > \frac{\lambda}{2\sqrt{NT}}T\right) \leq ANe^{-B(\log(p\vee N)M^{1/3})} \leq AN^{1-BM^{1/3}}.$$

Note also that the upper bound of the preceding probability becomes arbitrarily small for sufficiently large  $N$  and  $M$ , such that we may also conclude

$$\|D'\varepsilon\|_\infty = O_p\left(\frac{\lambda}{\sqrt{N}}\right). \quad (8.13)$$

□

#### 8.4 A Lower Bound on $\mathbb{P}(\mathcal{B})$

The following lemma shows that  $\kappa^2(\Psi_N, s_1, s_2)$  and  $\kappa^2(\Psi, s_1, s_2)$  are close if  $\Psi_N$  and  $\Psi$  are in some sense close.

**Lemma 6.** *Let  $A$  and  $B$  be two positive semidefinite  $(p+N) \times (p+N)$  matrices and  $\delta := \max_{1 \leq i, j \leq p+N} |A_{ij} - B_{ij}|$ . For any integers  $r_1 \in \{1, \dots, p\}$  and  $r_2 \in \{1, \dots, N\}$ , one has*

$$\kappa^2(B, r_1, r_2) \geq \kappa^2(A, r_1, r_2) - \delta 25(r_1 + r_2).$$

*Proof.* Let  $x$  be a  $(p+N) \times 1$  non-zero vector, satisfying  $\|x_{R^c}\|_1 \leq 4\|x_R\|_1$  for  $R = R_1 \cup (R_2 + p)$  where  $R_1 \subseteq \{1, \dots, p\}$  with  $|R_1| \leq r_1$ , and  $R_2 \subseteq \{1, \dots, N\}$  with  $|R_2| \leq r_2$ . Now,

$$\begin{aligned} |x'Ax - x'Bx| &= |x'(A-B)x| \leq \|x\|_1 \|(A-B)x\|_\infty \leq \|x\|_1^2 \delta = \delta (\|x_R\|_1 + \|x_{R^c}\|_1)^2 \\ &\leq \delta (\|x_R\|_1 + 4\|x_R\|_1)^2 = \delta 25\|x_R\|_1^2. \end{aligned}$$

Hence,

$$\frac{x' B x}{\frac{1}{r_1+r_2} \|x_R\|_1^2} \geq \frac{x' A x}{\frac{1}{r_1+r_2} \|x_R\|_1^2} - \delta 25(r_1 + r_2) \geq \kappa^2(A, r_1, r_2) - \delta 25(r_1 + r_2),$$

where the last inequality is true because of the definition of  $\kappa^2(A, r_1, r_2)$ . Minimising the left-hand side over non-zero  $x$  satisfying  $\|x_{R^c}\|_1 \leq 4\|x_R\|_1$  yields the claim.  $\square$

Define

$$\tilde{\mathcal{B}} = \left\{ \max_{1 \leq i, j \leq p+N} |\Psi_{N,ij} - \Psi_{ij}| \leq \frac{\kappa^2(\Psi, s_1, s_2)}{100E \left( \frac{\lambda}{\sqrt{NT}} \right)^{-\nu}} \right\}.$$

Setting  $A = \Psi$ ,  $B = \Psi_N$  it follows from Lemma 6 that  $\tilde{\mathcal{B}} \subseteq \mathcal{B}$ . Thus, we just need to find a lower bound on  $\mathbb{P}(\tilde{\mathcal{B}})$  in order to prove Theorem 1.

**Lemma 7.** *Let Assumptions 1, 2 and 3 hold. Assume that  $E \left( \frac{\lambda}{\sqrt{NT}} \right)^{-\nu} \lesssim \sqrt{N}$ . Then, there exist positive constants  $A, B$  such that*

$$\mathbb{P}(\mathcal{B}^c) \leq \mathbb{P}(\tilde{\mathcal{B}}^c) \leq A(p^2 + pN) \exp \left( -B \left\{ N / \left[ E \left( \frac{\lambda}{\sqrt{NT}} \right)^{-\nu} \right]^2 \right\}^{1/3} \right).$$

*Proof.* Since the lower right  $N \times N$  blocks of  $\Psi_N$  and  $\Psi$  are identical, it suffices to bound the entries of  $\frac{1}{NT} Z' Z - \frac{1}{NT} \mathbb{E}[Z' Z]$  and  $\frac{1}{T\sqrt{N}} Z' D$ . A typical element of  $\frac{1}{NT} Z' Z - \frac{1}{NT} \mathbb{E}[Z' Z]$  is of the form  $\frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T (z_{i,t,l} z_{i,t,k} - \mathbb{E}[z_{i,t,l} z_{i,t,k}])$  for some  $l, k \in \{1, \dots, p\}$ . By Assumption 3 we have for every  $\epsilon > 0$

$$\mathbb{P}(|z_{i,t,l} z_{i,t,k}| \geq \epsilon) \leq \mathbb{P}(|z_{i,t,l}| \geq \sqrt{\epsilon}) + \mathbb{P}(|z_{i,t,k}| \geq \sqrt{\epsilon}) \leq K e^{-C\epsilon}.$$

It follows from Lemma 2.2.1 in van der Vaart and Wellner (1996) that  $\|z_{i,t,l} z_{i,t,k}\|_{\psi_1} \leq (1+K)/C$ . Hence, by subadditivity of the Orlicz norm and Jensen's inequality

$$\left\| \frac{1}{T} \sum_{t=1}^T (z_{i,t,l} z_{i,t,k} - \mathbb{E}[z_{i,t,l} z_{i,t,k}]) \right\|_{\psi_1} \leq 2 \max_{1 \leq t \leq T} \|z_{i,t,l} z_{i,t,k}\|_{\psi_1} \leq \frac{2(1+K)}{C}.$$

Thus, by the definition of the Orlicz norm,  $\mathbb{E} \exp \left( \frac{C}{2(1+K)} \left| \frac{1}{T} \sum_{t=1}^T (z_{i,t,l} z_{i,t,k} - \mathbb{E}[z_{i,t,l} z_{i,t,k}]) \right| \right) \leq 2$ . Using independence across  $i$  (Assumption 1) to invoke Proposition 2 in Appendix F with  $D = \frac{C}{2(1+K)}$ ,  $\alpha = 1/3$  and  $C_1 = 2$  such that for every  $x \gtrsim \frac{1}{\sqrt{N}}$

$$\mathbb{P} \left( \left| \sum_{i=1}^N \frac{1}{T} \sum_{t=1}^T (z_{i,t,l} z_{i,t,k} - \mathbb{E}[z_{i,t,l} z_{i,t,k}]) \right| \geq Nx \right) \leq A e^{-B(x^2 N)^{1/3}}, \quad (8.14)$$

for positive constants  $A$  and  $B$ .

Next, consider  $\frac{1}{T\sqrt{N}} Z' D$ . A typical element can be written as  $\frac{1}{\sqrt{NT}} \sum_{t=1}^T z_{i,t,l}$  for some  $i \in \{1, \dots, N\}$  and  $l \in \{1, \dots, p\}$ . By Assumption 3, we have  $\mathbb{P}(|z_{i,t,l}| \geq \epsilon) \leq \frac{1}{2} K e^{-C\epsilon^2}$  for all



$\epsilon > 0$  and it follows from Lemma 2.2.1 in van der Vaart and Wellner (1996) that  $\|z_{i,t,l}\|_{\psi_2} \leq \left(\frac{1+K/2}{C}\right)^{1/2}$ . Hence,

$$\left\| \frac{1}{\sqrt{NT}} \sum_{t=1}^T z_{i,t,l} \right\|_{\psi_2} \leq \frac{1}{\sqrt{N}} \max_{1 \leq t \leq T} \|z_{i,t,l}\|_{\psi_2} \leq \frac{1}{\sqrt{N}} \left(\frac{1+K/2}{C}\right)^{1/2} =: \frac{C'}{\sqrt{N}}.$$

Thus, it follows by Markov's inequality, positivity and increasingness of  $\psi_2(x)$ , as well as  $1 \wedge \psi_2(x)^{-1} = 1 \wedge (e^{x^2} - 1)^{-1} \leq 2e^{-x^2}$  that for any  $x > 0$

$$\mathbb{P} \left( \left| \frac{1}{\sqrt{NT}} \sum_{t=1}^T z_{i,t,l} \right| > x \right) \leq 1 \wedge \frac{1}{e^{(x\sqrt{N}/C')^2} - 1} \leq 2e^{-\frac{Nx^2}{C'^2}} \leq Ae^{-Bx^2N}, \quad (8.15)$$

where the last estimate follows by choosing  $A$  and  $B$  sufficiently large/small for (8.14) and (8.15) both to be valid. Setting  $x = \frac{\kappa^2}{100E\left(\frac{\lambda}{\sqrt{NT}}\right)^{-\nu}} = \frac{\kappa^2}{100} \frac{1}{E\left(\frac{\lambda}{\sqrt{NT}}\right)^{-\nu}}$ , using that  $\frac{1}{E\left(\frac{\lambda}{\sqrt{NT}}\right)^{-\nu}} \gtrsim \frac{1}{\sqrt{N}}$  and  $\kappa^2$  being bounded away from 0 (Assumption 2), we have

$$\begin{aligned} \mathbb{P}(\mathcal{B}^c) &\leq \mathbb{P}(\tilde{\mathcal{B}}^c) = \mathbb{P} \left( \max_{1 \leq i, j \leq p+N} |\Psi_{N,ij} - \Psi_{ij}| > x \right) \\ &\leq A(p^2 + pN) \left[ \exp \left( -B \left\{ \left[ \frac{\kappa^2/100}{E\left(\frac{\lambda}{\sqrt{NT}}\right)^{-\nu}} \right]^2 N \right\}^{1/3} \right) \vee \exp \left( -B \left[ \frac{\kappa^2/100}{E\left(\frac{\lambda}{\sqrt{NT}}\right)^{-\nu}} \right]^2 N \right) \right] \\ &\leq A(p^2 + pN) \exp \left( -B \left\{ N / \left[ E\left(\frac{\lambda}{\sqrt{NT}}\right)^{-\nu} \right]^2 \right\}^{1/3} \right) \end{aligned}$$

where the last estimate has merged  $(\kappa^2/100)^{2/3}$  into  $B$ .  $\square$

*Proof of Theorem 1.* Theorem 1 follows by combining Lemmas 4, 5, and 7.  $\square$

## 8.5 Rates of Convergence

**Corollary 1.** *Let the conditions of Theorem 1 hold. For large enough  $M > 0$  and assuming  $\frac{(\log(p \vee N))^3 E^2 [(\log(p \vee N))^3 / T]^{-\nu}}{N} = o(1)$ , we have the following stochastic orders valid uniformly over  $\mathcal{F}(\nu, E)$ .*

$$\begin{aligned} \frac{1}{NT} \|\Pi(\hat{\gamma} - \gamma)\|^2 &= O_p \left( \frac{\lambda}{\sqrt{N}NT} E \left( \frac{\lambda}{\sqrt{NT}} \right)^{1-\nu} \right), \\ \|\hat{\alpha} - \alpha\|_1 &= O_p \left( \frac{1}{\sqrt{N}} E \left( \frac{\lambda}{\sqrt{NT}} \right)^{1-\nu} \right), \\ \|\hat{\eta} - \eta\|_1 &= O_p \left( E \left( \frac{\lambda}{\sqrt{NT}} \right)^{1-\nu} \right). \end{aligned}$$

*Proof of Corollary 1.* First note that  $Ap^{1-AM^{1/3}}$  and  $AN^{1-AM^{1/3}}$  become arbitrarily small for large enough  $M > 0$ . By  $\frac{(\log(p \vee N))^3 E^2 [(\log(p \vee N))^3 / T]^{-\nu}}{N} = o(1)$ ,

$$A(p^2 + pN) \exp \left( -A \left\{ \frac{N}{E^2 [(\log(p \vee N))^3 / T]^{-\nu}} \right\}^{1/3} \right) \rightarrow 0.$$

Thus the lower bound on the probability in Theorem 1 goes to one as  $N, T, p \rightarrow \infty$  for large enough  $M > 0$  and the conclusion follows from Theorem 1.  $\square$

## 9 Appendix C

### 9.1 Properties of the Nodewise Lasso

The following lemma gives the rates of the uniform prediction and estimation errors for nodewise regression. It is used in the proof of Lemma 9.

**Lemma 8.** *Let Assumptions 1, 3 and 4 hold. Let  $\lambda_{node} = \sqrt{16M(\log p)^3/N}$  for some  $M > 0$ . For  $M$  sufficiently large, we have*

$$\max_{j \in H_1} \frac{1}{NT} \|Z_{-j}(\hat{\phi}_j - \phi_j)\|^2 = O_p\left(\bar{G}\lambda_{node}^{2-\vartheta}\right) \quad (9.1)$$

$$\max_{j \in H_1} \|\hat{\phi}_j - \phi_j\|_1 = O_p\left(\bar{G}\lambda_{node}^{1-\vartheta}\right) \quad (9.2)$$

$$\max_{j \in H_1} \frac{1}{NT} \|Z'_{-j}\zeta_j\|_\infty = O_p(\lambda_{node}). \quad (9.3)$$

*Proof.* We say that a  $(p-1) \times (p-1)$  matrix  $A$  satisfies the *compatibility condition*  $CC(r)$  for some integer  $r \in \{1, \dots, p-1\}$  if

$$\kappa^2(A, r) := \min_{\substack{R \subseteq \{1, \dots, p-1\} \\ |R| \leq r}} \min_{\substack{\delta \in \mathbb{R}^{p-1} \setminus \{0\} \\ \|\delta_{R^c}\|_1 \leq 3\|\delta_R\|_1}} \frac{\delta' A \delta}{\frac{1}{r} \|\delta_R\|_1^2} > 0.$$

Define a  $(p-1) \times 1$  vector  $\phi_j^*$  such that

$$\phi_{j,k}^* := \phi_{j,k} 1_{\{|\phi_{j,k}| \geq \lambda_{node}\}} \quad k = 1, \dots, p-1$$

and its active set  $J_j^*$  as well as its sparsity index  $s_j^*$

$$J_j^* := \{k : \phi_{j,k}^* \neq 0, k = 1, \dots, p-1\} \quad 1 \leq s_j^* := |J_j^*| \leq p-1.$$

Consider the events

$$\mathcal{D} = \left\{ \max_{j \in H_1} \frac{1}{NT} \|Z'_{-j}\zeta_j\|_\infty \leq \frac{\lambda_{node}}{4} \right\},$$

$$\mathcal{E}_j = \left\{ \kappa^2\left(\frac{1}{NT} Z'_{-j} Z_{-j}, s_j^*\right) \geq \frac{\kappa^2(\Psi_{Z, -j, -j}, s_j^*)}{2} \right\},$$

and

$$\mathcal{F} = \left\{ \left\| \frac{1}{NT} Z' Z - \Psi_Z \right\|_\infty \leq \lambda_{node} \right\}.$$

Using the same technique as in Section 6.2.3 of Bühlmann and van de Geer (2011), we arrive at the following oracle inequality, which is almost the same as the one on the top of p111 of

Bühlmann and van de Geer (2011): for each  $j \in H_1$ , on  $\mathcal{D} \cap \mathcal{E}_j$

$$\begin{aligned} \frac{1}{NT} \|Z_{-j}(\hat{\phi}_j - \phi_j)\|^2 + \lambda_{node} \|\hat{\phi}_j - \phi_j\|_1 &\leq \frac{3}{NT} \|Z_{-j}(\phi_j^* - \phi_j)\|^2 + \frac{48\lambda_{node}^2 s_j^*}{\kappa^2(\frac{1}{NT} Z'_{-j} Z_{-j}, s_j^*)} + \lambda_{node} \|\phi_j^* - \phi_j\|_1 \\ &\leq \frac{3}{NT} \|Z_{-j}(\phi_j^* - \phi_j)\|^2 + \frac{96\lambda_{node}^2 s_j^*}{\kappa^2(\Psi_Z, s_j^*)} + \lambda_{node} \|\phi_j^* - \phi_j\|_1 \end{aligned} \quad (9.4)$$

where the second inequality is due to event  $\mathcal{E}_j$  and that  $\kappa^2(\Psi_{Z, -j, -j}, r) \geq \kappa^2(\Psi_Z, r)$  for all  $j = 1, \dots, p$  and  $r = 1, \dots, p-1$ .

We now bound the three terms on the right hand side of (9.4). Let  $b_j := \phi_j^* - \phi_j$ .

$$\begin{aligned} \frac{1}{NT} \|Z_{-j}(\phi_j^* - \phi_j)\|^2 &= b_j' \Psi_{Z, -j, -j} b_j + b_j' \left( \frac{1}{NT} Z'_{-j} Z_{-j} - \Psi_{Z, -j, -j} \right) b_j \\ &\leq \max_{\text{eval}}(\Psi_{Z, -j, -j}) \|b_j\|^2 + \left\| \frac{1}{NT} Z'_{-j} Z_{-j} - \Psi_{Z, -j, -j} \right\|_{\infty} \|b_j\|_1^2 \leq \max_{\text{eval}}(\Psi_Z) \|b_j\|^2 + \lambda_{node} \|b_j\|_1^2 \end{aligned}$$

where the last inequality holds on event  $\mathcal{F}$ . Note that

$$\|b_j\|^2 = \sum_{k=1}^{p-1} |\phi_{j,k}|^2 1\{|\phi_{j,k}| < \lambda_{node}\} \leq \lambda_{node}^{2-\vartheta} \sum_{k=1}^{p-1} |\phi_{j,k}|^{\vartheta} 1\{|\phi_{j,k}| < \lambda_{node}\} \leq G_j \lambda_{node}^{2-\vartheta}.$$

$$\|b_j\|_1 = \sum_{k=1}^{p-1} |\phi_{j,k}| 1\{|\phi_{j,k}| < \lambda_{node}\} \leq \lambda_{node}^{1-\vartheta} \sum_{k=1}^{p-1} |\phi_{j,k}|^{\vartheta} 1\{|\phi_{j,k}| < \lambda_{node}\} \leq G_j \lambda_{node}^{1-\vartheta}. \quad (9.5)$$

$$1 \leq s_j^* = \sum_{k=1}^{p-1} 1\{|\phi_{j,k}| \geq \lambda_{node}\} = \sum_{k=1}^{p-1} 1\{|\phi_{j,k}|^{\vartheta} \geq \lambda_{node}^{\vartheta}\} \leq G_j \lambda_{node}^{-\vartheta}. \quad (9.6)$$

Thus, for each  $j \in H_1$ , on  $\mathcal{D} \cap \mathcal{E}_j \cap \mathcal{F}$

$$\begin{aligned} &\frac{1}{NT} \|Z_{-j}(\hat{\phi}_j - \phi_j)\|^2 + \lambda_{node} \|\hat{\phi}_j - \phi_j\|_1 \\ &\leq \max_{\text{eval}}(\Psi_Z) G_j \lambda_{node}^{2-\vartheta} + G_j^2 \lambda_{node}^{3-2\vartheta} + \frac{96}{\kappa^2(\Psi_Z, s_j^*)} G_j \lambda_{node}^{2-\vartheta} + G_j \lambda_{node}^{2-\vartheta} \\ &= \left( \max_{\text{eval}}(\Psi_Z) + \frac{96}{\kappa^2(\Psi_Z, s_j^*)} + 1 \right) G_j \lambda_{node}^{2-\vartheta} + G_j^2 \lambda_{node}^{3-2\vartheta} \end{aligned}$$

from where we can extract two oracle inequalities

$$\begin{aligned} \frac{1}{NT} \|Z_{-j}(\hat{\phi}_j - \phi_j)\|^2 &\leq \left( \max_{\text{eval}}(\Psi_Z) + \frac{96}{\kappa^2(\Psi_Z, s_j^*)} + 1 \right) G_j \lambda_{node}^{2-\vartheta} + G_j^2 \lambda_{node}^{3-2\vartheta}, \\ \|\hat{\phi}_j - \phi_j\|_1 &\leq \left( \max_{\text{eval}}(\Psi_Z) + \frac{96}{\kappa^2(\Psi_Z, s_j^*)} + 1 \right) G_j \lambda_{node}^{1-\vartheta} + G_j^2 \lambda_{node}^{2-2\vartheta}. \end{aligned}$$

As the oracle inequalities in the above display are valid simultaneously on  $\mathcal{D} \cap (\cap_{j \in H_1} \mathcal{E}_j) \cap \mathcal{F}$

we conclude that

$$\max_{j \in H_1} \frac{1}{NT} \|Z_{-j}(\hat{\phi}_j - \phi_j)\|^2 \leq \left( \max_{\text{eval}}(\Psi_Z) + \frac{96}{\min_{j \in H_1} \kappa^2(\Psi_Z, s_j^*)} + 1 \right) \bar{G} \lambda_{node}^{2-\vartheta} + \bar{G}^2 \lambda_{node}^{3-2\vartheta},$$

$$\max_{j \in H_1} \|\hat{\phi}_j - \phi_j\|_1 \leq \left( \text{maxeval}(\Psi_Z) + \frac{96}{\min_{j \in H_1} \kappa^2(\Psi_Z, s_j^*)} + 1 \right) \bar{G} \lambda_{node}^{1-\vartheta} + \bar{G}^2 \lambda_{node}^{2-2\vartheta}, \quad (9.7)$$

on  $\mathcal{D} \cap (\cap_{j \in H_1} \mathcal{E}_j) \cap \mathcal{F}$ .

Next, we establish a lower bound on the probability of  $\mathcal{D} \cap (\cap_{j \in H_1} \mathcal{E}_j) \cap \mathcal{F}$ . Consider  $\mathcal{D}$  first. A typical element of  $Z'_{-j} \zeta_j$  is of the form  $\sum_{i=1}^N \sum_{t=1}^T z_{i,t,l} \zeta_{j,i,t}$  for some  $l \neq j$ . By (3.5), one has  $\frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T z_{i,t,l} \zeta_{j,i,t} = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T (z_{i,t,l} \zeta_{j,i,t} - \mathbb{E}[z_{i,t,l} \zeta_{j,i,t}])$  for  $l \neq j$ . By Assumptions 3 and 4(c), it holds for any  $\epsilon > 0$  that

$$\mathbb{P}(|z_{i,t,l} \zeta_{j,i,t}| > \epsilon) \leq \mathbb{P}(|z_{i,t,l}| > \sqrt{\epsilon}) + \mathbb{P}(|\zeta_{j,i,t}| > \sqrt{\epsilon}) \leq K e^{-C\epsilon}.$$

such that Lemma 2.2.1 in van der Vaart and Wellner (1996) yields that  $\|z_{i,t,l} \zeta_{j,i,t}\|_{\psi_1} \leq (1 + K)/C$ . Therefore, by Jensen's inequality and subadditivity of the Orlicz norm

$$\left\| \frac{1}{T} \sum_{t=1}^T (z_{i,t,l} \zeta_{j,i,t} - \mathbb{E}[z_{i,t,l} \zeta_{j,i,t}]) \right\|_{\psi_1} \leq 2 \max_{1 \leq t \leq T} \|z_{i,t,l} \zeta_{j,i,t}\|_{\psi_1} \leq \frac{2(1+K)}{C}.$$

Using the definition of the Orlicz norm  $\mathbb{E} \exp \left( \frac{C}{2(1+K)} \left| \frac{1}{T} \sum_{t=1}^T (z_{i,t,l} \zeta_{j,i,t} - \mathbb{E}[z_{i,t,l} \zeta_{j,i,t}]) \right| \right) \leq 2$ . Using independence across  $i$  (Assumption 1) to invoke Proposition 2 in Appendix F with  $D = C/(1+K)$ ,  $\alpha = 1/3$ ,  $C_1 = 2$  and  $\epsilon = \lambda_{node}/4 \gtrsim \frac{1}{\sqrt{N}}$ , we conclude (using  $h_1 \leq p$ )

$$\begin{aligned} \mathbb{P} \left( \max_{j \in H_1} \frac{1}{NT} \|Z'_{-j} \zeta_j\|_\infty > \epsilon \right) &\leq h_1 p \mathbb{P} \left( \left| \sum_{i=1}^N \frac{1}{T} \sum_{t=1}^T (z_{i,t,l} \zeta_{j,i,t} - \mathbb{E}[z_{i,t,l} \zeta_{j,i,t}]) \right| > \epsilon N \right) \\ &\leq A h_1 p e^{-B(\epsilon^2 N)^{1/3}} \leq A p^2 e^{-B M^{1/3} \log p} = A p^{2-BM^{1/3}} \end{aligned}$$

for positive constants  $A$  and  $B$ . The upper bound of the preceding probability becomes arbitrarily small for  $M$  sufficiently large such that

$$\max_{j \in H_1} \frac{1}{NT} \|Z'_{-j} \zeta_j\|_\infty = O_p(\lambda_{node}),$$

which is (9.3). In order to provide a lower bound on the probability of  $(\cap_{j \in H_1} \mathcal{E}_j)$  define the event

$$\tilde{\mathcal{E}}_j := \left\{ \max_{1 \leq l, k \leq p-1} \left| \left[ \frac{1}{NT} Z'_{-j} Z_{-j} \right]_{lk} - [\Psi_{Z, -j, -j}]_{lk} \right| \leq \frac{\kappa^2(\Psi_{Z, -j, -j}, s_j^*)}{32s_j^*} \right\} \subseteq \mathcal{E}_j$$

by Proposition 1 in Appendix F with  $A = \Psi_{Z, -j, -j}$ ,  $B = \frac{1}{NT} Z'_{-j} Z_{-j}$ ,  $r = s_j^*$  and  $\delta = \frac{\kappa^2(\Psi_{Z, -j, -j}, s_j^*)}{32s_j^*}$ . Observe that the relation

$$\begin{aligned} \max_{1 \leq l, k \leq p-1} \left| \left[ \frac{1}{NT} Z'_{-j} Z_{-j} \right]_{lk} - [\Psi_{Z, -j, -j}]_{lk} \right| &\leq \max_{1 \leq l, k \leq p} \left| \left[ \frac{1}{NT} Z' Z \right]_{lk} - [\Psi_Z]_{lk} \right| \\ &\leq \frac{\kappa^2(\Psi_Z, \max_{j \in H_1} s_j^*)}{32\bar{G} \lambda_{node}^{-\vartheta}} \leq \frac{\kappa^2(\Psi_{Z, -j, -j}, s_j^*)}{32s_j^*}, \end{aligned}$$

implies  $\mathcal{E} := \left\{ \max_{1 \leq l, k \leq p} \left| \left[ \frac{1}{NT} Z' Z \right]_{lk} - [\Psi_Z]_{lk} \right| \leq \frac{\kappa^2(\Psi_Z, \max_{j \in H_1} s_j^*)}{32 \bar{G} \lambda_{node}^{-\vartheta}} \right\} \subseteq \tilde{\mathcal{E}}_j \subseteq \mathcal{E}_j$  for all  $j \in H_1$  and hence  $\mathcal{E} \subseteq \cap_{j \in H_1} \mathcal{E}_j$ . It remains to provide a lower bound on  $\mathbb{P}(\mathcal{E})$ . A typical element of  $\frac{1}{NT} Z' Z - \Psi_Z$  is of the form  $\frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T (z_{i,t,l} z_{i,t,k} - \mathbb{E}[z_{i,t,l} z_{i,t,k}])$  for some  $l, k \in \{1, \dots, p\}$ . Invoking (8.14) with  $x = \frac{\kappa^2(\Psi_Z, \max_{j \in H_1} s_j^*)}{32 \bar{G} \lambda_{node}^{-\vartheta}} \gtrsim \frac{1}{\sqrt{N}}$  (using  $\bar{G} \lambda_{node}^{1-\vartheta} = O(\log^{3/2} p)$ , implied by Assumption 4(b))

$$\mathbb{P} \left( \left| \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T (z_{i,t,l} z_{i,t,k} - \mathbb{E}[z_{i,t,l} z_{i,t,k}]) \right| \geq x \right) \leq A e^{-B(x^2 N)^{1/3}},$$

for positive constants  $A$  and  $B$ . Therefore,

$$\mathbb{P}(\mathcal{E}^c) = \mathbb{P} \left( \max_{1 \leq l, k \leq p} \left| \left[ \frac{1}{NT} Z' Z \right]_{lk} - [\Psi_Z]_{lk} \right| \geq x \right) \leq p^2 A e^{-B(x^2 N)^{1/3}}.$$

The upper bound of the preceding probability becomes arbitrarily small for  $M$  sufficiently large (using  $\bar{G} \lambda_{node}^{1-\vartheta} = O(1)$ , implied by Assumption 4(b)). In a similar manner, invoke (8.14) with  $x = \lambda_{node} = \sqrt{\frac{16M(\log p)^3}{N}} \gtrsim \frac{1}{\sqrt{N}}$  ( $M > 0$ ),

$$\mathbb{P}(\mathcal{F}^c) = \mathbb{P} \left( \max_{1 \leq l, k \leq p} \left| \left[ \frac{1}{NT} Z' Z \right]_{lk} - [\Psi_Z]_{lk} \right| \geq x \right) \leq A p^2 e^{-B(x^2 N)^{1/3}} = A p^{2-BM^{1/3}},$$

for positive constants  $A$  and  $B$ , letting  $B$  absorb the extra constants. The upper bound of the preceding probability becomes arbitrarily small for sufficiently large  $N$  and  $M$ . We also have

$$\left\| \frac{Z' Z}{NT} - \Psi_Z \right\|_{\infty} = O_p(\lambda_{node}) = O_p \left( \sqrt{\frac{(\log p)^3}{N}} \right). \quad (9.8)$$

Lastly, use Assumption 4(b) in the display (9.7) to get the claimed orders.  $\square$

## 9.2 Proof of Lemma 9

Lemma 9 below is used as a stepping stone towards the establishing asymptotically gaussian inference as it provides the rate at which  $\hat{\Theta}_Z$  approaches  $\Theta_Z$  uniformly over  $H_1$

**Lemma 9.** *Let Assumptions 1, 3 and 4 hold. Define  $\lambda_{node} = \sqrt{16M(\log p)^3/N}$  for some*

$M > 0$ . Then, for  $M$  sufficiently large,

$$\max_{j \in H_1} |\hat{\tau}_j^2 - \tau_j^2| = O_p \left( \bar{G}^{1/2} \left[ \frac{(\log p)^3}{N} \right]^{\frac{2-\vartheta}{4}} \right) \quad (9.9)$$

$$\max_{j \in H_1} \frac{1}{\hat{\tau}_j^2} = O_p(1) \quad (9.10)$$

$$\max_{j \in H_1} \left| \frac{1}{\hat{\tau}_j^2} - \frac{1}{\tau_j^2} \right| = O_p \left( \bar{G}^{1/2} \left[ \frac{(\log p)^3}{N} \right]^{\frac{2-\vartheta}{4}} \right) \quad (9.11)$$

$$\max_{j \in H_1} \left\| \hat{\Theta}_{Z,j} - \Theta_{Z,j} \right\|_1 = O_p \left( \bar{G} \left[ \frac{(\log p)^3}{N} \right]^{\frac{1-\vartheta}{2}} \right) \quad (9.12)$$

$$\max_{j \in H_1} \left\| \hat{\Theta}_{Z,j} - \Theta_{Z,j} \right\| = O_p \left( \bar{G}^{1/2} \left[ \frac{(\log p)^3}{N} \right]^{\frac{2-\vartheta}{4}} \right) \quad (9.13)$$

$$\max_{j \in H_1} \left\| \hat{\Theta}_{Z,j} \right\|_1 = O_p \left( \bar{G}^{1/2} \left[ \frac{(\log p)^3}{N} \right]^{-\frac{\vartheta}{4}} \right) \quad (9.14)$$

Note that for  $H_1 = \{1, \dots, p\}$ , (9.12) provides an upper bound on the induced  $\ell_\infty$ -distance between  $\hat{\Theta}_Z$  and  $\Theta_Z$ . However, we only need to control this distance for those indices corresponding to the parameters whose joint limit distribution is sought. On the other hand, it should be stressed that the uniformity over  $H_1$  of the above results is crucial in establishing the limiting gaussian inference and providing a feasible estimator of the covariance matrix of the parameter estimates. In case one is only interested in one entry of  $\gamma$ ,  $H_1$  reduces to a singleton if this entry is in  $\alpha$ . If this entry is in  $\eta$ , Lemma 9 is actually superfluous as the lower right hand corners of  $\hat{\Theta}$  and  $\Theta$  are identical.

*Proof of Lemma 9.* Recall (13.3) and use  $z_j = Z_{-j}\phi_j + \zeta_j$ :

$$\hat{\tau}_j^2 = \frac{1}{NT} \zeta_j' \zeta_j + \frac{1}{NT} \zeta_j' Z_{-j} \phi_j - \frac{1}{NT} (\hat{\phi}_j - \phi_j)' Z_{-j}' \zeta_j - \frac{1}{NT} (\hat{\phi}_j - \phi_j)' Z_{-j}' Z_{-j} \phi_j.$$

Thus,

$$\begin{aligned} \max_{j \in H_1} |\hat{\tau}_j^2 - \tau_j^2| &\leq \max_{j \in H_1} \left| \frac{1}{NT} \zeta_j' \zeta_j - \tau_j^2 \right| + \max_{j \in H_1} \left| \frac{1}{NT} \zeta_j' Z_{-j} \phi_j \right| \\ &\quad + \max_{j \in H_1} \left| \frac{1}{NT} (\hat{\phi}_j - \phi_j)' Z_{-j}' \zeta_j \right| + \max_{j \in H_1} \left| \frac{1}{NT} (\hat{\phi}_j - \phi_j)' Z_{-j}' Z_{-j} \phi_j \right|. \end{aligned} \quad (9.15)$$

Consider the first term on the right of the inequality in (9.15). By Assumption 4(c), we have for all  $\epsilon > 0$ ,  $\mathbb{P}(|\zeta_{j,i,t}^2| \geq \epsilon) = \mathbb{P}(|\zeta_{j,i,t}| \geq \sqrt{\epsilon}) \leq \frac{1}{2} K e^{-C\epsilon}$ . It follows from Lemma 2.2.1 in van der Vaart and Wellner (1996) that  $\|\zeta_{j,i,t}^2\|_{\psi_1} \leq (1 + K/2)/C$ . Therefore, by Jensen's inequality and subadditivity of the Orlicz norm

$$\left\| \frac{1}{T} \sum_{t=1}^T \left( \zeta_{j,i,t}^2 - \mathbb{E}[\zeta_{j,i,t}^2] \right) \right\|_{\psi_1} \leq 2 \max_{1 \leq t \leq T} \|\zeta_{j,i,t}^2\|_{\psi_1} \leq \frac{2+K}{C}.$$

Using the definition of the Orlicz norm,  $\mathbb{E} \exp \left( \frac{C}{2+K} \left| \frac{1}{T} \sum_{t=1}^T (\zeta_{j,i,t}^2 - \mathbb{E}[\zeta_{j,i,t}^2]) \right| \right) \leq 2$ . Using independence across  $i = 1, \dots, N$  (Assumption 1) to invoke Proposition 2 in Appendix F with  $D = C/(2+K)$ ,  $\alpha = 1/3$  and  $C_1 = 2$  for  $x \gtrsim \frac{1}{\sqrt{N}}$ ,

$$\mathbb{P} \left( \left| \frac{1}{N} \sum_{i=1}^N \frac{1}{T} \sum_{t=1}^T (\zeta_{j,i,t}^2 - \mathbb{E}[\zeta_{j,i,t}^2]) \right| \geq x \right) \leq A e^{-B(x^2 N)^{1/3}},$$

for positive constants  $A$  and  $B$ . Setting  $x = \sqrt{\frac{M(\log h_1)^3}{N}}$  for some  $M > 0$ , we have

$$\begin{aligned} & \mathbb{P} \left( \max_{j \in H_1} \left| \frac{1}{N} \sum_{i=1}^N \frac{1}{T} \sum_{t=1}^T (\zeta_{j,i,t}^2 - \mathbb{E}[\zeta_{j,i,t}^2]) \right| \geq \sqrt{\frac{M(\log h_1)^3}{N}} \right) \\ & \leq \sum_{j \in H_1} \mathbb{P} \left( \left| \frac{1}{N} \sum_{i=1}^N \frac{1}{T} \sum_{t=1}^T (\zeta_{j,i,t}^2 - \mathbb{E}[\zeta_{j,i,t}^2]) \right| \geq \sqrt{\frac{M(\log h_1)^3}{N}} \right) \leq A h_1^{1-BM^{1/3}}. \end{aligned}$$

Recognising that the upper bound of the preceding probability becomes arbitrarily small for sufficiently large  $N$  and  $M$ , we have

$$\max_{j \in H_1} \left| \frac{1}{NT} \zeta_j' \zeta_j - \tau_j^2 \right| = O_p \left( \sqrt{\frac{(\log h_1)^3}{N}} \right) = O_p(\lambda_{node}).$$

Now consider the second term on the right of the inequality in (9.15). Recall that

$$C = \begin{pmatrix} 1 & -\phi_{1,2} & \cdots & -\phi_{1,p} \\ -\phi_{2,1} & 1 & \cdots & -\phi_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ -\phi_{p,1} & -\phi_{p,2} & \cdots & 1 \end{pmatrix}$$

such that  $C_j$  is the  $j$ th row of  $C$  but written as a  $p \times 1$  vector. Then

$$\begin{aligned} & \max_{j \in H_1} \|\phi_j\|_1 = \max_{j \in H_1} \|\phi_j^* - \phi_j - \phi_j^*\|_1 \leq \max_{j \in H_1} \|\phi_j^* - \phi_j\|_1 + \max_{j \in H_1} \|\phi_j^*\|_1 \leq \bar{G} \lambda_{node}^{1-\vartheta} + \max_{j \in H_1} \|\phi_j^*\|_1 \\ & \leq \bar{G} \lambda_{node}^{1-\vartheta} + \max_{j \in H_1} \sqrt{s_j^*} \|\phi_j^*\| \leq \bar{G} \lambda_{node}^{1-\vartheta} + \max_{j \in H_1} \sqrt{s_j^*} \|\phi_j\| \leq \bar{G} \lambda_{node}^{1-\vartheta} + \max_{j \in H_1} \sqrt{s_j^*} \|C_j\| \\ & \leq \bar{G} \lambda_{node}^{1-\vartheta} + \max_{j \in H_1} \sqrt{s_j^*} \sqrt{\frac{C_j' \Psi_Z C_j}{\text{mineval}(\Psi_Z)}} = \bar{G} \lambda_{node}^{1-\vartheta} + \max_{j \in H_1} \frac{\sqrt{s_j^*} \sqrt{\Psi_{Z,j,j} - \Psi_{Z,j,-j} \Psi_{Z,-j,-j}^{-1} \Psi_{Z,-j,j}}}{\sqrt{\text{mineval}(\Psi_Z)}} \\ & \leq \bar{G} \lambda_{node}^{1-\vartheta} + \max_{j \in H_1} \frac{\sqrt{s_j^*} \sqrt{\Psi_{Z,j,j}}}{\sqrt{\text{mineval}(\Psi_Z)}} \leq \bar{G} \lambda_{node}^{1-\vartheta} + \max_{j \in H_1} \frac{\sqrt{s_j^*} \sqrt{\text{maxeval}(\Psi_Z)}}{\sqrt{\text{mineval}(\Psi_Z)}} = O(\bar{G}^{1/2} \lambda_{node}^{-\vartheta/2}) \end{aligned} \tag{9.16}$$

where the second inequality is due to (9.5), the second equality is due to (3.4), the seventh inequality is due to that Assumption 4(a) implies that  $\Psi_{Z,-j,-j}^{-1}$  is positive definite for all  $j \in H_1$ , and the last equality is due to (9.6) and Assumption 4(b)). Now,

$$\max_{j \in H_1} \left| \frac{1}{NT} \zeta_j' Z_{-j} \phi_j \right| \leq \max_{j \in H_1} \left( \left\| \frac{1}{NT} \zeta_j' Z_{-j} \right\|_{\infty} \|\phi_j\|_1 \right) = O_p(\lambda_{node}) O(\bar{G}^{1/2} \lambda_{node}^{-\vartheta/2}) = O_p(\bar{G}^{1/2} \lambda_{node}^{1-\vartheta/2}),$$

where the first equality is due to (9.3).

The third term in (9.15) is bounded as

$$\max_{j \in H_1} \left| \frac{1}{NT} (\hat{\phi}_j - \phi_j)' Z'_{-j} \zeta_j \right| \leq \max_{j \in H_1} \left( \left\| \hat{\phi}_j - \phi_j \right\|_1 \left\| \frac{1}{NT} Z'_{-j} \zeta_j \right\|_\infty \right) = O_p(\bar{G} \lambda_{node}^{2-\vartheta}),$$

where the equality is due to (9.2) and (9.3).

To bound the fourth term on the right of the inequality in (9.15), recall (13.5) and manipulate to get  $\frac{1}{NT} Z'_{-j} Z_{-j} (\hat{\phi}_j - \phi_j) = \frac{1}{NT} Z'_{-j} \zeta_j - \lambda_{node} w_j$ . Thus,

$$\left\| \frac{1}{NT} (\hat{\phi}_j - \phi_j)' Z'_{-j} Z_{-j} \right\|_\infty \leq \left\| \frac{1}{NT} Z'_{-j} \zeta_j \right\|_\infty + \lambda_{node} \|w_j\|_\infty = O_p(\lambda_{node}),$$

where the equality is due to (9.3). Thus,

$$\max_{j \in H_1} \left| \frac{1}{NT} (\hat{\phi}_j - \phi_j)' Z'_{-j} Z_{-j} \phi_j \right| \leq \max_{j \in H_1} \left\| \frac{1}{NT} (\hat{\phi}_j - \phi_j)' Z'_{-j} Z_{-j} \right\|_\infty \max_{j \in H_1} \|\phi_j\|_1 = O_p(\bar{G}^{1/2} \lambda_{node}^{1-\vartheta/2}),$$

where the last equality is due to (9.16). Summing up all four terms on the right of the inequality in (9.15), we get

$$\max_{j \in H_1} |\hat{\tau}_j^2 - \tau_j^2| \leq O_p(\lambda_{node}) + O_p(\bar{G}^{1/2} \lambda_{node}^{1-\vartheta/2}) + O_p(\bar{G} \lambda_{node}^{2-\vartheta}) = O_p(\bar{G}^{1/2} \lambda_{node}^{1-\vartheta/2}) = o_p(1),$$

where the first equality is due to that  $O_p(\bar{G}^{1/2} \lambda_{node}^{1-\vartheta/2})$  dominates  $O_p(\bar{G} \lambda_{node}^{2-\vartheta})$  by Assumption 4(b), and the second equality is also due to Assumption 4(b). This establishes (9.9).

We now prove (9.10). We first recall

$$\tau_j^2 = \mathbb{E} \left[ \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T (z_{i,t,j} - z'_{i,t,-j} \phi_j)^2 \right] = \Psi_{Z,j,j} - \Psi_{Z,j,-j} \Psi_{Z,-j,-j}^{-1} \Psi_{Z,-j,j} = \frac{1}{\Theta_{Z,j,j}},$$

Furthermore,

$$\Theta_{Z,j,j} \equiv \frac{e_j' \Theta_Z e_j}{\|e_j\|^2} \leq \max_{\delta \in \mathbb{R}^p \setminus \{0\}} \frac{\delta' \Theta_Z \delta}{\|\delta\|^2} = \text{maxeval}(\Theta_Z) = \frac{1}{\text{mineval}(\Psi_Z)}.$$

The preceding inequality is uniform in  $j$ . Thus,  $\min_{j \in H_1} \tau_j^2 \geq \text{mineval}(\Psi_Z)$ , which is uniformly bounded away from zero by Assumption 4(a). Therefore,

$$\min_{j \in H_1} \hat{\tau}_j^2 = \min_{j \in H_1} (\hat{\tau}_j^2 - \tau_j^2 + \tau_j^2) \geq \min_{j \in H_1} \tau_j^2 - \max_{j \in H_1} |\hat{\tau}_j^2 - \tau_j^2| \geq \text{mineval}(\Psi_Z) - o_p(1).$$

Hence, we conclude that  $\min_{j \in H_1} \hat{\tau}_j^2$  is bounded away from zero for  $N$  large enough and  $\max_{j \in H_1} \frac{1}{\hat{\tau}_j^2} = O_p(1)$  which establishes (9.10).

Hence,

$$\max_{j \in H_1} \left| \frac{1}{\hat{\tau}_j^2} - \frac{1}{\tau_j^2} \right| \leq \frac{\max_{j \in H_1} |\tau_j^2 - \hat{\tau}_j^2|}{\min_{j \in H_1} \tau_j^2} \cdot \max_{j \in H_1} \frac{1}{\tau_j^2} = \max_{j \in H_1} |\tau_j^2 - \hat{\tau}_j^2| O(1) O_p(1) = O_p(\bar{G}^{1/2} \lambda_{node}^{1-\vartheta/2}),$$

which establishes (9.11).



We can now bound  $\max_{j \in H_1} \|\hat{\Theta}_{Z,j} - \Theta_{Z,j}\|_1$ . Use the definition of  $C_j$  and (3.3) to recognise that  $\Theta_{Z,j} = C_j \Theta_{Z,j,j} = C_j / \tau_j^2$ .

$$\begin{aligned} \max_{j \in H_1} \|\hat{\Theta}_{Z,j} - \Theta_{Z,j}\|_1 &= \max_{j \in H_1} \left\| \frac{\hat{C}_j}{\hat{\tau}_j^2} - \frac{C_j}{\tau_j^2} \right\|_1 = \max_{j \in H_1} \left| \frac{1}{\hat{\tau}_j^2} - \frac{1}{\tau_j^2} \right| + \max_{j \in H_1} \left\| \frac{\hat{\phi}_j}{\hat{\tau}_j^2} - \frac{\phi_j}{\tau_j^2} \right\|_1 \\ &\leq \max_{j \in H_1} \left| \frac{1}{\hat{\tau}_j^2} - \frac{1}{\tau_j^2} \right| + \max_{j \in H_1} \left\| \frac{\hat{\phi}_j}{\hat{\tau}_j^2} - \frac{\phi_j}{\tau_j^2} \right\|_1 + \max_{j \in H_1} \left\| \frac{\phi_j}{\hat{\tau}_j^2} - \frac{\phi_j}{\tau_j^2} \right\|_1 \\ &= \max_{j \in H_1} \left| \frac{1}{\hat{\tau}_j^2} - \frac{1}{\tau_j^2} \right| + \max_{j \in H_1} \frac{1}{\hat{\tau}_j^2} \|\hat{\phi}_j - \phi_j\|_1 + \max_{j \in H_1} \|\phi_j\|_1 \left| \frac{1}{\hat{\tau}_j^2} - \frac{1}{\tau_j^2} \right| = O_p(\bar{G} \lambda_{node}^{1-\vartheta}), \end{aligned}$$

which establishes (9.12). Next, we bound  $\max_{j \in H_1} \|\hat{\Theta}_{Z,j} - \Theta_{Z,j}\|$ . Since

$$\begin{aligned} \left| (\hat{C}_j - C_j)' \frac{Z'Z}{NT} (\hat{C}_j - C_j) - (\hat{C}_j - C_j)' \Psi_Z (\hat{C}_j - C_j) \right| &\leq \left\| \frac{Z'Z}{NT} - \Psi_Z \right\|_\infty \|\hat{C}_j - C_j\|_1^2, \\ \max_{j \in H_1} \left| (\hat{C}_j - C_j)' \Psi_Z (\hat{C}_j - C_j) \right| &\leq \max_{j \in H_1} \left| (\hat{C}_j - C_j)' \frac{Z'Z}{NT} (\hat{C}_j - C_j) \right| + \left\| \frac{Z'Z}{NT} - \Psi_Z \right\|_\infty \max_{j \in H_1} \|\hat{C}_j - C_j\|_1^2. \end{aligned} \quad (9.17)$$

Consider the first term on the right hand side of (9.17).

$$\max_{j \in H_1} \left| (\hat{C}_j - C_j)' \frac{Z'Z}{NT} (\hat{C}_j - C_j) \right| = \max_{j \in H_1} \frac{1}{NT} \left\| Z(\hat{C}_j - C_j) \right\|^2 = \max_{j \in H_1} \frac{1}{NT} \left\| Z_{-j}(\hat{\phi}_j - \phi_j) \right\|^2 = O_p(\bar{G} \lambda_{node}^{2-\vartheta}),$$

where the last equality is due to (9.1). Next, consider the second term on the right of the inequality (9.17). We have

$$\left\| \frac{Z'Z}{NT} - \Psi_Z \right\|_\infty \max_{j \in H_1} \|\hat{C}_j - C_j\|_1^2 = \left\| \frac{Z'Z}{NT} - \Psi_Z \right\|_\infty \max_{j \in H_1} \|\hat{\phi}_j - \phi_j\|_1^2 = O_p(\bar{G}^2 \lambda_{node}^{3-2\vartheta}),$$

where the first equality is due to the definitions of  $\hat{C}_j$  and  $C_j$ , and the second equality is due to (9.2) and (9.8). Adding up the two terms, we have

$$\max_{j \in H_1} \left| (\hat{C}_j - C_j)' \Psi_Z (\hat{C}_j - C_j) \right| \leq O_p(\bar{G} \lambda_{node}^{2-\vartheta}) + O_p(\bar{G}^2 \lambda_{node}^{3-2\vartheta}) = O_p(\bar{G} \lambda_{node}^{2-\vartheta}),$$

where the last equality is due to Assumption 4(b). Since  $\max_{j \in H_1} |(\hat{C}_j - C_j)' \Psi_Z (\hat{C}_j - C_j)| \geq \text{mineval}(\Psi_Z) \max_{j \in H_1} \|\hat{C}_j - C_j\|^2$  and  $\text{mineval}(\Psi_Z)$  is uniformly bounded away from zero we have  $\max_{j \in H_1} \|\hat{\phi}_j - \phi_j\| = \max_{j \in H_1} \|\hat{C}_j - C_j\| = O_p(\bar{G}^{1/2} \lambda_{node}^{1-\vartheta/2})$ . Then,

$$\begin{aligned} \max_{j \in H_1} \|\hat{\Theta}_{Z,j} - \Theta_{Z,j}\| &= \max_{j \in H_1} \left\| \frac{\hat{C}_j}{\hat{\tau}_j^2} - \frac{C_j}{\tau_j^2} \right\| \leq \max_{j \in H_1} \left| \frac{1}{\hat{\tau}_j^2} - \frac{1}{\tau_j^2} \right| + \max_{j \in H_1} \left\| \frac{\hat{\phi}_j}{\hat{\tau}_j^2} - \frac{\phi_j}{\tau_j^2} \right\| \\ &\leq \max_{j \in H_1} \left| \frac{1}{\hat{\tau}_j^2} - \frac{1}{\tau_j^2} \right| + \max_{j \in H_1} \frac{1}{\hat{\tau}_j^2} \|\hat{\phi}_j - \phi_j\| + \max_{j \in H_1} \|\phi_j\| \left| \frac{1}{\hat{\tau}_j^2} - \frac{1}{\tau_j^2} \right| = O_p(\bar{G}^{1/2} \lambda_{node}^{1-\vartheta/2}), \end{aligned}$$

where in the last equality we have used that  $\max_{j \in H_1} \|\phi_j\| = O(1)$ , which follows from inspecting the arguments in (9.16). We have hence established (9.13). Finally, recall that  $\Theta_{Z,j} = C_j \Theta_{Z,j,j} = C_j / \tau_j^2$ . Thus,

$$\max_{j \in H_1} \|\Theta_{Z,j}\|_1 = \max_{j \in H_1} |1/\tau_j^2| + \max_{j \in H_1} \|\phi_j\|_1 \max_{j \in H_1} 1/\tau_j^2 = O(\bar{G}^{1/2} \lambda_{node}^{-\vartheta/2}). \quad (9.18)$$

Therefore,

$$\max_{j \in H_1} \|\hat{\Theta}_{Z,j}\|_1 \leq \max_{j \in H_1} \|\hat{\Theta}_{Z,j} - \Theta_{Z,j}\|_1 + \max_{j \in H_1} \|\Theta_{Z,j}\|_1 = O_p(\bar{G}\lambda_{node}^{1-\vartheta}) + O(\bar{G}^{1/2}\lambda_{node}^{-\vartheta/2}) = O_p(\bar{G}^{1/2}\lambda_{node}^{-\vartheta/2}),$$

where the last equality is due to Assumption 4(b).  $\square$

## 10 Appendix D

### 10.1 Proof of Theorem 2

*Proof of Theorem 2.* The following assumption is implied by Assumption 5.<sup>9</sup> However, as Assumption 5 is much simpler, we have chosen to use the latter in the main text even though it is slightly less general than the following assumption. Note again how the assumptions simplify when either  $h_1$  or  $h_2$  equals 0.

**Assumption 6.**

$$\begin{aligned} (a) \quad (i) \quad & \frac{h_1^2 \bar{G}^2 \left[ \frac{(\log p)^3}{N} \right]^{-\vartheta} (\log(p \vee T))^5}{N} = o(1); \\ (ii) \quad & \frac{h_1 h_2 \bar{G} \left[ \frac{(\log p)^3}{N} \right]^{-\vartheta/2} (\log(p \vee N \vee T))^3}{N} = o(1); \\ (iii) \quad & \frac{h_1 \bar{G}^2 \left[ \frac{(\log p)^3}{N} \right]^{-\vartheta} (\log p)^3 (\log(p \vee N))^3}{N} = o(1); \\ (iv) \quad & \frac{h_1 \bar{G} \left[ \frac{(\log p)^3}{N} \right]^{-\vartheta/2} (\log(N \vee T))^2 (\log p)^2}{NT} = o(1); \\ (v) \quad & \frac{(\log(N \vee T))^2 1_{\{h_2 \neq 0\}}}{T} = o(1). \end{aligned}$$

(b) Let

$$a := \frac{\left[ E \left( \frac{(\log(p \vee N))^3}{T} \right)^{-\nu/2} \right] (\log(p \vee N))^3}{NT}.$$

$$\begin{aligned} (i) \quad & h_1^2 \bar{G}^2 \left[ \frac{(\log p)^3}{N} \right]^{-\vartheta} \left( 1 \vee \sqrt{\frac{(\log(p \vee T))^7}{N}} \right) a = o(1); \\ (ii) \quad & h_1 \bar{G} \left[ \frac{(\log p)^3}{N} \right]^{-\vartheta/2} \log(p \vee N \vee T) a = o(1); \\ (iii) \quad & h_1 h_2 \bar{G} \left[ \frac{(\log p)^3}{N} \right]^{-\vartheta/2} (\log(p \vee N \vee T))^2 a = o(1); \\ (iv) \quad & \sqrt{h_1 h_2 \bar{G} \left[ \frac{(\log p)^3}{N} \right]^{-\vartheta/2} N \log(p \vee N \vee T)} a = o(1); \end{aligned}$$

---

<sup>9</sup>To be precise, Assumption 5(a) implies Assumption 6(a) by recognising that  $h_1 \geq 1$  if  $h_1 \neq 0$ , and  $\bar{G} \left[ \frac{(\log p)^3}{N} \right]^{-\vartheta/2} \geq s_j^* \geq 1$ . Assumption 5(b) implies Assumption 6(b) by recognising that  $\sqrt{\frac{(\log(p \vee T))^7}{N}} = o(1)$  and  $\sqrt{\frac{(\log(N \vee T))^3}{T}} = o(1)$ , implied by Assumption 5(a) provided  $h_1 \neq 0$  and  $h_2 \neq 0$ , respectively. Last, Assumption 5(c) implies Assumption 6(c).

$$(v) \quad N h_2^2 \left( 1 \vee \sqrt{\frac{(\log N)^3}{T}} \right) a = o(1).$$

(c)

$$\frac{(h_1 \vee h_2)(\log(p \vee N))^3 \left[ E^2 \left( \frac{(\log(p \vee N))^3}{T} \right)^{-\nu} \right] b}{N} = o(1),$$

$$\text{where } b := \left[ \left( \bar{G} \left[ \frac{(\log p)^3}{N} \right]^{-\vartheta/2} \log(p \vee N) \vee (\log p)^3 \right) 1\{h_1 \neq 0\} \right] \vee \left[ \log(p \vee N) 1\{h_2 \neq 0\} \right].$$

(d) *mineval*( $\Sigma_{\Pi_\varepsilon}$ ) is uniformly bounded away from zero and *maxeval*( $\Sigma_{1,N}$ ) is uniformly bounded from above.

We show that

$$t = \frac{\rho' S (\tilde{\gamma} - \gamma)}{\sqrt{\rho' \hat{\Theta} \hat{\Sigma}_{\Pi_\varepsilon} \hat{\Theta}' \rho}} \xrightarrow{d} N(0, 1).$$

To this end, note that by (3.2) one may write  $t = t_1 + t_2$ , where

$$t_1 = \frac{\rho' \hat{\Theta} S^{-1} \Pi' \varepsilon}{\sqrt{\rho' \hat{\Theta} \hat{\Sigma}_{\Pi_\varepsilon} \hat{\Theta}' \rho}} \quad \text{and} \quad t_2 = \frac{-\rho' \Delta}{\sqrt{\rho' \hat{\Theta} \hat{\Sigma}_{\Pi_\varepsilon} \hat{\Theta}' \rho}}.$$

Defining

$$t'_1 = \frac{\rho' \Theta S^{-1} \Pi' \varepsilon}{\sqrt{\rho' \Theta \Sigma_{\Pi_\varepsilon} \Theta' \rho}}$$

it suffices to show that  $t'_1 \xrightarrow{d} N(0, 1)$ ,  $t'_1 - t_1 = o_p(1)$ , and  $t_2 = o_p(1)$ . In the sequel we first show that  $t_1 - t'_1 = o_p(1)$ , then  $t'_1 \xrightarrow{d} N(0, 1)$  and finally  $t_2 = o_p(1)$ . To show that  $t_1 - t'_1 = o_p(1)$ , it suffices to show that the denominators as well as the numerators of  $t_1$  and  $t'_1$  are asymptotically equivalent since

$$\rho' \Theta \Sigma_{\Pi_\varepsilon} \Theta' \rho \geq \text{mineval}(\Sigma_{\Pi_\varepsilon}) (\text{mineval}(\Theta))^2 = \frac{\text{mineval}(\Sigma_{\Pi_\varepsilon})}{(\text{maxeval}(\Psi))^2} \quad (10.1)$$

which is uniformly bounded away from zero by Assumptions 4(a) and 6(d).

### 10.1.1 Denominators of $t_1$ and $t'_1$

We first show that the denominators of  $t_1$  and  $t'_1$  are asymptotically equivalent, i.e.,

$$|\rho' \hat{\Theta} \hat{\Sigma}_{\Pi_\varepsilon} \hat{\Theta}' \rho - \rho' \Theta \Sigma_{\Pi_\varepsilon} \Theta' \rho| = o_p(1). \quad (10.2)$$

Write

$$\left| (\rho'_1, \rho'_2) \begin{pmatrix} \hat{\Theta}_Z \hat{\Sigma}_{1,N} \hat{\Theta}'_Z & \hat{\Theta}_Z \hat{\Sigma}_{2,N} \\ \hat{\Sigma}'_{2,N} \hat{\Theta}'_Z & \hat{\Sigma}_{3,N} \end{pmatrix} \begin{pmatrix} \rho_1 \\ \rho_2 \end{pmatrix} - (\rho'_1, \rho'_2) \begin{pmatrix} \Theta_Z \Sigma_{1,N} \Theta'_Z & \Theta_Z \Sigma_{2,N} \\ \Sigma'_{2,N} \Theta'_Z & \Sigma_{3,N} \end{pmatrix} \begin{pmatrix} \rho_1 \\ \rho_2 \end{pmatrix} \right|$$

$$\leq |\rho'_1 \hat{\Theta}_Z \hat{\Sigma}_{1,N} \hat{\Theta}'_Z \rho_1 - \rho'_1 \Theta_Z \Sigma_{1,N} \Theta'_Z \rho_1| \quad (10.3)$$

$$+ 2|\rho'_1 \hat{\Theta}_Z \hat{\Sigma}_{2,N} \rho_2 - \rho'_1 \Theta_Z \Sigma_{2,N} \rho_2| \quad (10.4)$$

$$+ |\rho'_2 \hat{\Sigma}_{3,N} \rho_2 - \rho'_2 \Sigma_{3,N} \rho_2|. \quad (10.5)$$

To establish (10.2), we show that (10.3), (10.4) and (10.5) are  $o_p(1)$ , respectively.

**(10.3) is  $o_p(1)$ :**

Define  $\tilde{\Sigma}_{1,N} := \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \varepsilon_{i,t}^2 z_{i,t} z'_{i,t}$ . To show that (10.3) is  $o_p(1)$ , it suffices to show that

$$|\rho'_1 \hat{\Theta}_Z \hat{\Sigma}_{1,N} \hat{\Theta}'_Z \rho_1 - \rho'_1 \hat{\Theta}_Z \tilde{\Sigma}_{1,N} \hat{\Theta}'_Z \rho_1| = o_p(1) \quad (10.6)$$

$$|\rho'_1 \hat{\Theta}_Z \tilde{\Sigma}_{1,N} \hat{\Theta}'_Z \rho_1 - \rho'_1 \hat{\Theta}_Z \Sigma_{1,N} \hat{\Theta}'_Z \rho_1| = o_p(1) \quad (10.7)$$

$$|\rho'_1 \hat{\Theta}_Z \Sigma_{1,N} \hat{\Theta}'_Z \rho_1 - \rho'_1 \Theta_Z \Sigma_{1,N} \Theta'_Z \rho_1| = o_p(1). \quad (10.8)$$

We prove (10.6) first. Note that

$$|\rho'_1 \hat{\Theta}_Z \hat{\Sigma}_{1,N} \hat{\Theta}'_Z \rho_1 - \rho'_1 \hat{\Theta}_Z \tilde{\Sigma}_{1,N} \hat{\Theta}'_Z \rho_1| \leq \|\hat{\Sigma}_{1,N} - \tilde{\Sigma}_{1,N}\|_\infty \|\hat{\Theta}'_Z \rho_1\|_1^2.$$

First,

$$\|\hat{\Theta}'_Z \rho_1\|_1 = \left\| \sum_{j \in H_1} \hat{\Theta}_{Z,j} \rho_{1j} \right\|_1 \leq \sum_{j \in H_1} |\rho_{1j}| \|\hat{\Theta}_{Z,j}\|_1 = O_p(h_1^{1/2} \bar{G}^{1/2} \lambda_{node}^{-\vartheta/2}), \quad (10.9)$$

where the last equality is due to (9.14). We now bound  $\|\hat{\Sigma}_{1,N} - \tilde{\Sigma}_{1,N}\|_\infty$ . Since  $\hat{\varepsilon}_{i,t} = y_{i,t} - z'_{i,t} \hat{\alpha} - \hat{\eta}_i = \varepsilon_{i,t} - z'_{i,t}(\hat{\alpha} - \alpha) - (\hat{\eta}_i - \eta_i) =: \varepsilon_{i,t} - \pi_{i,t}(\hat{\gamma} - \gamma)$ , substituting for  $\hat{\varepsilon}_{i,t}$ , we have

$$\begin{aligned} \|\hat{\Sigma}_{1,N} - \tilde{\Sigma}_{1,N}\|_\infty &= \left\| \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \hat{\varepsilon}_{i,t}^2 z_{i,t} z'_{i,t} - \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \varepsilon_{i,t}^2 z_{i,t} z'_{i,t} \right\|_\infty \\ &\leq 2 \left\| \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T z_{i,t} z'_{i,t} \varepsilon_{i,t} \pi'_{i,t} (\hat{\gamma} - \gamma) \right\|_\infty + \left\| \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T z_{i,t} z'_{i,t} [\pi'_{i,t} (\hat{\gamma} - \gamma)]^2 \right\|_\infty. \end{aligned} \quad (10.10)$$

Consider the first term of (10.10). A typical element of  $\frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T z_{i,t} z'_{i,t} \varepsilon_{i,t} \pi'_{i,t} (\hat{\gamma} - \gamma)$  is

$$\begin{aligned} \frac{1}{NT} \sum_{j=1}^{NT} z_{j,l} z_{j,k} \varepsilon_j \pi'_j (\hat{\gamma} - \gamma) &\leq \frac{1}{NT} \left( \sum_{j=1}^{NT} z_{j,l}^2 z_{j,k}^2 \varepsilon_j^2 \right)^{1/2} \left( \sum_{j=1}^{NT} [\pi'_j (\hat{\gamma} - \gamma)]^2 \right)^{1/2} \\ &= \left( \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T z_{i,t,l}^2 z_{i,t,k}^2 \varepsilon_{i,t}^2 \right)^{1/2} \left( \frac{1}{NT} \|\Pi(\hat{\gamma} - \gamma)\|^2 \right)^{1/2} \end{aligned} \quad (10.11)$$

for some  $l, k \in \{1, \dots, p\}$ , where the inequality is due to Cauchy-Schwarz inequality. Use independence across  $i$  (Assumption 1) and subgaussianity (Assumption 3) to invoke Proposition 3 in Appendix F, such that

$$\max_{1 \leq l \leq p} \max_{1 \leq k \leq p} \left| \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T (z_{i,t,l}^2 z_{i,t,k}^2 \varepsilon_{i,t}^2 - \mathbb{E}[z_{i,t,l}^2 z_{i,t,k}^2 \varepsilon_{i,t}^2]) \right| = O_p \left( \sqrt{\frac{(\log(p^2 T))^7}{N}} \right)$$

and

$$\max_{1 \leq l \leq p} \max_{1 \leq k \leq p} \max_{1 \leq i \leq N} \max_{1 \leq t \leq T} \mathbb{E}[z_{i,t,l}^2 z_{i,t,k}^2 \varepsilon_{i,t}^2] \leq A = O(1)$$

for some positive constant  $A$ . Then, by the triangle inequality,

$$\max_{1 \leq l \leq p} \max_{1 \leq k \leq p} \left| \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T z_{i,t,l}^2 z_{i,t,k}^2 \varepsilon_{i,t}^2 \right| = O_p \left( \sqrt{\frac{(\log(p \vee T))^7}{N}} \right) + O(1). \quad (10.12)$$

Combining (10.11) and (10.12), we have

$$\left\| \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T z_{i,t} z'_{i,t} \varepsilon_{i,t} \pi'_{i,t} (\hat{\gamma} - \gamma) \right\|_{\infty} = O_p \left( \frac{(\log(p \vee T))^{7/4}}{N^{1/4}} \vee 1 \right) \left( \frac{1}{NT} \|\Pi(\hat{\gamma} - \gamma)\|^2 \right)^{1/2}. \quad (10.13)$$

We now consider the second term of (10.10). A typical element of  $\frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T z_{i,t} z'_{i,t} [\pi'_{i,t} (\hat{\gamma} - \gamma)]^2$  is  $\frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T z_{i,t,l} z_{i,t,k} [\pi'_{i,t} (\hat{\gamma} - \gamma)]^2 \leq \max_{1 \leq i \leq N} \max_{1 \leq t \leq T} |z_{i,t,l} z_{i,t,k}| \frac{1}{NT} \|\Pi(\hat{\gamma} - \gamma)\|^2$  for some  $l, k \in \{1, \dots, p\}$ . Recall that we have proved in the proof of Lemma 7 that  $\|z_{i,t,l} z_{i,t,k}\|_{\psi_1} \leq (1+K)/C$ . Using the definition of the Orlicz norm, we have  $\mathbb{E} e^{\frac{C}{1+K} |z_{i,t,l} z_{i,t,k}|} \leq 2$ . Using Markov's inequality, we have for any  $\epsilon > 0$

$$\mathbb{P} \left( \max_{1 \leq l \leq p} \max_{1 \leq k \leq p} \max_{1 \leq i \leq N} \max_{1 \leq t \leq T} |z_{i,t,l} z_{i,t,k}| \geq \epsilon \right) \leq \sum_{l=1}^p \sum_{k=1}^p \sum_{i=1}^N \sum_{t=1}^T \frac{\mathbb{E} e^{\frac{C}{1+K} |z_{i,t,l} z_{i,t,k}|}}{e^{\frac{C}{1+K} \epsilon}} \leq 2NTp^2 e^{-\frac{C}{1+K} \epsilon}.$$

Set  $\epsilon = M \log(p^2 NT)$  for some  $M > 0$  and note that the upper bound of the preceding probability becomes arbitrarily small for  $N$  and  $M$  sufficiently large. Thus,

$$\max_{1 \leq l \leq p} \max_{1 \leq k \leq p} \max_{1 \leq i \leq N} \max_{1 \leq t \leq T} |z_{i,t,l} z_{i,t,k}| = O_p(\log(p^2 NT))$$

and we get

$$\left\| \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T z_{i,t} z'_{i,t} [\pi'_{i,t} (\hat{\gamma} - \gamma)]^2 \right\|_{\infty} = O_p(\log(p \vee N \vee T)) \frac{1}{NT} \|\Pi(\hat{\gamma} - \gamma)\|^2. \quad (10.14)$$

Combining (10.13) and (10.14), conclude

$$\begin{aligned} & \left\| \hat{\Sigma}_{1,N} - \tilde{\Sigma}_{1,N} \right\|_{\infty} \\ &= O_p \left( \frac{(\log(p \vee T))^{7/4}}{N^{1/4}} \vee 1 \right) \left( \frac{1}{NT} \|\Pi(\hat{\gamma} - \gamma)\|^2 \right)^{1/2} + O_p(\log(p \vee N \vee T)) \frac{1}{NT} \|\Pi(\hat{\gamma} - \gamma)\|^2. \end{aligned}$$

Therefore, combining the preceding rates with (10.9) one gets

$$\begin{aligned}
& |\rho'_1 \hat{\Theta}_Z \hat{\Sigma}_{1,N} \hat{\Theta}'_Z \rho_1 - \rho'_1 \hat{\Theta}_Z \tilde{\Sigma}_{1,N} \hat{\Theta}'_Z \rho_1| \\
&= O_p(h_1 \bar{G} \lambda_{node}^{-\vartheta}) O_p \left( \frac{(\log(p \vee T))^{7/4}}{N^{1/4}} \vee 1 \right) \left[ \frac{1}{NT} \|\Pi(\hat{\gamma} - \gamma)\|^2 \right]^{1/2} \\
&\quad + O_p(h_1 \bar{G} \lambda_{node}^{-\vartheta}) O_p(\log(p \vee N \vee T)) \frac{1}{NT} \|\Pi(\hat{\gamma} - \gamma)\|^2 \\
&= o_p(1),
\end{aligned}$$

where the last equality is also due to Assumption 6(b)(i)-(ii), which establishes (10.6).

Next, turn to (10.7). Note that

$$|\rho'_1 \hat{\Theta}_Z \tilde{\Sigma}_{1,N} \hat{\Theta}'_Z \rho_1 - \rho'_1 \hat{\Theta}_Z \Sigma_{1,N} \hat{\Theta}'_Z \rho_1| \leq \|\tilde{\Sigma}_{1,N} - \Sigma_{1,N}\|_\infty \|\hat{\Theta}'_Z \rho_1\|_1^2.$$

Given (10.9), we only need to consider  $\|\tilde{\Sigma}_{1,N} - \Sigma_{1,N}\|_\infty$ . Using independence across  $i$  (Assumption 1) and subgaussianity (Assumption 3) to invoke Proposition 3 in Appendix F such that

$$\|\tilde{\Sigma}_{1,N} - \Sigma_{1,N}\|_\infty = \max_{1 \leq l \leq p} \max_{1 \leq k \leq p} \left| \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T (z_{i,t,l} z_{i,t,k} \varepsilon_{it}^2 - \mathbb{E}[z_{i,t,l} z_{i,t,k} \varepsilon_{it}^2]) \right| = O_p \left( \sqrt{\frac{(\log(p^2 T))^5}{N}} \right). \quad (10.15)$$

Thus,

$$|\rho'_1 \hat{\Theta}_Z \tilde{\Sigma}_{1,N} \hat{\Theta}'_Z \rho_1 - \rho'_1 \hat{\Theta}_Z \Sigma_{1,N} \hat{\Theta}'_Z \rho_1| = O_p \left( \sqrt{\frac{(\log(p \vee T))^5}{N}} h_1 \bar{G} \lambda_{node}^{-\vartheta} \right) = o_p(1),$$

where the last equality is due to Assumption 6(a)(i), establishing (10.7).

To prove (10.8) invoke Lemma 10 in Appendix F:

$$\begin{aligned}
& |\rho'_1 \hat{\Theta}_Z \Sigma_{1,N} \hat{\Theta}'_Z \rho_1 - \rho'_1 \Theta_Z \Sigma_{1,N} \Theta'_Z \rho_1| \leq \|\Sigma_{1,N}\|_\infty \|(\hat{\Theta}'_Z - \Theta'_Z) \rho_1\|_1^2 + 2\|\Sigma_{1,N} \Theta'_Z \rho_1\| \|(\hat{\Theta}'_Z - \Theta'_Z) \rho_1\| \\
&\leq \|\Sigma_{1,N}\|_\infty \|(\hat{\Theta}'_Z - \Theta'_Z) \rho_1\|_1^2 + 2\text{maxeval}(\Sigma_{1,N}) \|\Theta'_Z \rho_1\| \|(\hat{\Theta}'_Z - \Theta'_Z) \rho_1\|.
\end{aligned}$$

First, note that  $\|\Sigma_{1,N}\|_\infty$  is uniformly bounded as every entry is an average of uniformly bounded population moments (see Proposition 3 in Appendix F).

$$\begin{aligned}
& \|(\hat{\Theta}'_Z - \Theta'_Z) \rho_1\|_1 \leq \sum_{j \in H_1} \|\hat{\Theta}_{Z,j} - \Theta_{Z,j}\|_1 |\rho_{1j}| \leq \max_{j \in H_1} \|\hat{\Theta}_{Z,j} - \Theta_{Z,j}\|_1 \sqrt{h_1} \\
&= O_p \left( \bar{G} \left[ \frac{(\log p)^3}{N} \right]^{\frac{1-\vartheta}{2}} \sqrt{h_1} \right) = o_p(1),
\end{aligned} \quad (10.16)$$

where the first equality is due to (9.12), and the last equality is due to Assumption 6(a)(i). Next,  $\|\Theta'_Z \rho_1\| \leq \text{maxeval}(\Theta_Z) \|\rho_1\| \leq \text{maxeval}(\Theta_Z) = 1/\text{mineval}(\Psi_Z)$ , which is uniformly bounded

from above by Assumption 4(a). Furthermore,

$$\begin{aligned} \|(\hat{\Theta}'_Z - \Theta'_Z)\rho_1\| &= \left\| \sum_{j \in H_1} (\hat{\Theta}_{Z,j} - \Theta_{Z,j})\rho_{1j} \right\| \leq \sum_{j \in H_1} \|\hat{\Theta}_{Z,j} - \Theta_{Z,j}\| |\rho_{1j}| \\ &\leq \max_{j \in H_1} \|\hat{\Theta}_{Z,j} - \Theta_{Z,j}\| \sqrt{h_1} = O_p \left( \bar{G}^{1/2} \left[ \frac{(\log p)^3}{N} \right]^{\frac{2-\vartheta}{4}} \sqrt{h_1} \right) = o_p(1), \end{aligned}$$

where the second last equality is due to (9.13), and the last equality is due to (10.16). Thus, we have established (10.8) concluding the proof of (10.3) is  $o_p(1)$ .

**(10.4) is  $o_p(1)$ :**

Define  $\tilde{\Sigma}_{2,N} := \frac{1}{\sqrt{NT}} \sum_{i=1}^N \sum_{t=1}^T \varepsilon_{i,t}^2 z_{i,t} d'_{i,t}$ . It suffices to show

$$|\rho'_1 \hat{\Theta}_Z \hat{\Sigma}_{2,N} \rho_2 - \rho'_1 \hat{\Theta}_Z \tilde{\Sigma}_{2,N} \rho_2| = o_p(1) \quad (10.17)$$

$$|\rho'_1 \hat{\Theta}_Z \tilde{\Sigma}_{2,N} \rho_2 - \rho'_1 \hat{\Theta}_Z \Sigma_{2,N} \rho_2| = o_p(1) \quad (10.18)$$

$$|\rho'_1 \hat{\Theta}_Z \Sigma_{2,N} \rho_2 - \rho'_1 \Theta_Z \Sigma_{2,N} \rho_2| = o_p(1). \quad (10.19)$$

Consider (10.17) first. Note that

$$\begin{aligned} |\rho'_1 \hat{\Theta}_Z \hat{\Sigma}_{2,N} \rho_2 - \rho'_1 \hat{\Theta}_Z \tilde{\Sigma}_{2,N} \rho_2| &\leq \left\| \rho'_1 \hat{\Theta}_Z (\hat{\Sigma}_{2,N} - \tilde{\Sigma}_{2,N}) \right\|_{\infty} \|\rho_2\|_1 \\ &\leq \|\rho'_1 \hat{\Theta}_Z\|_1 \|\hat{\Sigma}_{2,N} - \tilde{\Sigma}_{2,N}\|_{\infty} \sqrt{h_2} = O_p \left( \sqrt{h_1 h_2 \bar{G} \lambda_{node}^{-\vartheta}} \right) \|\hat{\Sigma}_{2,N} - \tilde{\Sigma}_{2,N}\|_{\infty}, \end{aligned}$$

where the last equality is due to (10.9). In addition,

$$\begin{aligned} \|\hat{\Sigma}_{2,N} - \tilde{\Sigma}_{2,N}\|_{\infty} &= \left\| \frac{1}{\sqrt{NT}} \sum_{i=1}^N \sum_{t=1}^T \varepsilon_{i,t}^2 z_{i,t} d'_{i,t} - \frac{1}{\sqrt{NT}} \sum_{i=1}^N \sum_{t=1}^T \varepsilon_{i,t}^2 z_{i,t} d'_{i,t} \right\|_{\infty} \\ &\leq 2 \left\| \frac{1}{\sqrt{NT}} \sum_{i=1}^N \sum_{t=1}^T z_{i,t} d'_{i,t} \varepsilon_{i,t} \pi'_{i,t} (\hat{\gamma} - \gamma) \right\|_{\infty} + \left\| \frac{1}{\sqrt{NT}} \sum_{i=1}^N \sum_{t=1}^T z_{i,t} d'_{i,t} [\pi'_{i,t} (\hat{\gamma} - \gamma)]^2 \right\|_{\infty} \end{aligned} \quad (10.20)$$

Consider the first term of (10.20). A typical element of  $\frac{1}{\sqrt{NT}} \sum_{i=1}^N \sum_{t=1}^T z_{i,t} d'_{i,t} \varepsilon_{i,t} \pi'_{i,t} (\hat{\gamma} - \gamma)$  is

$$\begin{aligned} \frac{1}{\sqrt{NT}} \sum_{j=1}^{NT} z_{j,l} d_{j,k} \varepsilon_j \pi'_j (\hat{\gamma} - \gamma) &\leq \frac{1}{\sqrt{NT}} \left( \sum_{j=1}^{NT} z_{j,l}^2 d_{j,k}^2 \varepsilon_j^2 \right)^{1/2} \left( \sum_{j=1}^{NT} [\pi'_j (\hat{\gamma} - \gamma)]^2 \right)^{1/2} \\ &= \left( \frac{1}{T} \sum_{i=1}^N \sum_{t=1}^T z_{i,t,l}^2 d_{i,t,k}^2 \varepsilon_{i,t}^2 \right)^{1/2} \frac{1}{\sqrt{NT}} \|\Pi(\hat{\gamma} - \gamma)\| = \left( \frac{1}{T} \sum_{t=1}^T z_{k,t,l}^2 \varepsilon_{k,t}^2 \right)^{1/2} \frac{1}{\sqrt{NT}} \|\Pi(\hat{\gamma} - \gamma)\| \end{aligned}$$

for some  $l \in \{1, \dots, p\}$  and  $k \in \{1, \dots, N\}$  where the inequality is due to Cauchy-Schwarz inequality. By subgaussianity, Assumption 3, we can use the same technique as in (12.3) in Proposition 3 in Appendix F to prove  $\mathbb{E} e^{D \left| \frac{1}{T} \sum_{t=1}^T z_{i,t,l}^2 \varepsilon_{i,t}^2 \right|^{1/2}} \leq BT$  for positive constants  $D, B$ .

Using Markov's inequality, we have for  $\epsilon > 0$

$$\mathbb{P} \left( \max_{1 \leq l \leq p} \max_{1 \leq k \leq N} \left| \frac{1}{T} \sum_{t=1}^T z_{k,t,l}^2 \varepsilon_{k,t}^2 \right| \geq \epsilon \right) \leq \sum_{l=1}^p \sum_{k=1}^N \frac{\mathbb{E} e^{D \left| \frac{1}{T} \sum_{t=1}^T z_{k,t,l}^2 \varepsilon_{k,t}^2 \right|^{1/2}}}{e^{D\epsilon^{1/2}}} \leq BpNT e^{-D\epsilon^{1/2}}.$$

Set  $\epsilon = M(\log(pNT))^2$  for some  $M > 0$  and note that the upper bound of the preceding probability becomes arbitrarily small for  $N$  and  $M$  sufficiently large. Thus,  $\max_{1 \leq l \leq p} \max_{1 \leq k \leq N} \left| \frac{1}{T} \sum_{t=1}^T z_{k,t,l}^2 \varepsilon_{k,t}^2 \right| = O_p((\log(pNT))^2)$ . Therefore,

$$\begin{aligned} \left\| \frac{1}{\sqrt{NT}} \sum_{i=1}^N \sum_{t=1}^T z_{i,t} d'_{i,t} \varepsilon_{i,t} \pi'_{i,t} (\hat{\gamma} - \gamma) \right\|_{\infty} &\leq \left( \max_{1 \leq l \leq p} \max_{1 \leq k \leq N} \frac{1}{T} \sum_{t=1}^T z_{k,t,l}^2 \varepsilon_{k,t}^2 \right)^{1/2} \frac{1}{\sqrt{NT}} \|\Pi(\hat{\gamma} - \gamma)\| \\ &\leq O_p(\log(pNT)) \frac{1}{\sqrt{NT}} \|\Pi(\hat{\gamma} - \gamma)\|. \end{aligned} \quad (10.21)$$

Now consider the second term of (10.20). A typical element of  $\frac{1}{\sqrt{NT}} \sum_{i=1}^N \sum_{t=1}^T z_{i,t} d'_{i,t} [\pi'_{i,t} (\hat{\gamma} - \gamma)]^2$  is

$$\begin{aligned} \frac{1}{\sqrt{NT}} \sum_{i=1}^N \sum_{t=1}^T z_{i,t,l} d_{i,t,k} [\pi'_{i,t} (\hat{\gamma} - \gamma)]^2 &\leq \max_{1 \leq i \leq N} \max_{1 \leq t \leq T} \sqrt{N} |z_{i,t,l} d_{i,t,k}| \frac{1}{NT} \|\Pi(\hat{\gamma} - \gamma)\|^2 \\ &\leq \max_{1 \leq t \leq T} \sqrt{N} |z_{k,t,l}| \frac{1}{NT} \|\Pi(\hat{\gamma} - \gamma)\|^2 \end{aligned}$$

for some  $l \in \{1, \dots, p\}$ ,  $k \in \{1, \dots, N\}$ . Using Markov's inequality, we have for any  $\epsilon > 0$

$$\mathbb{P} \left( \max_{1 \leq l \leq p} \max_{1 \leq k \leq N} \max_{1 \leq t \leq T} |z_{k,t,l}| \geq \epsilon \right) \leq \sum_{l=1}^p \sum_{k=1}^N \sum_{t=1}^T \mathbb{P}(|z_{k,t,l}| \geq \epsilon) \leq pNT \frac{K}{2} e^{-C\epsilon^2}.$$

Set  $\epsilon = \sqrt{M \log(pNT)}$  for some  $M > 0$  to see that the upper bound of the preceding probability becomes arbitrarily small for  $N$  and  $M$  sufficiently large. Thus,  $\max_{1 \leq l \leq p} \max_{1 \leq k \leq N} \max_{1 \leq t \leq T} |z_{k,t,l}| = O_p(\sqrt{\log(pNT)})$ . In total,

$$\begin{aligned} \left\| \frac{1}{\sqrt{NT}} \sum_{i=1}^N \sum_{t=1}^T z_{i,t} d'_{i,t} [\pi'_{i,t} (\hat{\gamma} - \gamma)]^2 \right\|_{\infty} &\leq \max_{1 \leq l \leq p} \max_{1 \leq k \leq N} \max_{1 \leq t \leq T} \sqrt{N} |z_{k,t,l}| \frac{1}{NT} \|\Pi(\hat{\gamma} - \gamma)\|^2 \\ &= O_p(\sqrt{N \log(pNT)}) \frac{1}{NT} \|\Pi(\hat{\gamma} - \gamma)\|^2. \end{aligned} \quad (10.22)$$

Therefore, combining (10.21) and (10.22)

$$\begin{aligned} |\rho'_1 \hat{\Theta}_Z \hat{\Sigma}_{2,N} \rho_2 - \rho'_1 \hat{\Theta}_Z \tilde{\Sigma}_{2,N} \rho_2| &\leq \|\hat{\Sigma}_{2,N} - \tilde{\Sigma}_{2,N}\|_{\infty} O_p(\sqrt{h_1 h_2 \bar{G} \lambda_{node}^{-\vartheta}}) \\ &= O_p(\sqrt{h_1 h_2 \bar{G} \lambda_{node}^{-\vartheta}} \log(pNT)) \frac{1}{\sqrt{NT}} \|\Pi(\hat{\gamma} - \gamma)\| + O_p(\sqrt{h_1 h_2 \bar{G} \lambda_{node}^{-\vartheta}} N \log(pNT)) \frac{1}{NT} \|\Pi(\hat{\gamma} - \gamma)\|^2 \\ &= o_p(1), \end{aligned}$$

where the last equality is due to Assumption 6(b)(iii)-(iv), which establishes (10.17).

Next, turn to (10.18). Note that

$$|\rho'_1 \hat{\Theta}_Z \tilde{\Sigma}_{2,N} \rho_2 - \rho'_1 \hat{\Theta}_Z \Sigma_{2,N} \rho_2| \leq \|\tilde{\Sigma}_{2,N} - \Sigma_{2,N}\|_{\infty} \|\hat{\Theta}'_Z \rho_1\|_1 \sqrt{h_2}.$$



Given (10.9), it suffices to consider

$$\begin{aligned}\|\tilde{\Sigma}_{2,N} - \Sigma_{2,N}\|_\infty &= \max_{1 \leq l \leq p} \max_{1 \leq k \leq N} \left| \frac{1}{\sqrt{NT}} \sum_{i=1}^N \sum_{t=1}^T (z_{i,t,l} d_{i,t,k} \varepsilon_{i,t}^2 - \mathbb{E}[z_{i,t,l} d_{i,t,k} \varepsilon_{i,t}^2]) \right| \\ &= \max_{1 \leq l \leq p} \max_{1 \leq k \leq N} \left| \frac{1}{\sqrt{NT}} \sum_{t=1}^T (z_{k,t,l} \varepsilon_{k,t}^2 - \mathbb{E}[z_{k,t,l} \varepsilon_{k,t}^2]) \right|.\end{aligned}$$

By subgaussianity, Assumption 3, we can use the same technique as in (12.3) in Proposition 3 in Appendix F to prove  $\mathbb{E} e^{D|\frac{1}{T} \sum_{t=1}^T (z_{k,t,l} \varepsilon_{k,t}^2 - \mathbb{E}[z_{k,t,l} \varepsilon_{k,t}^2])|^{2/3}} \leq BT$  for some positive constant  $B$ .

Using Markov's inequality, we have for any  $\epsilon > 0$

$$\begin{aligned}\mathbb{P} \left( \max_{1 \leq l \leq p} \max_{1 \leq k \leq N} \left| \frac{1}{T} \sum_{t=1}^T (z_{k,t,l} \varepsilon_{k,t}^2 - \mathbb{E}[z_{k,t,l} \varepsilon_{k,t}^2]) \right| \geq \epsilon \right) &\leq \sum_{l=1}^p \sum_{k=1}^N \mathbb{P} \left( \left| \frac{1}{T} \sum_{t=1}^T (z_{k,t,l} \varepsilon_{k,t}^2 - \mathbb{E}[z_{k,t,l} \varepsilon_{k,t}^2]) \right| \geq \epsilon \right) \\ &\leq \sum_{l=1}^p \sum_{k=1}^N \frac{\mathbb{E} e^{D|\frac{1}{T} \sum_{t=1}^T (z_{k,t,l} \varepsilon_{k,t}^2 - \mathbb{E}[z_{k,t,l} \varepsilon_{k,t}^2])|^{2/3}}}{e^{D\epsilon^{2/3}}} \leq BpNTe^{-D\epsilon^{2/3}}.\end{aligned}$$

Set  $\epsilon = \sqrt{M(\log(pNT))^3}$  for some  $M > 0$  and note that the upper bound of the preceding probability becomes arbitrarily small for  $N$  and  $M$  sufficiently large. Thus,

$$\max_{1 \leq l \leq p} \max_{1 \leq k \leq N} \left| \frac{1}{T} \sum_{t=1}^T (z_{k,t,l} \varepsilon_{k,t}^2 - \mathbb{E}[z_{k,t,l} \varepsilon_{k,t}^2]) \right| = O_p \left( \sqrt{(\log(pNT))^3} \right)$$

and so

$$\|\tilde{\Sigma}_{2,N} - \Sigma_{2,N}\|_\infty = \frac{1}{\sqrt{N}} \max_{1 \leq l \leq p} \max_{1 \leq k \leq N} \left| \frac{1}{T} \sum_{t=1}^T (z_{k,t,l} \varepsilon_{k,t}^2 - \mathbb{E}[z_{k,t,l} \varepsilon_{k,t}^2]) \right| = O_p \left( \sqrt{\frac{(\log(pNT))^3}{N}} \right). \quad (10.23)$$

In total,

$$|\rho'_1 \hat{\Theta}_Z \tilde{\Sigma}_{2,N} \rho_2 - \rho'_1 \hat{\Theta}_Z \Sigma_{2,N} \rho_2| = O_p \left( \sqrt{\frac{(\log(p \vee N \vee T))^3 h_1 h_2 \bar{G} \lambda_{node}^{-\vartheta}}{N}} \right) = o_p(1),$$

where the last equality is due to Assumption 6(a)(ii), establishing (10.18).

We now establish (10.19).

$$\begin{aligned}|\rho'_1 \hat{\Theta}_Z \Sigma_{2,N} \rho_2 - \rho'_1 \Theta_Z \Sigma_{2,N} \rho_2| &\leq \|\Sigma_{2,N}\|_\infty \|(\hat{\Theta}'_Z - \Theta'_Z) \rho_1\|_1 \sqrt{h_2} \\ &= \|\Sigma_{2,N}\|_\infty O_p(\bar{G} \lambda_{node}^{1-\vartheta} \sqrt{h_1 h_2}) = O(1/\sqrt{N}) O_p(\bar{G} \lambda_{node}^{1-\vartheta} \sqrt{h_1 h_2}) = o_p(1),\end{aligned}$$

where the first equality is due to (10.16), the second equality is due to the definition of  $\Sigma_{2,N}$  and (12.1), and the last equality is due to Assumption 6(a)(ii) and 4b). Thus, we have established (10.19), concluding the proof that (10.4) is  $o_p(1)$ .

(10.5) is  $o_p(1)$ :

We now prove that (10.5) is  $o_p(1)$ . First,

$$|\rho'_2 \hat{\Sigma}_{3,N} \rho_2 - \rho'_2 \Sigma_{3,N} \rho_2| \leq \|\hat{\Sigma}_{3,N} - \Sigma_{3,N}\|_\infty h_2 \leq h_2 (\|\hat{\Sigma}_{3,N} - \tilde{\Sigma}_{3,N}\|_\infty + \|\tilde{\Sigma}_{3,N} - \Sigma_{3,N}\|_\infty),$$

where  $\tilde{\Sigma}_{3,N} := \frac{1}{T} \sum_{i=1}^N \sum_{t=1}^T \varepsilon_{i,t}^2 d_{i,t} d'_{i,t}$ . We consider  $\|\hat{\Sigma}_{3,N} - \tilde{\Sigma}_{3,N}\|_\infty$  first.

$$\begin{aligned} \|\hat{\Sigma}_{3,N} - \tilde{\Sigma}_{3,N}\|_\infty &= \left\| \frac{1}{T} \sum_{i=1}^N \sum_{t=1}^T \hat{\varepsilon}_{i,t}^2 d_{i,t} d'_{i,t} - \frac{1}{T} \sum_{i=1}^N \sum_{t=1}^T \varepsilon_{i,t}^2 d_{i,t} d'_{i,t} \right\|_\infty \\ &\leq 2 \left\| \frac{1}{T} \sum_{i=1}^N \sum_{t=1}^T d_{i,t} d'_{i,t} \varepsilon_{i,t} \pi'_{i,t} (\hat{\gamma} - \gamma) \right\|_\infty + \left\| \frac{1}{T} \sum_{i=1}^N \sum_{t=1}^T d_{i,t} d'_{i,t} [\pi'_{i,t} (\hat{\gamma} - \gamma)]^2 \right\|_\infty. \end{aligned} \quad (10.24)$$

Consider the first term of (10.24). A typical element of  $\frac{1}{T} \sum_{i=1}^N \sum_{t=1}^T d_{i,t} d'_{i,t} \varepsilon_{i,t} \pi'_{i,t} (\hat{\gamma} - \gamma)$  is

$$\begin{aligned} \frac{1}{T} \sum_{j=1}^{NT} d_{j,l} d_{j,k} \varepsilon_j \pi'_j (\hat{\gamma} - \gamma) &\leq \frac{1}{T} \left( \sum_{j=1}^{NT} d_{j,l}^2 d_{j,k}^2 \varepsilon_j^2 \right)^{1/2} \left( \sum_{j=1}^{NT} [\pi'_j (\hat{\gamma} - \gamma)]^2 \right)^{1/2} \\ &= \left( \frac{1}{T} \sum_{i=1}^N \sum_{t=1}^T d_{i,t}^2 d_{i,t,k}^2 \varepsilon_{i,t}^2 \right)^{1/2} \frac{1}{\sqrt{T}} \|\Pi(\hat{\gamma} - \gamma)\| = \left( \frac{1}{T} \sum_{t=1}^T \varepsilon_{k,t}^2 \right)^{1/2} \frac{1}{\sqrt{T}} \|\Pi(\hat{\gamma} - \gamma)\| \end{aligned}$$

for some  $l, k \in \{1, \dots, N\}$ , where the inequality is due to Cauchy-Schwarz inequality. By Assumption 3 we have  $\mathbb{P}(|\varepsilon_{i,t}^2| \geq \epsilon) \leq \mathbb{P}(|\varepsilon_{i,t}| \geq \epsilon^{1/2}) \leq \frac{1}{2} K e^{-C\epsilon}$  for every  $\epsilon > 0$ . It follows from Lemma 2.2.1 in van der Vaart and Wellner (1996) that  $\|\varepsilon_{i,t}^2\|_{\psi_1} \leq (1 + K/2)/C$  for all  $i$  and  $t$ . Hence, by subadditivity of the Orlicz norm and Jensen's inequality,  $\|\varepsilon_{i,t}^2 - \mathbb{E}[\varepsilon_{i,t}^2]\|_{\psi_1} \leq 2\|\varepsilon_{i,t}^2\|_{\psi_1} \leq (2 + K)/C$ . Using the definition of the Orlicz norm, we have  $\mathbb{E} \exp(\frac{C}{2+K} |\varepsilon_{i,t}^2 - \mathbb{E}[\varepsilon_{i,t}^2]|) \leq 2$ . Use independence of  $\varepsilon_{i,t}$  across  $t$  to invoke Proposition 2 in Appendix F for  $D = \frac{C}{2+K}$  and  $\alpha = 1/3$  to conclude

$$\mathbb{P} \left( \left| \sum_{t=1}^T (\varepsilon_{i,t}^2 - \mathbb{E}[\varepsilon_{i,t}^2]) \right| \geq T\epsilon \right) \leq A e^{-B(\epsilon^2 T)^{1/3}},$$

for positive constants  $A$  and  $B$ . Setting  $\epsilon = \sqrt{\frac{M(\log N)^3}{T}}$  for some  $M > 0$  ( $\epsilon \gtrsim \frac{1}{\sqrt{T}}$ ), one has

$$\mathbb{P} \left( \max_{1 \leq k \leq N} \left| \sum_{t=1}^T (\varepsilon_{k,t}^2 - \mathbb{E}[\varepsilon_{k,t}^2]) \right| \geq T\epsilon \right) \leq \sum_{k=1}^N \mathbb{P} \left( \left| \sum_{t=1}^T (\varepsilon_{k,t}^2 - \mathbb{E}[\varepsilon_{k,t}^2]) \right| \geq T\epsilon \right) \leq A N^{1-BM^{1/3}}.$$

The upper bound of the preceding probability becomes arbitrarily small for  $N$  and  $M$  sufficiently large. Hence,

$$\max_{1 \leq k \leq N} \left| \frac{1}{T} \sum_{t=1}^T (\varepsilon_{k,t}^2 - \mathbb{E}[\varepsilon_{k,t}^2]) \right| = O_p \left( \sqrt{\frac{(\log N)^3}{T}} \right). \quad (10.25)$$

Furthermore, since  $\max_{1 \leq k \leq N} \max_{1 \leq t \leq T} \mathbb{E}[\varepsilon_{k,t}^2] \leq \max_{1 \leq k \leq N} \max_{1 \leq t \leq T} \|\varepsilon_{k,t}^2\|_{\psi_1} \leq (1 + K/2)/C = O(1)$

$$\begin{aligned} \max_{1 \leq k \leq N} \left| \frac{1}{T} \sum_{t=1}^T \varepsilon_{k,t}^2 \right| &\leq \max_{1 \leq k \leq N} \left| \frac{1}{T} \sum_{t=1}^T (\varepsilon_{k,t}^2 - \mathbb{E}[\varepsilon_{k,t}^2]) \right| + \max_{1 \leq k \leq N} \max_{1 \leq t \leq T} \mathbb{E}[\varepsilon_{k,t}^2] = O_p \left( \sqrt{\frac{(\log N)^3}{T}} \right) + O(1). \end{aligned} \quad (10.26)$$

Therefore,

$$\left\| \frac{1}{T} \sum_{i=1}^N \sum_{t=1}^T d_{i,t} d'_{i,t} \varepsilon_{i,t} \pi'_{i,t} (\hat{\gamma} - \gamma) \right\|_{\infty} = O_p \left( \frac{(\log N)^{3/4}}{T^{1/4}} \vee 1 \right) \frac{1}{\sqrt{T}} \|\Pi(\hat{\gamma} - \gamma)\|. \quad (10.27)$$

Now consider the second term of (10.24). A typical element of  $\frac{1}{T} \sum_{i=1}^N \sum_{t=1}^T d_{i,t} d'_{i,t} [\pi'_{i,t} (\hat{\gamma} - \gamma)]^2$  is

$$\frac{1}{T} \sum_{i=1}^N \sum_{t=1}^T d_{i,t,l} d_{i,t,k} [\pi'_{i,t} (\hat{\gamma} - \gamma)]^2 \leq \max_{1 \leq i \leq N} \max_{1 \leq t \leq T} |d_{i,t,l} d_{i,t,k}| \frac{1}{T} \|\Pi(\hat{\gamma} - \gamma)\|^2 = \frac{1}{T} \|\Pi(\hat{\gamma} - \gamma)\|^2, \quad (10.28)$$

uniformly over  $l, k \in \{1, \dots, N\}$ . Combining (10.27) and (10.28), we have

$$\left\| \hat{\Sigma}_{3,N} - \tilde{\Sigma}_{3,N} \right\|_{\infty} = O_p \left( \frac{(\log N)^{3/4}}{T^{1/4}} \vee 1 \right) \frac{1}{\sqrt{T}} \|\Pi(\hat{\gamma} - \gamma)\| + \frac{1}{T} \|\Pi(\hat{\gamma} - \gamma)\|^2. \quad (10.29)$$

Next, consider  $\left\| \tilde{\Sigma}_{3,N} - \Sigma_{3,N} \right\|_{\infty}$ .

$$\begin{aligned} \left\| \tilde{\Sigma}_{3,N} - \Sigma_{3,N} \right\|_{\infty} &= \max_{1 \leq l \leq N} \max_{1 \leq k \leq N} \left| \frac{1}{T} \sum_{i=1}^N \sum_{t=1}^T (\varepsilon_{i,t}^2 d_{i,t,l} d_{i,t,k} - \mathbb{E}[\varepsilon_{i,t}^2 d_{i,t,l} d_{i,t,k}]) \right| \\ &= \max_{1 \leq k \leq N} \left| \frac{1}{T} \sum_{t=1}^T (\varepsilon_{k,t}^2 - \mathbb{E}[\varepsilon_{k,t}^2]) \right| = O_p \left( \sqrt{\frac{(\log N)^3}{T}} \right), \end{aligned} \quad (10.30)$$

where the last equality is due to (10.25). Summing up (10.29) and (10.30) yields

$$\begin{aligned} &|\rho'_2 \hat{\Sigma}_{3,N} \rho_2 - \rho'_2 \Sigma_{3,N} \rho_2| \\ &= h_2 O_p \left( \frac{(\log N)^{3/4}}{T^{1/4}} \vee 1 \right) \frac{1}{\sqrt{T}} \|\Pi(\hat{\gamma} - \gamma)\| + h_2 \frac{1}{T} \|\Pi(\hat{\gamma} - \gamma)\|^2 + O_p \left( h_2 \sqrt{\frac{(\log N)^3}{T}} \right) \\ &= o_p(1), \end{aligned}$$

where the last equality is due to Assumptions 6(b)(v), which, in turns, implies that (10.5) is  $o_p(1)$ .

Thus, we have proved (10.2). (3.8) then follows trivially since the conclusions of Theorem 1 and Corollary 1 are uniform over the set  $\mathcal{F}(\nu, E)$  and the true parameter vector only entered the above arguments when these results were used.

### 10.1.2 Numerators of $t_1$ and $t'_1$

We now show that the numerators of  $t_1$  and  $t'_1$  are asymptotically equivalent, i.e.,

$$|\rho' \hat{\Theta} S^{-1} \Pi' \varepsilon - \rho' \Theta S^{-1} \Pi' \varepsilon| = o_p(1). \quad (10.31)$$

Note that

$$\begin{aligned} |\rho' \hat{\Theta} S^{-1} \Pi' \varepsilon - \rho' \Theta S^{-1} \Pi' \varepsilon| &\leq \|\rho'(\hat{\Theta} - \Theta)\|_1 \|S^{-1} \Pi' \varepsilon\|_\infty = \|\rho'_1(\hat{\Theta}_Z - \Theta_Z)\|_1 \|S^{-1} \Pi' \varepsilon\|_\infty \\ &= O_p(\bar{G} \lambda_{node}^{1-\vartheta} \sqrt{h_1}) \left( \frac{1}{\sqrt{NT}} \|Z' \varepsilon\|_\infty \vee \frac{1}{\sqrt{T}} \|D' \varepsilon\|_\infty \right) = O_p(\bar{G} \lambda_{node}^{1-\vartheta} \sqrt{h_1}) O_p(\sqrt{(\log(p \vee N))^3}) = o_p(1), \end{aligned}$$

where the second equality is due to (10.16), and the third equality is due to (8.12) and (8.13), and the last equality is due to Assumption 6(a)(iii).

### 10.1.3 $t'_1 \xrightarrow{d} N(0, 1)$

We now prove that  $t'_1$  is asymptotically distributed as a standard normal by verifying (i)-(iii) of Theorem 5 in Appendix F. Note that

$$t'_1 := \frac{\rho' \Theta S^{-1} \Pi' \varepsilon}{\sqrt{\rho' \Theta \Sigma_{\Pi \varepsilon} \Theta' \rho}} = \frac{\rho' \Theta S^{-1} \sum_{i=1}^N \sum_{t=1}^T \begin{pmatrix} z_{i,t} \varepsilon_{i,t} \\ d_{i,t} \varepsilon_{i,t} \end{pmatrix}}{\sqrt{\rho' \Theta \Sigma_{\Pi \varepsilon} \Theta' \rho}} = \frac{\rho' \Theta S^{-1} \sum_{j=1}^k \begin{pmatrix} z_j \varepsilon_j \\ d_j \varepsilon_j \end{pmatrix}}{\sqrt{\rho' \Theta \Sigma_{\Pi \varepsilon} \Theta' \rho}},$$

where  $k := NT$ . In the proof of Lemma 5, we have shown that  $t'_1$  is a martingale difference array with variance

$$\text{var}(t'_1) = \mathbb{E}[t_1'^2] = \frac{\rho' \Theta S^{-1} \mathbb{E}[\Pi' \varepsilon \varepsilon' \Pi] S^{-1} \Theta' \rho}{\rho' \Theta \Sigma_{\Pi \varepsilon} \Theta' \rho} = 1$$

where we have used the definition of  $\Sigma_{\Pi \varepsilon}$ . We have already shown in (10.1) that the denominator of  $t'_1$  is uniformly bounded away from zero. Thus, verifying that  $t'_1$  satisfies (i) and (ii) of Theorem 5 in Appendix F is equivalent to verifying that the numerator of  $t'_1$  satisfies (i) and (ii) of Theorem 5. First, note that

$$\|\rho'_1 \Theta_Z\|_1 = \left\| \sum_{j \in H_1} \rho_{1j} \Theta'_{Z,j} \right\|_1 \leq \sum_{j \in H_1} |\rho_{1j}| \|\Theta'_{Z,j}\|_1 = O(\sqrt{h_1 \bar{G} \lambda_{node}^{-\vartheta}}), \quad (10.32)$$

where the last equality is due to (9.18). Next,

$$\begin{aligned} \left| \rho' \Theta S^{-1} \begin{pmatrix} z_{i,t} \varepsilon_{i,t} \\ d_{i,t} \varepsilon_{i,t} \end{pmatrix} \right| &\leq \left| \rho'_1 \Theta_Z \frac{z_{i,t} \varepsilon_{i,t}}{\sqrt{NT}} \right| + \left| \frac{\rho_{2,i} \varepsilon_{i,t}}{\sqrt{T}} \right| \leq \|\rho'_1 \Theta_Z\|_1 \max_{1 \leq l \leq p} \left| \frac{z_{i,t,l} \varepsilon_{i,t}}{\sqrt{NT}} \right| + \frac{\|\rho_2\|_\infty |\varepsilon_{i,t}|}{\sqrt{T}} \\ &\lesssim \sqrt{h_1 \bar{G} \lambda_{node}^{-\vartheta}} \max_{1 \leq l \leq p} \left| \frac{z_{i,t,l} \varepsilon_{i,t}}{\sqrt{NT}} \right| + \frac{\|\rho_2\|_\infty |\varepsilon_{i,t}|}{\sqrt{T}}, \end{aligned}$$

where the last inequality due to (10.32). We have already shown in the proof of Lemma 5 that  $z_{i,t,l} \varepsilon_{i,t}$  has uniformly bounded  $\psi_1$ -Orlicz norm. The same is the case for  $\varepsilon_{i,t}$ . Hence,

$$\begin{aligned} \left\| \sqrt{h_1 \bar{G} \lambda_{node}^{-\vartheta}} \max_{1 \leq l \leq p} \left| \frac{z_{i,t,l} \varepsilon_{i,t}}{\sqrt{NT}} \right| + \frac{\|\rho_2\|_\infty |\varepsilon_{i,t}|}{\sqrt{T}} \right\|_{\psi_1} &\leq \sqrt{\frac{h_1 \bar{G} \lambda_{node}^{-\vartheta}}{NT}} \left\| \max_{1 \leq l \leq p} z_{i,t,l} \varepsilon_{i,t} \right\|_{\psi_1} + \frac{\|\rho_2\|_\infty}{\sqrt{T}} \|\varepsilon_{i,t}\|_{\psi_1} \\ &\lesssim \sqrt{\frac{h_1 \bar{G} \lambda_{node}^{-\vartheta}}{NT} \log(1+p)} \max_{1 \leq l \leq p} \|z_{i,t,l} \varepsilon_{i,t}\|_{\psi_1} + \frac{\|\rho_2\|_\infty}{\sqrt{T}} \|\varepsilon_{i,t}\|_{\psi_1} \lesssim \sqrt{\frac{h_1 \bar{G} \lambda_{node}^{-\vartheta}}{NT} \log(1+p)} + \frac{\|\rho_2\|_\infty}{\sqrt{T}}, \end{aligned}$$

for all  $i$  and  $T$ , where the first rate inequality is due to Lemma 2.2.2 in van der Vaart and Wellner (1996). Using Lemma 2.2.2 in van der Vaart and Wellner (1996) one more time,

$$\left\| \max_{1 \leq i \leq N} \max_{1 \leq t \leq T} \left| \rho' \Theta S^{-1} \begin{pmatrix} z_{i,t} \varepsilon_{i,t} \\ d_{i,t} \varepsilon_{i,t} \end{pmatrix} \right| \right\|_{\psi_1} \lesssim \log(1 + NT) \left[ \sqrt{\frac{h_1 \bar{G} \lambda_{node}^{-\vartheta}}{NT}} \log(1 + p) + \frac{\|\rho_2\|_\infty}{\sqrt{T}} \right] = o(1),$$

where the last equality is due to Assumption 6(a)(iv)-(v). Since  $\|U\|_{L_r} \leq r! \|U\|_{\psi_1}$  for any random variable  $U$  (van der Vaart and Wellner (1996), p95), we conclude that (i) and (ii) of Theorem 5 are satisfied.

We now verify (iii) of Theorem 5. That is,

$$\frac{\sum_{j=1}^{k_N} \left[ \rho' \Theta S^{-1} \begin{pmatrix} z_j \varepsilon_j \\ d_j \varepsilon_j \end{pmatrix} \right]^2}{\rho' \Theta \Sigma_{\Pi \varepsilon} \Theta' \rho} = \frac{\rho' \Theta \begin{pmatrix} \tilde{\Sigma}_{1,N} & \tilde{\Sigma}_{2,N} \\ \tilde{\Sigma}'_{2,N} & \tilde{\Sigma}_{3,N} \end{pmatrix} \Theta' \rho}{\rho' \Theta \Sigma_{\Pi \varepsilon} \Theta' \rho} \xrightarrow{p} 1.$$

Since we have already shown in (10.1) that the denominator of  $t'_1$  is uniformly bounded away from zero, it suffices to show

$$\left| \rho' \Theta \begin{pmatrix} \tilde{\Sigma}_{1,N} & \tilde{\Sigma}_{2,N} \\ \tilde{\Sigma}'_{2,N} & \tilde{\Sigma}_{3,N} \end{pmatrix} \Theta' \rho - \rho' \Theta \Sigma_{\Pi \varepsilon} \Theta' \rho \right| = o_p(1). \quad (10.33)$$

The left-hand side of (10.33) can be bounded by

$$\left| \rho' \Theta \begin{pmatrix} \tilde{\Sigma}_{1,N} & \tilde{\Sigma}_{2,N} \\ \tilde{\Sigma}'_{2,N} & \tilde{\Sigma}_{3,N} \end{pmatrix} \Theta' \rho - \rho' \Theta \Sigma_{\Pi \varepsilon} \Theta' \rho \right| \leq |\rho'_1 \Theta_Z \tilde{\Sigma}_{1,N} \Theta'_Z \rho_1 - \rho'_1 \Theta_Z \Sigma_{1,N} \Theta'_Z \rho_1| \quad (10.34)$$

$$+ 2|\rho'_1 \Theta_Z \tilde{\Sigma}_{2,N} \rho_2 - \rho'_1 \Theta_Z \Sigma_{2,N} \rho_2| \quad (10.35)$$

$$+ |\rho'_2 \tilde{\Sigma}_{3,N} \rho_2 - \rho'_2 \Sigma_{3,N} \rho_2|. \quad (10.36)$$

Thus, we establish that (10.34), (10.35) and (10.36) are  $o_p(1)$ . Consider (10.34) first.

$$\begin{aligned} |\rho'_1 \Theta_Z \tilde{\Sigma}_{1,N} \Theta'_Z \rho_1 - \rho'_1 \Theta_Z \Sigma_{1,N} \Theta'_Z \rho_1| &\leq \|\tilde{\Sigma}_{1,N} - \Sigma_{1,N}\|_\infty \|\Theta'_Z \rho_1\|_1^2 \\ &= O_p \left( \sqrt{\frac{(\log(p^2 T))^5}{N}} \right) O(h_1 \bar{G} \lambda_{node}^{-\vartheta}) = o_p(1) \end{aligned}$$

where the first equality is due to (10.32) and (10.15), and the last equality is due to Assumption 6(a)(i). Now consider (10.35).

$$\begin{aligned} |\rho'_1 \Theta_Z \tilde{\Sigma}_{2,N} \rho_2 - \rho'_1 \Theta_Z \Sigma_{2,N} \rho_2| &\leq \|\tilde{\Sigma}_{2,N} - \Sigma_{2,N}\|_\infty \|\Theta'_Z \rho_1\|_1 \|\rho_2\|_1 \\ &= O_p \left( \sqrt{\frac{(\log(pNT))^3 h_1 h_2 \bar{G} \lambda_{node}^{-\vartheta}}{N}} \right) = o_p(1), \end{aligned}$$

where the first equality is due to (10.23), and the last equality is due to Assumption 6(a)(ii). Finally, consider (10.36).

$$|\rho'_2 \tilde{\Sigma}_{3,N} \rho_2 - \rho'_2 \Sigma_{3,N} \rho_2| \leq \|\tilde{\Sigma}_{3,N} - \Sigma_{3,N}\|_\infty \|\rho_2\|_1^2 = O_p\left(\sqrt{\frac{(\log N)^3}{T}}\right) O(h_2) = o_p(1),$$

where the first equality is due to (10.30), and the last equality is due to Assumption 6(b)(v). Therefore, we have established (10.33) and  $t'_1$  is asymptotically standard gaussian.

#### 10.1.4 $t_2 = o_p(1)$

Last, we prove that  $t_2 = o_p(1)$ . Since the denominator of  $t_2$  is bounded away from zero by a positive constant with probability approaching one by (10.1) and (10.2), it suffices to show  $\rho' \Delta = o_p(1)$ .

$$\begin{aligned} |\rho' \Delta| &= \left| \sum_{j \in H} \rho_j \Delta_j \right| \leq \sqrt{h} \max_{j \in H} |\Delta_j| \leq \sqrt{h} \|S(\hat{\gamma} - \gamma)\|_1 \max_{j \in H} \|\hat{\Theta}'_j \Psi_N - \mathbf{I}'_{p+N,j}\|_\infty \\ &= \sqrt{h} \|S(\hat{\gamma} - \gamma)\|_1 \left( \max_{j \in H_1} \left\| \begin{pmatrix} \frac{1}{NT} Z' Z \hat{\Theta}_{Z,j} - e_j \\ \frac{1}{T\sqrt{N}} D' Z \hat{\Theta}_{Z,j} \end{pmatrix} \right\|_\infty \vee \max_{i \in H_2} \left\| \begin{pmatrix} \frac{1}{T\sqrt{N}} Z' D e_i \\ 0 \end{pmatrix} \right\|_\infty \right) \\ &= \sqrt{h} \|S(\hat{\gamma} - \gamma)\|_1 \left( \max_{j \in H_1} \left( \left\| \frac{1}{NT} Z' Z \hat{\Theta}_{Z,j} - e_j \right\|_\infty \vee \left\| \frac{1}{T\sqrt{N}} D' Z \hat{\Theta}_{Z,j} \right\|_\infty \right) \vee \max_{i \in H_2} \left\| \frac{1}{T\sqrt{N}} Z' D \right\|_\infty \right) \\ &\leq \sqrt{h} \|S(\hat{\gamma} - \gamma)\|_1 \left( \max_{j \in H_1} \left( \left\| \frac{1}{NT} Z' Z \hat{\Theta}_{Z,j} - e_j \right\|_\infty \vee \|\hat{\Theta}_{Z,j}\|_1 \left\| \frac{1}{T\sqrt{N}} D' Z \right\|_\infty \right) \vee \max_{i \in H_2} \left\| \frac{1}{T\sqrt{N}} Z' D \right\|_\infty \right) \end{aligned}$$

where  $\hat{\Theta}_j$  is the  $j$ th row of  $\hat{\Theta}$  but written as a  $(p+N) \times 1$  vector, and  $\mathbf{I}_{p+N,j}$  is the  $j$ th row of  $\mathbf{I}_{p+N}$  but written as a  $(p+N) \times 1$  vector. Note that

$$\max_{j \in H_1} \left\| \frac{1}{NT} Z' Z \hat{\Theta}_{Z,j} - e_j \right\|_\infty \leq \max_{j \in H_1} \frac{\lambda_{node}}{\hat{\tau}_j^2} = O_p(\lambda_{node}),$$

where the inequality is due to the extended KKT conditions (13.6), and the equality is due to (9.10). Recall that by (8.15) we have that for every  $\epsilon > 0$

$$\mathbb{P} \left( \max_{1 \leq i \leq N} \max_{1 \leq l \leq p} \left| \frac{1}{\sqrt{NT}} \sum_{t=1}^T z_{i,t,l} \right| \geq \epsilon \right) \leq \sum_{i=1}^N \sum_{l=1}^p \mathbb{P} \left( \left| \frac{1}{\sqrt{NT}} \sum_{t=1}^T z_{i,t,l} \right| \geq \epsilon \right) \leq ApNe^{-B\epsilon^2 N},$$

for positive constants  $A, B$ . Setting  $\epsilon = \sqrt{\frac{M \log(pN)}{N}}$  ( $M > 0$ ) makes the upper bound of the preceding inequality arbitrarily small for sufficiently large  $N$  and  $M$ , such that

$$\|\hat{\Theta}_{Z,j}\|_1 \left\| \frac{1}{T\sqrt{N}} D' Z \right\|_\infty = O_p \left( \sqrt{\frac{\bar{G} \lambda_{node}^{-\vartheta} \log(pN)}{N}} \right).$$

Thus,  $|\rho' \Delta| = o_p(1)$  by Assumption 6(c). For later reference,

$$\sup_{\gamma \in \mathcal{F}(\nu, E)} |\rho' \Delta| = o_p(1) \tag{10.37}$$

by the same reasoning leading to the uniform validity of (3.8).  $\square$

## 11 Appendix E

### 11.1 Proof of Theorem 3

*Proof of Theorem 3.* For every  $\epsilon > 0$ , define

$$\begin{aligned} A_{1,N} &:= \left\{ \sup_{\gamma \in \mathcal{F}(\nu, E)} |\rho' \Delta| < \epsilon \right\} & A_{2,N} &:= \left\{ \sup_{\gamma \in \mathcal{F}(\nu, E)} \left| \frac{\sqrt{\rho' \hat{\Theta} \hat{\Sigma}_{\Pi\epsilon} \hat{\Theta}' \rho}}{\sqrt{\rho' \Theta \Sigma_{\Pi\epsilon} \Theta' \rho}} - 1 \right| < \epsilon \right\} \\ A_{3,N} &:= \{ |\rho' \hat{\Theta} S^{-1} \Pi' \epsilon - \rho' \Theta S^{-1} \Pi' \epsilon| < \epsilon \}. \end{aligned}$$

By (10.37), (3.8), (10.1) and (10.31), the probabilities of the preceding three events all tend to one. Thus, for every  $t \in \mathbb{R}$ ,

$$\begin{aligned} & \left| \mathbb{P} \left( \frac{\rho' S(\tilde{\gamma} - \gamma)}{\sqrt{\rho' \hat{\Theta} \hat{\Sigma}_{\Pi\epsilon} \hat{\Theta}' \rho}} \leq t \right) - \Phi(t) \right| \\ & \leq \left| \mathbb{P} \left( \frac{\rho' \hat{\Theta} S^{-1} \Pi' \epsilon}{\sqrt{\rho' \hat{\Theta} \hat{\Sigma}_{\Pi\epsilon} \hat{\Theta}' \rho}} - \frac{\rho' \Delta}{\sqrt{\rho' \hat{\Theta} \hat{\Sigma}_{\Pi\epsilon} \hat{\Theta}' \rho}} \leq t, A_{1,N}, A_{2,N}, A_{3,N} \right) - \Phi(t) \right| + \mathbb{P} \left( \bigcup_{i=1}^3 A_{i,N}^c \right). \end{aligned}$$

We consider  $\mathbb{P} \left( \frac{\rho' \hat{\Theta} S^{-1} \Pi' \epsilon}{\sqrt{\rho' \hat{\Theta} \hat{\Sigma}_{\Pi\epsilon} \hat{\Theta}' \rho}} - \frac{\rho' \Delta}{\sqrt{\rho' \hat{\Theta} \hat{\Sigma}_{\Pi\epsilon} \hat{\Theta}' \rho}} \leq t, A_{1,N}, A_{2,N}, A_{3,N} \right)$  first.

$$\begin{aligned} & \mathbb{P} \left( \frac{\rho' \hat{\Theta} S^{-1} \Pi' \epsilon}{\sqrt{\rho' \hat{\Theta} \hat{\Sigma}_{\Pi\epsilon} \hat{\Theta}' \rho}} - \frac{\rho' \Delta}{\sqrt{\rho' \hat{\Theta} \hat{\Sigma}_{\Pi\epsilon} \hat{\Theta}' \rho}} \leq t, A_{1,N}, A_{2,N}, A_{3,N} \right) \\ & \leq \mathbb{P} \left( \frac{\rho' \Theta S^{-1} \Pi' \epsilon}{\sqrt{\rho' \Theta \Sigma_{\Pi\epsilon} \Theta' \rho}} \leq t(1 + \epsilon) + \frac{\epsilon + \epsilon}{\sqrt{\rho' \Theta \Sigma_{\Pi\epsilon} \Theta' \rho}} \right) \leq \mathbb{P} \left( \frac{\rho' \Theta S^{-1} \Pi' \epsilon}{\sqrt{\rho' \Theta \Sigma_{\Pi\epsilon} \Theta' \rho}} \leq t(1 + \epsilon) + 2D\epsilon \right) \end{aligned}$$

for some positive constant  $D$ , where the first and second inequalities are due to the fact that  $\rho' \Theta \Sigma_{\Pi\epsilon} \Theta' \rho$  is uniformly bounded away from zero, see (10.1). Since the last inequality in the above does not depend on  $\gamma$ ,

$$\begin{aligned} & \sup_{\gamma \in \mathcal{F}(\nu, E)} \mathbb{P} \left( \frac{\rho' \hat{\Theta} S^{-1} \Pi' \epsilon}{\sqrt{\rho' \hat{\Theta} \hat{\Sigma}_{\Pi\epsilon} \hat{\Theta}' \rho}} - \frac{\rho' \Delta}{\sqrt{\rho' \hat{\Theta} \hat{\Sigma}_{\Pi\epsilon} \hat{\Theta}' \rho}} \leq t, A_{1,N}, A_{2,N}, A_{3,N} \right) \\ & \leq \mathbb{P} \left( \frac{\rho' \Theta S^{-1} \Pi' \epsilon}{\sqrt{\rho' \Theta \Sigma_{\Pi\epsilon} \Theta' \rho}} \leq t(1 + \epsilon) + 2D\epsilon \right). \end{aligned}$$

By the asymptotic normality of  $t'_1$ , for  $N$  sufficiently large,

$$\sup_{\gamma \in \mathcal{F}(\nu, E)} \mathbb{P} \left( \frac{\rho' \hat{\Theta} S^{-1} \Pi' \epsilon}{\sqrt{\rho' \hat{\Theta} \hat{\Sigma}_{\Pi\epsilon} \hat{\Theta}' \rho}} - \frac{\rho' \Delta}{\sqrt{\rho' \hat{\Theta} \hat{\Sigma}_{\Pi\epsilon} \hat{\Theta}' \rho}} \leq t, A_{1,N}, A_{2,N}, A_{3,N} \right) \leq \Phi(t(1 + \epsilon) + 2D\epsilon) + \epsilon.$$

As the above arguments are valid for every  $\epsilon > 0$ , we can use the continuity of  $q \mapsto \Phi(q)$  to conclude that for every  $\delta > 0$ , one can choose  $\epsilon$  sufficiently small such that

$$\sup_{\gamma \in \mathcal{F}(\nu, E)} \mathbb{P} \left( \frac{\rho' \hat{\Theta} S^{-1} \Pi' \epsilon}{\sqrt{\rho' \hat{\Theta} \hat{\Sigma}_{\Pi\epsilon} \hat{\Theta}' \rho}} - \frac{\rho' \Delta}{\sqrt{\rho' \hat{\Theta} \hat{\Sigma}_{\Pi\epsilon} \hat{\Theta}' \rho}} \leq t, A_{1,N}, A_{2,N}, A_{3,N} \right) \leq \Phi(t) + \delta + \epsilon. \quad (11.1)$$

We next find a lower bound for  $\mathbb{P}\left(\frac{\rho'\hat{\Theta}S^{-1}\Pi'\varepsilon}{\sqrt{\rho'\hat{\Theta}\hat{\Sigma}_{\Pi\varepsilon}\hat{\Theta}'\rho}} - \frac{\rho'\Delta}{\sqrt{\rho'\hat{\Theta}\hat{\Sigma}_{\Pi\varepsilon}\hat{\Theta}'\rho}} \leq t, A_{1,N}, A_{2,N}, A_{3,N}\right)$ .

$$\begin{aligned} & \mathbb{P}\left(\frac{\rho'\hat{\Theta}S^{-1}\Pi'\varepsilon}{\sqrt{\rho'\hat{\Theta}\hat{\Sigma}_{\Pi\varepsilon}\hat{\Theta}'\rho}} - \frac{\rho'\Delta}{\sqrt{\rho'\hat{\Theta}\hat{\Sigma}_{\Pi\varepsilon}\hat{\Theta}'\rho}} \leq t, A_{1,N}, A_{2,N}, A_{3,N}\right) \\ & \geq \mathbb{P}\left(\frac{\rho'\Theta S^{-1}\Pi'\varepsilon}{\sqrt{\rho'\Theta\Sigma_{\Pi\varepsilon}\Theta'\rho}} \leq t(1-\epsilon) - \frac{\epsilon+\epsilon}{\sqrt{\rho'\Theta\Sigma_{\Pi\varepsilon}\Theta'\rho}}, A_{1,N}, A_{2,N}, A_{3,N}\right) \\ & \geq \mathbb{P}\left(\frac{\rho'\Theta S^{-1}\Pi'\varepsilon}{\sqrt{\rho'\Theta\Sigma_{\Pi\varepsilon}\Theta'\rho}} \leq t(1-\epsilon) - 2D\epsilon, A_{1,N}, A_{2,N}, A_{3,N}\right) \\ & \geq \mathbb{P}\left(\frac{\rho'\Theta S^{-1}\Pi'\varepsilon}{\sqrt{\rho'\Theta\Sigma_{\Pi\varepsilon}\Theta'\rho}} \leq t(1-\epsilon) - 2D\epsilon\right) + \mathbb{P}(\cap_{i=1}^3 A_{i,N}) - 1 \end{aligned}$$

for some positive constant  $D$ , where the first and second inequalities are due to the fact that  $\rho'\Theta\Sigma_{\Pi\varepsilon}\Theta'\rho$  is uniformly bounded away from zero, see (10.1). Since the last inequality in the above display does not depend on  $\gamma$ , and  $\mathbb{P}(\cap_{i=1}^3 A_{i,N})$  can be made arbitrarily close to one for sufficiently large  $N$ ,

$$\begin{aligned} & \inf_{\gamma \in \mathcal{F}(\nu, E)} \mathbb{P}\left(\frac{\rho'\hat{\Theta}S^{-1}\Pi'\varepsilon}{\sqrt{\rho'\hat{\Theta}\hat{\Sigma}_{\Pi\varepsilon}\hat{\Theta}'\rho}} - \frac{\rho'\Delta}{\sqrt{\rho'\hat{\Theta}\hat{\Sigma}_{\Pi\varepsilon}\hat{\Theta}'\rho}} \leq t, A_{1,N}, A_{2,N}, A_{3,N}\right) \\ & \geq \mathbb{P}\left(\frac{\rho'\Theta S^{-1}\Pi'\varepsilon}{\sqrt{\rho'\Theta\Sigma_{\Pi\varepsilon}\Theta'\rho}} \leq t(1-\epsilon) - 2D\epsilon\right) - \epsilon. \end{aligned}$$

By the asymptotic normality of  $t'_1$ , for  $N$  sufficiently large,

$$\inf_{\gamma \in \mathcal{F}(\nu, E)} \mathbb{P}\left(\frac{\rho'\hat{\Theta}S^{-1}\Pi'\varepsilon}{\sqrt{\rho'\hat{\Theta}\hat{\Sigma}_{\Pi\varepsilon}\hat{\Theta}'\rho}} - \frac{\rho'\Delta}{\sqrt{\rho'\hat{\Theta}\hat{\Sigma}_{\Pi\varepsilon}\hat{\Theta}'\rho}} \leq t, A_{1,N}, A_{2,N}, A_{3,N}\right) \geq \Phi(t(1-\epsilon) - 2D\epsilon) - 2\epsilon.$$

As the above arguments are valid for every  $\epsilon > 0$ , we can use the continuity of  $q \mapsto \Phi(q)$  to conclude that for every  $\delta > 0$ , one can choose  $\epsilon$  sufficiently small such that

$$\inf_{\gamma \in \mathcal{F}(\nu, E)} \mathbb{P}\left(\frac{\rho'\hat{\Theta}S^{-1}\Pi'\varepsilon}{\sqrt{\rho'\hat{\Theta}\hat{\Sigma}_{\Pi\varepsilon}\hat{\Theta}'\rho}} - \frac{\rho'\Delta}{\sqrt{\rho'\hat{\Theta}\hat{\Sigma}_{\Pi\varepsilon}\hat{\Theta}'\rho}} \leq t, A_{1,N}, A_{2,N}, A_{3,N}\right) \geq \Phi(t) - \delta - 2\epsilon. \quad (11.2)$$

Thus, by (11.1), (11.2) and the fact that  $\sup_{\gamma \in \mathcal{F}(\nu, E)} \mathbb{P}(\cup_{i=1}^3 A_{i,N}^c) = \mathbb{P}(\cup_{i=1}^3 A_{i,N}^c) = o(1)$ , we have proved (4.1) (the uniformity over  $t \in \mathbb{R}$  follows from the fact that  $\Phi(t)$  is continuous). To see (4.2), note that

$$\begin{aligned} & \mathbb{P}\left(\alpha_j \notin \left[\tilde{\alpha}_j - z_{1-\delta/2} \frac{\tilde{\sigma}_{\alpha,j}}{\sqrt{NT}}, \tilde{\alpha}_j + z_{1-\delta/2} \frac{\tilde{\sigma}_{\alpha,j}}{\sqrt{NT}}\right]\right) = \mathbb{P}\left(\left|\frac{\sqrt{NT}(\tilde{\alpha}_j - \alpha_j)}{\tilde{\sigma}_{z,j}}\right| > z_{1-\delta/2}\right) \\ & \leq 1 - \mathbb{P}\left(\frac{\sqrt{NT}(\tilde{\alpha}_j - \alpha_j)}{\tilde{\sigma}_{z,j}} \leq z_{1-\delta/2}\right) + \mathbb{P}\left(\frac{\sqrt{NT}(\tilde{\alpha}_j - \alpha_j)}{\tilde{\sigma}_{z,j}} \leq -z_{1-\delta/2}\right). \end{aligned}$$

Thus, taking the supremum over  $\gamma \in \mathcal{F}(\nu, E)$  and letting  $N$  tend to infinity yields (4.2) via (4.1). The proof is the same for (4.3). Next, we turn to (4.4).

$$\begin{aligned} & \sqrt{NT} \sup_{\gamma \in \mathcal{F}(\nu, E)} \text{diam}\left(\left[\tilde{\alpha}_j - z_{1-\delta/2} \frac{\tilde{\sigma}_{\alpha,j}}{\sqrt{NT}}, \tilde{\alpha}_j + z_{1-\delta/2} \frac{\tilde{\sigma}_{\alpha,j}}{\sqrt{NT}}\right]\right) \\ & = 2z_{1-\delta/2} \left(\sqrt{[\Theta_Z \Sigma_{1,N} \Theta_Z]_{jj}} + o_p(1)\right) \leq 2z_{1-\delta/2} \left(\frac{\sqrt{\text{maxeval}(\Sigma_{1,N})}}{\text{mineval}(\Psi_Z)} + o_p(1)\right) = O_p(1), \end{aligned}$$



where the first equality is due to (3.8), and the last equality is due to Assumptions 4(a) and 6(d). Similarly, we can prove (4.5):

$$\begin{aligned} \sqrt{T} \sup_{\gamma \in \mathcal{F}(\nu, E)} \text{diam} \left( \left[ \tilde{\eta}_i - z_{1-\delta/2} \frac{\tilde{\sigma}_{\eta,i}}{\sqrt{T}}, \tilde{\eta}_i + z_{1-\delta/2} \frac{\tilde{\sigma}_{\eta,i}}{\sqrt{T}} \right] \right) &= 2z_{1-\delta/2} \left( \sqrt{[\Sigma_{3,N}]_{ii}} + o_p(1) \right) \\ &= 2z_{1-\delta/2} \left( \left[ \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\varepsilon_{i,t}^2] \right]^{1/2} + o_p(1) \right) = O_p(1), \end{aligned}$$

where the third equality follows from the arguments above (10.26).  $\square$

## 12 Appendix F

In this appendix we prove some auxiliary results used throughout the previous appendices.

**Proposition 1.** *Let  $A$  and  $B$  be two positive semidefinite  $(p-1) \times (p-1)$  matrices and  $\delta := \max_{1 \leq l, k \leq p-1} |A_{lk} - B_{lk}|$ . For any integer  $r \in \{1, \dots, p-1\}$ , one has*

$$\kappa^2(B, r) \geq \kappa^2(A, r) - \delta 16r.$$

*Proof.* The proof is exactly the same as that of Lemma 6.  $\square$

**Theorem 4 (Fan et al. (2012)).** *Let  $\alpha \in (0, 1)$ . Assume that  $(X_i, \mathcal{F}_i)_{i=1}^n$  is a sequence of supermartingale differences satisfying  $\sup_i \mathbb{E}[e^{|X_i|^{\frac{2\alpha}{1-\alpha}}}] \leq C_1$  for some constant  $C_1 \in (0, \infty)$ . Define  $S_k := \sum_{i=1}^k X_i$ . Then, for all  $\epsilon > 0$ ,*

$$\mathbb{P} \left( \max_{1 \leq k \leq n} S_k \geq n\epsilon \right) \leq C(\alpha, n, \epsilon) e^{-(\epsilon/4)^{2\alpha} n^\alpha},$$

where

$$C(\alpha, n, \epsilon) := 2 + 35C_1 \left[ \frac{1}{16^{1-\alpha}(n\epsilon^2)^\alpha} + \frac{1}{n\epsilon^2} \left( \frac{3(1-\alpha)}{2\alpha} \right)^{\frac{1-\alpha}{\alpha}} \right].$$

The preceding theorem is not exactly the same as Theorem 2.1 in Fan et al. (2012), but taken from the proof of Theorem 2.1 in Fan et al. (2012). This theorem generalises Theorem 3.2 in Lesigne and Volny (2001).

**Proposition 2.** *Let  $\alpha \in (0, 1)$ . Assume that  $(X_i, \mathcal{F}_i)_{i=1}^n$  is a sequence of martingale differences satisfying  $\sup_i \mathbb{E}[e^{D|X_i|^{\frac{2\alpha}{1-\alpha}}}] \leq C_1$  for some positive constant  $D$ . ( $C_1$  could change with the sample size  $n$ .) Then, for all  $\epsilon \gtrsim \frac{1}{\sqrt{n}}$ ,*

$$\mathbb{P} \left( \left| \sum_{i=1}^n X_i \right| \geq n\epsilon \right) \leq AC_1 e^{-K(\epsilon^2 n)^\alpha},$$

for positive constants  $A$  and  $K$ .

*Proof.* This proposition is a simple adaptation of the preceding theorem. Note that for some positive constant  $D$ ,

$$\mathbb{P}\left(\sum_{i=1}^n X_i \geq n\epsilon\right) = \mathbb{P}\left(\sum_{i=1}^n D^{\frac{1-\alpha}{2\alpha}} X_i \geq nD^{\frac{1-\alpha}{2\alpha}} \epsilon\right) = \mathbb{P}\left(\sum_{i=1}^n Y_i \geq n\delta\right),$$

where  $Y_i := D^{\frac{1-\alpha}{2\alpha}} X_i$  and  $\delta := D^{\frac{1-\alpha}{2\alpha}} \epsilon$ . Now  $(Y_i)_{i=1}^n$  is a sequence of martingale differences satisfying  $\sup_i \mathbb{E}[e^{|Y_i|^{\frac{2\alpha}{1-\alpha}}}] \leq C_1$ . Invoking the preceding theorem, we have

$$\mathbb{P}\left(\sum_{i=1}^n Y_i \geq n\delta\right) \leq C(\alpha, n, \delta) e^{-(\delta/4)^{2\alpha} n^\alpha}.$$

$(-Y_i)_{i=1}^n$  is also a sequence of martingale differences satisfying the same exponential moment condition. Thus,

$$\begin{aligned} \mathbb{P}\left(\left|\sum_{i=1}^n X_i\right| \geq n\epsilon\right) &= \mathbb{P}\left(\left|\sum_{i=1}^n Y_i\right| \geq n\delta\right) \leq 2C(\alpha, n, \delta) e^{-(\delta/4)^{2\alpha} n^\alpha} \\ &= 2C(\alpha, n, D^{\frac{1-\alpha}{2\alpha}} \epsilon) e^{-(D^{\frac{1-\alpha}{2\alpha}} \epsilon/4)^{2\alpha} n^\alpha} \leq AC_1 e^{-K\epsilon^{2\alpha} n^\alpha}, \end{aligned}$$

for positive constants  $A, K$ , where the last inequality used that if  $\epsilon \gtrsim \frac{1}{\sqrt{n}}$  then  $2C(\alpha, n, D^{\frac{1-\alpha}{2\alpha}} \epsilon) \leq AC_1$  for some positive constant  $A$ .  $\square$

**Proposition 3.** Suppose we have random variables  $Z_{l,i,t,j}$  uniformly subgaussian for  $l = 1, \dots, L$  ( $L \geq 2$  fixed),  $i = 1, \dots, N$ ,  $t = 1, \dots, T$  and  $j = 1, \dots, p$ .  $Z_{l_1, i_1, t_1, j_1}$  and  $Z_{l_2, i_2, t_2, j_2}$  are independent as long as  $i_1 \neq i_2$  regardless of the values of other subscripts. Then,

$$\max_{1 \leq j \leq p} \max_{1 \leq t \leq T} \max_{1 \leq i \leq N} \mathbb{E} \left| \prod_{l=1}^L Z_{l,i,t,j} \right| \leq A = O(1), \quad (12.1)$$

for some positive constant  $A$  and

$$\max_{1 \leq j \leq p} \left| \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \left( \prod_{l=1}^L Z_{l,i,t,j} - \mathbb{E} \left[ \prod_{l=1}^L Z_{l,i,t,j} \right] \right) \right| = O_p \left( \sqrt{\frac{(\log(pT))^{L+1}}{N}} \right). \quad (12.2)$$

*Proof.* For every  $\epsilon \geq 0$ ,  $\mathbb{P}(|\prod_{l=1}^L Z_{l,i,t,j}| \geq \epsilon) \leq \sum_{l=1}^L \mathbb{P}(|Z_{l,i,t,j}| \geq \epsilon^{1/L}) \leq L \frac{K}{2} e^{-C\epsilon^{2/L}}$  for positive constants  $K, C$ . Next, using Hölder's inequality, we have

$$\max_{1 \leq j \leq p} \max_{1 \leq t \leq T} \max_{1 \leq i \leq N} \mathbb{E} \left| \prod_{l=1}^L Z_{l,i,t,j} \right| \leq \max_{1 \leq j \leq p} \max_{1 \leq t \leq T} \max_{1 \leq i \leq N} \prod_{l=1}^L \left( \mathbb{E} |Z_{l,i,t,j}|^L \right)^{\frac{1}{L}}.$$

Uniform subgaussianity implies that  $\left( \mathbb{E} |Z_{l,i,t,j}|^L \right)^{\frac{1}{L}}$  is uniformly bounded. That is,  $\left( \mathbb{E} |Z_{l,i,t,j}|^L \right)^{\frac{1}{L}} \leq L! \|Z_{l,i,t,j}\|_{\psi_1} \leq L! (\log 2)^{-1/2} \|Z_{l,i,t,j}\|_{\psi_2} \leq L! (\log 2)^{-1/2} \left( \frac{1+K/2}{C} \right)^{1/2}$ , where the first two inequalities are taken from p95 of van der Vaart and Wellner (1996), and the third inequality is due to Lemma 2.2.1 in van der Vaart and Wellner (1996). (12.1) then follows.

For every  $\epsilon \geq 0$ ,

$$\begin{aligned} \mathbb{P} \left( \left| \frac{1}{T} \sum_{t=1}^T \left( \prod_{l=1}^L Z_{l,i,t,j} - \mathbb{E} \left[ \prod_{l=1}^L Z_{l,i,t,j} \right] \right) \right| \geq \epsilon \right) &\leq \mathbb{P} \left( \max_{1 \leq t \leq T} \left| \prod_{l=1}^L Z_{l,i,t,j} \right| \geq \epsilon - A \right) \\ &\leq \sum_{t=1}^T \mathbb{P} \left( \left| \prod_{l=1}^L Z_{l,i,t,j} \right| \geq \epsilon - A \wedge \epsilon \right) \leq \frac{L}{2} T K e^{-C(\epsilon - A \wedge \epsilon)^{2/L}} \leq \frac{L}{2} T K e^{-C[\epsilon^{2/L} - (A \wedge \epsilon)^{2/L}]} \leq T K' e^{-C\epsilon^{2/L}}, \end{aligned}$$

for  $K' = \frac{L}{2} K e^{CA^{2/L}}$  and where the second last inequality is due to subadditivity of the concave function:  $(x+y)^{2/L} \leq x^{2/L} + y^{2/L}$  for  $x, y \geq 0$ ,  $L \geq 3$ . (The inequality  $(x+y)^{2/L} \leq x^{2/L} + y^{2/L}$  for  $L = 2$  is trivial.) Let  $X_{i,j}$  denote  $\frac{1}{T} \sum_{t=1}^T \left( \prod_{l=1}^L Z_{l,i,t,j} - \mathbb{E} \left[ \prod_{l=1}^L Z_{l,i,t,j} \right] \right)$ . Consider some positive constant  $D < C$ .

$$\begin{aligned} \mathbb{E} \left[ e^{D|X_{i,j}|^{2/L}} \right] &= \int_{x \in \mathbb{R}} \int_0^{|x|^{2/L}} D e^{Ds} ds P(dx) + 1 = \int_0^\infty D e^{Ds} \mathbb{P}(|X_{i,j}| > s^{L/2}) ds + 1 \\ &\leq \int_0^\infty T K' D e^{(D-C)s} ds + 1 = \frac{T K' D}{C-D} + 1 \leq B T, \end{aligned} \tag{12.3}$$

for some positive constant  $B$ , where the second equality is by Fubini's theorem. Then we can use independence across  $i$  to invoke Proposition 2 in Appendix F with  $\alpha = \frac{1}{L+1}$  and  $C_1 = B T$ , for  $\epsilon \gtrsim \frac{1}{\sqrt{N}}$ ,

$$\mathbb{P} \left( \left| \sum_{i=1}^N \frac{1}{T} \sum_{t=1}^T \left( \prod_{l=1}^L Z_{l,i,t,j} - \mathbb{E} \left[ \prod_{l=1}^L Z_{l,i,t,j} \right] \right) \right| \geq N \epsilon \right) \leq A' T e^{-K(\epsilon^2 N)^{\frac{1}{L+1}}}$$

for positive constants  $A'$  and  $K$ . Setting  $\epsilon = \sqrt{\frac{M(\log(pT))^{L+1}}{N}} \left( \gtrsim \frac{1}{\sqrt{N}} \right)$  for some  $M > 0$ , we have

$$\mathbb{P} \left( \max_{1 \leq j \leq p} \left| \sum_{i=1}^N \frac{1}{T} \sum_{t=1}^T \left( \prod_{l=1}^L Z_{l,i,t,j} - \mathbb{E} \left[ \prod_{l=1}^L Z_{l,i,t,j} \right] \right) \right| \geq N \epsilon \right) \leq p A' T e^{-K(\epsilon^2 N)^{\frac{1}{L+1}}} = A' (pT)^{1-KM^{\frac{1}{L+1}}}.$$

The upper bound of the preceding probability becomes arbitrarily small for  $T$  and  $M$  sufficiently large. Hence (12.2) follows.  $\square$

**Lemma 10.** *Let  $A$  be a symmetric  $p \times p$  matrix, and  $\hat{v}$  and  $v \in \mathbb{R}^p$ . Then*

$$|\hat{v}' A \hat{v} - v' A v| \leq \|A\|_\infty \|\hat{v} - v\|_1^2 + 2 \|A v\| \|\hat{v} - v\|.$$

*Proof.* See Lemma 6.1 in the working-paper version of van de Geer et al. (2014).  $\square$

**Theorem 5 (McLeish (1974)).** *Let  $\{X_{n,i}, i = 1, \dots, k_n\}$  be a martingale difference array with respect to the triangular array of  $\sigma$ -algebras  $\{\mathcal{F}_{n,i}, i = 0, \dots, k_n\}$  (i.e.,  $X_{n,i}$  is  $\mathcal{F}_{n,i}$ -measurable and  $\mathbb{E}[X_{n,i} | \mathcal{F}_{n,i-1}] = 0$  almost surely for all  $n$  and  $i$ ) satisfying  $\mathcal{F}_{n,i-1} \subseteq \mathcal{F}_{n,i}$  for all  $n \geq 1$ . Assume,*

(i)  $\max_{i \leq k_n} |X_{n,i}|$  is uniformly bounded in  $L_2$  norm,

(ii)  $\max_{i \leq k_n} |X_{n,i}| \xrightarrow{p} 0$ , and

(iii)  $\sum_{i=1}^{k_n} X_{n,i}^2 \xrightarrow{p} 1$ .

Then,  $S_n = \sum_{i=1}^{k_n} X_{n,i} \xrightarrow{d} N(0, 1)$  as  $n \rightarrow \infty$ .

## 13 Appendix G

### 13.1 Construction of $\hat{\Theta}$

In this subsection we show how  $\hat{\Theta}_Z$  is constructed by nodewise regressions. First, define

$$\hat{\phi}_j = \underset{\delta \in \mathbb{R}^{p-1}}{\operatorname{argmin}} \left\{ \frac{1}{NT} \|z_j - Z_{-j}\delta\|^2 + 2\lambda_{node} \|\delta\|_1 \right\}, \quad j = 1, \dots, p, \quad (13.1)$$

where  $z_j$  is the  $j$ th column of  $Z$ ,  $Z_{-j}$  is the  $NT \times (p-1)$  submatrix of  $Z$  with  $Z$ 's  $j$ th column removed, and the  $(p-1) \times 1$  vector  $\hat{\phi}_j = \{\hat{\phi}_{j,l} : l = 1, \dots, p, l \neq j\}$ . Thus,  $\hat{\phi}_j$  is the Lasso estimator resulting from regressing  $z_j$  on  $Z_{-j}$ . Next, define

$$\hat{C} = \begin{pmatrix} 1 & -\hat{\phi}_{1,2} & \cdots & -\hat{\phi}_{1,p} \\ -\hat{\phi}_{2,1} & 1 & \cdots & -\hat{\phi}_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ -\hat{\phi}_{p,1} & -\hat{\phi}_{p,2} & \cdots & 1 \end{pmatrix}$$

and  $\hat{\tau}_j^2 = \frac{1}{NT} \|z_j - Z_{-j}\hat{\phi}_j\|^2 + \lambda_{node} \|\hat{\phi}_j\|_1$  as well as  $\hat{T}^2 = \operatorname{diag}(\hat{\tau}_1^2, \dots, \hat{\tau}_p^2)$ . Finally, we set  $\hat{\Theta}_Z = \hat{T}^{-2} \hat{C}$ . Let  $\hat{C}_j$  denote the  $j$ th row of  $\hat{C}$  and let  $\hat{\Theta}_{Z,j}$  denote the  $j$ th row of  $\hat{\Theta}_Z$  but both written as a  $p \times 1$  vectors. Then,  $\hat{\Theta}_{Z,j} = \hat{C}_j / \hat{\tau}_j^2$ . For any  $j = 1, \dots, p$ , the KKT condition for a minimum in (13.1) are

$$-\frac{1}{NT} Z'_{-j} (z_j - Z_{-j} \hat{\phi}_j) + \lambda_{node} w_j = 0, \quad (13.2)$$

where  $w_j$  is the subdifferential of  $\|x\|_1$  evaluated at  $\hat{\phi}_j$ . Using this, the definition of  $\hat{\tau}_j$ , and  $\hat{\phi}_j' w_j = \|\hat{\phi}_j\|_1$  yields

$$\hat{\tau}_j^2 = \frac{1}{NT} (z_j - Z_{-j} \hat{\phi}_j)' (z_j - Z_{-j} \hat{\phi}_j) + \lambda_{node} \|\hat{\phi}_j\|_1 = \frac{1}{NT} (z_j - Z_{-j} \hat{\phi}_j)' z_j. \quad (13.3)$$

Thus, by the definition of  $\hat{\Theta}_{Z,j}$ , and as  $\hat{\tau}_j^2$  is bounded away from zero (we shall later argue rigorously for this)

$$\frac{1}{NT} z_j' Z \hat{\Theta}_{Z,j} = 1. \quad (13.4)$$

Furthermore, the KKT conditions (13.2) can also be written as

$$\frac{1}{NT} Z'_{-j} (z_j - Z_{-j} \hat{\phi}_j) = \lambda_{node} w_j, \quad (13.5)$$

which implies  $\frac{1}{NT}Z'_{-j}Z\hat{\Theta}_{Z,j} = \lambda_{node}w_j/\hat{\tau}_j^2$ . Combining with (13.4) yields

$$\left\| \frac{1}{NT}Z'Z\hat{\Theta}_{Z,j} - e_j \right\|_{\infty} \leq \frac{\lambda_{node}}{\hat{\tau}_j^2}, \quad (13.6)$$

where  $e_j$  is the  $j$ th basis vector of  $\mathbb{R}^p$ . Together with an oracle inequality for  $\|\hat{\gamma} - \gamma\|_1$ , (13.6) provides an upper bound on the  $j$ th entry of  $\Delta$  in (3.2).