# NETWORK REGRESSION AND SUPERVISED CENTRALITY ESTIMATION

JUNHUI CAI University of Pennsylvania

DAN YANG The University of Hong Kong

## WU ZHU Tsinghua University

### HAIPENG SHEN The University of Hong Kong

## LINDA ZHAO University of Pennsylvania

The centrality in a network is a popular metric for agents' network positions and is often used in regression models to model the network effect on an outcome variable of interest. In empirical studies, researchers often adopt a two-stage procedure to first estimate the centrality and then infer the network effect using the estimated centrality. Despite its prevalent adoption, this two-stage procedure lacks theoretical backing and can fail in both estimation and inference. We, therefore, propose a unified framework, under which we prove the shortcomings of the two-stage in centrality estimation and the undesirable consequences in the regression. We then propose a novel supervised network centrality estimation (SuperCENT) methodology that simultaneously yields superior estimations of the centrality and the network effect and provides valid and narrower confidence intervals than those from the two-stage. We showcase the superiority of SuperCENT in predicting the currency risk premium based on the global trade network.

KEYWORDS: Hub and authority centrality, Network regression inference, Measurement error, Hyperlink Induced Topic Search (HITS) Algorithm, Global trade network, Currency risk premium.

Cai: Department of Statistics and Data Science, The Wharton School, University of Pennsylvania, Philadelphia, PA 19104, USA (email: junhui@wharton.upenn.edu); Yang: Innovation and Information Management, Faculty of Business and Economics, The University of Hong Kong, Hong Kong, China (email: dyanghku@hku.hk); Zhu: Department of Finance, School of Economics and Management, Tsinghua University, Beijing, 100084, China (email: zhuwu@sem.tsinghua.edu.cn); Shen: Innovation and Information Management, Faculty of Business and Economics, The University of Hong Kong, China (email: haipeng@hku.hk); Zhao: Department of Statistics and Data Science, The Wharton School, University of Pennsylvania, Philadelphia, PA 19104, USA (email: Izhao@wharton.upenn.edu).

### 1. INTRODUCTION

In economics, finance, operations, sociology, and many other disciplines, agents (individual persons, firms, industries, and countries, etc.) are usually connected through different relationships. The collection of the agents and their relationships is usually represented by a network. One extremely influential concept of networks is the agents' positions, because they can induce a wide range of behaviors, including individuals' decisions on education and human capital investment (Jackson et al., 2017), information sharing and advertisement (Banerjee et al., 2019, Breza and Chandrasekhar, 2019), coalition for exchange or cooperation (Elliott and Golub, 2019), firms' investment decision-making (Allen et al., 2019), the identification of banks that are too-connected-to-fail (Gofman, 2017), and stock returns (Ahern, 2013, Richmond, 2019).

An agent's position or importance is usually captured by the so-called *centrality*, which measures how *central* the agent is in comparison to the others in a network. Since the position induces the agent's behavior and thus shapes certain outcomes, the network centrality is often used as an intermediary to measure the network effects on the outcome of interest (Ahern, 2013, Shao et al., 2018, Richmond, 2019, Allen et al., 2019, Banerjee et al., 2019, Bovet and Makse, 2019). There are many kinds of definitions of centrality (Jackson, 2010, Kolaczyk, 2010), among which we focus on the *hub* and *authority centralities* (Kleinberg, 1999), of which the *eigenvector* centrality is a special case. See a brief overview of centralities and concrete examples on the implication of centralities in Section 2.

The value of network centralities is, therefore, two-fold: first, the centralities have natural implications on the importance of agents; second, the centralities are often used as regressors in a regression to model the *network effects* on some outcome of interest. In practice, the centralities are not directly observable while the network is. Hence, the two-fold value of centralities leads to two goals of this article:

(G1) Estimate centralities from an observed network;

(G2) Estimate and conduct inference of the network effects through the centralities.

In empirical studies including many of the above-cited, these two goals are usually achieved in sequential order, to which we refer as the *two-stage* procedure throughout this article. Stage 1 targets goal (G1) solely, which estimates centralities by performing the singular value decomposition (SVD) on the adjacency matrix that represents the network. Stage 2 aims at goal (G2) next, which estimates the network effects by regressing the outcome on the *estimated* centralities from Stage 1 and conducts inference using the naive confidence intervals from the regression, an *ad-hoc* inference, ignoring the centrality estimation error. The drawbacks of such two-stage procedure are that: Stage 1 only uses the information from the network to estimate centralities, without resorting to the auxiliary information from the regression on the centralities, which leads to an inaccurate estimation of the centralities due to large observational errors in the network (see more evidence in Section 2); Stage 2 is contingent on Stage 1 – regressing the outcome on the inaccurately estimated centralities further results in an inaccurate estimation of the regression coefficients, thereby invalidating the inference. Consequently, these two stages do not achieve the two goals at their best.

Motivated by the two goals and the shortcomings of the two-stage procedure, we first propose a *unified* framework that encapsulates two models for the corresponding two goals: one network generation model based on centralities for goal (G1) and one network regression model for the dependency of the outcome on the centralities for goal (G2). We further propose a novel *supervised network centrality estimation (SuperCENT)* methodology that accomplishes both (G1) and (G2) *simultaneously*, instead of sequentially. SuperCENT exploits information from the two models – the network regression model contains auxiliary information on the centrality in addition to the network, providing *supervision* to the centrality estimation. The supervision effect can improve the centrality estimation, which in turn benefits the network regression. In

other words, the centrality estimation and the network regression complement and empower each other.

Under the unified framework, we derive the theoretical convergence rates of the centralities and the regression coefficients estimators via the two-stage and SuperCENT, as well as their asymptotic distributions, which can be used to construct confidence intervals. Comparing the two methods, SuperCENT universally dominates the two-stage theoretically and empirically in terms of centrality estimation (G1) as well as the estimation and inference of the network regression coefficients (G2). We summarize our contributions of this article in the following.

- 1. To the best of our knowledge, we are the first to provide a unified framework to study properties of centrality estimation, centrality inference, and the subsequent network regression analysis when a noisy network is observed.
- 2. We demonstrate that the common practice of two-stage can be problematic. For centrality estimation (G1), the two-stage centrality estimates using SVD in Stage 1 are inconsistent under large noise in the network. The same inconsistency phenomenon appears when we estimate the true underlying network. This finding of inconsistency extends the phase transition phenomenon of the singular vectors (Shabalin and Nobel, 2013) and eigenvectors (Shen et al., 2016) to the network centrality problem. For the network regression (G2), the centrality coefficients estimates are *biased* and *inconsistent* given the inconsistent centrality estimates from Stage 1 and the ad-hoc inference can be either *conservative* or *invalid* depending on the network noise level.
- 3. We show theoretically and empirically that the proposed SuperCENT dominates the twostage universally. For (G1), SuperCENT yields superior estimations of both the centralities and true network over the two-stage. For (G2), SuperCENT can mitigate the coefficients estimation bias and thus boost the estimation accuracy under large network noise thanks to the superior centrality estimation. In addition, SuperCENT provides confidence intervals that are *valid* and *narrower* than the ad-hoc two-stage confidence intervals.
- 4. Lastly, we apply SuperCENT and the two-stage to predict the currency risk premium, based on an economic theory on the relationship between a country's currency risk premium and its importance within the global trade network. We show that a long-short trading strategy based on the SuperCENT centrality estimates yields a return three times as high as those by the two-stage procedure. Furthermore, SuperCENT can verify the economic theory via a rigorous statistical test while the two-stage fails.

As a concrete manifestation of our contributions, we perform a toy experiment as follows. The network is constructed by perturbing a true network with noise of different level  $\sigma_a$ . The regression model includes both the hub and authority centralities, u and v, with their corresponding coefficients  $\beta_u$  and  $\beta_v$ , as well as other covariates<sup>1</sup>. Figure 1 shows the performance of the two-stage procedure and SuperCENT in terms of estimation of u and  $\beta_u$  as well as the coverage probability and the width of confidence interval (CI) for the hub centrality coefficient  $CI_{\beta_u}$  with varying network noise levels  $\sigma_a$ . Figure 1A shows the error of the hub centrality estimation. As the noise level increases, the two-stage estimate becomes increasingly inaccurate and eventually orthogonal to the truth, corroborating our inconsistency finding of the two-stage under large network noise. On the contrary, SuperCENT estimates u very accurately until the network noise becomes very large. Even then, the accuracy of the SuperCENT estimate does

<sup>&</sup>lt;sup>1</sup>The detail of the network model is:  $\mathbf{A} = \mathbf{u}\mathbf{v}^{\top} + \mathbf{E} \in \mathbb{R}^{256 \times 256}$ , where the true hub centrality is  $\mathbf{u}$ , the true authority is  $\mathbf{v}$ , both of which are scaled to have norm  $\sqrt{256}$ , and  $\mathbf{E}$  contains i.i.d. normal random variables with mean zero and variance  $\sigma_a^2$ . For the regression model:  $\mathbf{y} = \mathbf{X}\beta_x + 16\mathbf{u} + \mathbf{v} + \mathbf{\epsilon} \in \mathbb{R}^{256 \times 1}$  where  $\beta_x = (1,3,5)^{\top}$ , the covariate matrix  $\mathbf{X}$  consists of a column of 1's and two columns whose entries follow N(0,1) independently, and  $\mathbf{\epsilon}$  follows  $N(\mathbf{0}, \sigma_y^2 \mathbf{I}_{256})$  with  $\sigma_y^2 = 2^{-4}$ . We vary the network noise level  $\sigma_a \in 2^{1,1.5,\ldots,5}$  and each configuration is repeated 500 times.





FIGURE 1.—Comparison between the two-stage and SuperCENT in terms of the estimations of u and  $\beta_u$  as well as the coverage probability and the width of  $CI_{\beta_u}$  varying the network noise level  $\sigma_a$ . The performance of the two-stage is shown in the red dashed line and SuperCENT in the green solid line. Subfigure A shows the estimation error of the hub centrality  $\sin(\angle(\hat{u}, u))$ , i.e., sine of the angle between the true hub centrality u and the estimate  $\hat{u}$ . As the noise level increases, the two-stage hub centrality estimate using SVD becomes increasingly inaccurate and eventually orthogonal to the truth. SuperCENT estimates the centrality very accurately until  $\sigma_a$  becomes really large. Even then, the accuracy of the SuperCENT estimate does not deteriorate to zero thanks to the auxiliary information from the regression. Subfigure B shows the bias in estimation of  $\beta_u$ , i.e.,  $\hat{\beta}_u - \beta_u$ . The inconsistency of the centrality estimation bias. Subfigures C and D show the coverage probability and the width of the 95% confidence interval of the hub centrality coefficient,  $CI_{\beta_u}$ , respectively. The ad-hoc two-stage confidence interval is either conservative or invalid depending on the network noise level  $\sigma_a$  while SuperCENT provides a valid and narrower confidence interval until  $\sigma_a$  becomes unreasonably large.

not deteriorate to zero thanks to the auxiliary information from the regression. Figure 1B exhibits the estimation bias of  $\hat{\beta}_u$ : the two-stage estimate suffers from a severe attenuation bias while SuperCENT alleviates the bias. Figures 1C-D illustrate the coverage probability and the width of the confidence interval of  $\beta_u$ . Depending on the noise level in the network model, the ad-hoc two-stage inference has different undesirable consequences: when the noise is small, the ad-hoc two-stage confidence interval is still valid but conservative; when the noise is large, the ad-hoc two-stage confidence interval is invalid and wider than necessary. In contrast, SuperCENT provides a valid and narrower confidence interval until the network noise becomes unreasonably large.

Our paper contributes to several lines of inquiry in the network and econometrics literature on network modeling, network regression with centralities, covariate-assisted network modeling, network effect modeling, and measurement error. First, the proposed unified framework unites the literature on the noisy network and network regression with centralities. Most existing network literature focuses on only one of the two aspects in our unified framework. On one

hand, in the presence of noisy network generation, there are many empirical studies (Lakhina et al., 2003, Banerjee et al., 2013, Breza et al., 2020, Zhu and Yang, 2020) and many that try to estimate or recover the true network without involving centrality (Butts, 2003, Handcock and Gile, 2010, Chandrasekhar and Lewis, 2011, Le et al., 2018, Newman, 2018, Rohe, 2019). On the other hand, many researchers focus on the network regression model with centralities while ignoring the noise of the centrality estimation inherited from the noise of the network (Ahern, 2013, Hochberg et al., 2007, Shao et al., 2018, Richmond, 2019, Allen et al., 2019, Liu, 2019, Banerjee et al., 2019, Bovet and Makse, 2019, Fogli and Veldkamp, 2021).

Our unified framework is also connected to the line of research related to network with covariates supervision (Newman and Clauset, 2016, Zhang et al., 2016, Li et al., 2016, Graham, 2017, Binkiewicz et al., 2017, Yan et al., 2019, Ma et al., 2020, Chen et al., 2021). One major difference is that SuperCENT uses both the covariates and the response to supervise, instead of only the covariates. And they focus mostly on network formation or community detection.

In econometrics, there has been a significant effort to model the network effect on an outcome of interest through regression. See De Paula (2017) for a review on the econometrics of network models. One popular approach follows the pioneer work of Manski (1993), the "reflection model" (Lee, 2007, Bramoullé et al., 2009, Lee et al., 2010, Hsieh and Lee, 2016, Zhu et al., 2017). This approach handles the network effect through the observed adjacency matrix itself, not through centralities like ours. There is also a recent surge of literature in network recovery based on the reflection model (De Paula et al., 2019, Battaglini et al., 2021). Literature on this approach mainly focuses on the issue of identifiability, while our work attends to both the estimation and inference of the network effect. Nevertheless, it is possible to incorporate our model into the reflection model by assuming a low-rank structure on the true underlying network. Another popular approach assumes that the outcome depends on individual fixed effects, which are only estimable by imposing constraints or penalties, and the role of the network is cast through the Laplacian matrix, such that connected nodes share similar individual fixed effects (Jochmans and Weidner, 2019, Li et al., 2019, Le and Li, 2020). This approach emphasizes the network homophily, while ours concentrates on nodes' position or importance in the network using centralities.

Lastly, our methodology further contributes to the measurement error literature. Most literature concerns a regression setup where the covariates are observed with errors, which leads to bias in the coefficient estimation (Garber and Klepper, 1980, Griliches, 1986, Pischke, 2007, Wooldridge, 2015, Abel, 2017). We extend it to the network regression problem. Specifically, the two-stage procedure resembles the measurement error problem where the estimated centralities that are used as covariates in the regression of Stage 2 contain estimation error. Nevertheless, the derivation of the two-stage bias is not a trivial extension of the classical results because it involves the asymptotic joint distributions of the two-stage centrality estimators. Furthermore, SuperCENT corrects the bias problem in the regression coefficient estimation induced by the estimation errors and provides valid inference for the regression coefficients.

The rest of this article is organized as follows. Section 2 reviews the concept and properties of networks, network centralities, noisy networks, and provides concrete examples of network centralities and their effects on the outcome variables. Section 3 formally introduces the unified framework and makes connections with the existing work. Details of the proposed SuperCENT methodology are provided in Section 4. Theoretical properties are studied in Section 5 and the simulation study is demonstrated in Section 6. Section 7 presents the case study of the global trade network centralities and their relationships with risk premiums. Section 8 concludes with a summary and future work. Some concrete mathematical expressions, a special case of an undirected network with the eigenvector centrality, more simulation results, additional information of the case study, and the proofs are delegated to the supplementary materials. A new R package called SuperCENT implements the methods (https://jh-cai.com/SuperCENT).

### 2. NETWORK AND NETWORK CENTRALITY

In this section, we provide background knowledge and a literature review on networks and network centralities. We review how network centralities are used in various fields, particularly via the above-mentioned two-stage procedure, and point out the challenge in centralities estimation due to observational errors in networks. Sophisticated readers may skip this section.

In a network, the nodes are agents involved in a network of relationships and the edges represent the relationships between the nodes. The edges can be directed or undirected depending on whether the relationships are reciprocal. For a directed network  $\mathcal{G}$  composed of n nodes V and a set of weighted directed edges  $E \subseteq V \times V$ , it can be represented by an *asymmetric* adjacency matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  where  $a_{ij}$  is the weighted edge from node i to j.

A full description of the network depends on all the nodes and edges, whose information is too much to be thoroughly understood in empirical analysis. To feasibly analyze the network, researchers usually resort to dimension reduction or low-rank approximation to extract the characteristics and structural properties. Centralities are common low-rank summaries of the network information. Depending on goals and domain knowledge, researchers have used multiple versions of centralities. We refer readers to Chapter 2 of Jackson (2010), Chapter 4 of Kolaczyk (2010) or Chapter 1 of Graham and De Paula (2020) for a comprehensive introduction to centralities.

In this article, we focus on the *hub* and *authority centralities* for directed networks (Kleinberg, 1999, Benzi et al., 2013), which can be reduced to eigenvector centrality for undirected networks. For directed networks, there is a distinction between the giver and the recipient, such as the citee-citor in citation networks or webpage networks, the exporter-importer in trade networks, and the investor-investee in investment networks. The concept of "hubs and authorities" originated from web searching. Intuitively, the hub centrality of a web page depends on the total level of authority centrality of the web pages it links to, while the authority centrality of a web page depends on the total level of hub centrality of the web pages it receives links from. Similar intuition can be applied to citation networks where the hub centrality of a paper reveals the quality of a survey paper while an authoritative paper is one that is cited a lot by well-respected survey papers.

Let us use a toy example to further illuminate the intuition behind the hub and authority centralities. Consider a citation network where each paper is a node and an edge from Paper A to Paper B indicates Paper A cites Paper B. Figure 2a shows an example of the adjacency matrix of such network. Figures 2b and 2c show the same network with different node sizes: the node sizes in Figure 2b are proportional to the hub centralities while those in Figure 2c are proportional to the authority centralities.

To understand the hub centrality, note that Papers 1 and 4 are the major citors: they both cite three papers with Paper 2 being the common one. Except for the common one, Paper 4 cites Papers 5 and 6, which are only cited by Paper 4, and Paper 1 cites Papers 4 and 3, among which Paper 3 is cited twice. Therefore, compared with Paper 4, Paper 1 cites the same number of papers with one being cited more than the others. This makes the hub centrality of Paper 1 larger than that of Paper 4. One can think of Paper 1 as a better survey paper than Paper 4. Paper 7 cites only one paper, which makes its hub centrality smaller than Papers 1 and 4. The rest of the papers have small hub centrality since they do not cite other papers. As for the authority centrality, attention should be given to citees. Papers 2 and 3 both have two citations, but Paper 2 is cited by Papers 1 and 4 while Paper 3 is cited by Papers 1 and 7. Observe that Paper 4 as a hub is more influential than Paper 7. So the authority centrality for Paper 2 is the highest, followed by Paper 3. For the same reason, Paper 4 has higher authority centrality than Papers 5 and 6, since Paper 4 is cited by Paper 1 while Papers 5 and 6 are cited by Paper 4.



(a) The adjacency matrix of the citation network. For readability, we denote 0 as  $\cdot$ .

(b) Node size by hub centralities.

(c) Node size by authority centralities.

FIGURE 2.—A toy network to illustrate the hub and authority centrality.

One can obtain the hub and authority centralities using an iterative method. Let  $u_i$  denote the hub centrality and  $v_i$  denote the authority centrality for node *i*. Kleinberg (1999) proposes to initiate each  $u_i$  and  $v_i$  with certain nonzero value and then iteratively update as follows:

$$\boldsymbol{v}^{(k)} \leftarrow \mathbf{A}^{\top} \boldsymbol{u}^{(k-1)}$$
 and  $\boldsymbol{u}^{(k)} \leftarrow \mathbf{A} \boldsymbol{v}^{(k)}$  for  $k = 1, 2, 3, \dots$  (1)

This iterative algorithm is shown to converge under some regularity conditions with proper normalization; and the hub and authority centralities are the final  $u^{(\infty)}$  and  $v^{(\infty)}$  after convergence. This iterative algorithm is also well known as the power method to compute the leading left and right singular vectors of A (Van Loan and Golub, 1996). Due to the equivalence between singular value decomposition (SVD) and eigen-decomposition, the hub centrality is the leading left singular vector of A and the leading eigen-vector of  $AA^{\top}$  while the authority centrality is the leading right singular vector of A and the leading eigen-vector of  $A^{\top}A$ .

The hub and authority centralities and their variants are widely used in many fields to study how network positions affect a particular outcome of interest. In practice, as mentioned in Section 1, one needs to estimate the centralities and then use them as regressors in the subsequent regression to estimate the network effects. To achieve these two goals, the naive two-stage procedure has been widely used in the empirical studies although it lacks statistical justifications. In the following, we showcase some concrete examples on how centralities and the two-stage procedure are used in portfolio management, finance, and social media.

In portfolio management, recent research shows that, for a trade network of firms or industries or countries, a strategy that shorts portfolios of nodes with high centralities and longs those with low centralities yields significant excess return, where they use the two-stage procedure to provide empirical evidence (Hochberg et al., 2007, Ahern, 2013, Richmond, 2019). An accurate centrality estimate, therefore, can significantly boost the investment return. As a matter of fact, our case study in Section 7 shows that such a long-short strategy based on the centrality estimated using SuperCENT yields return three times as high as the existing method. In finance, financial institutions such as banks are usually linked through debt or equity, and thus an adversarial shock to one institution can be propagated to others via the debt or equity network (See Elliott et al. (2014), Glasserman and Young (2016), Vohra et al. (2020) and references therein). When a central institution is subject to a severe adversarial shock, the shock will propagate and the impact will be significantly amplified, resulting in a systemic risk for the whole economy. Thus, identifying the central institutions in the network is the key for policy-makers to impose additional supervision to mitigate concerns of systemic risk. To support their claims, many of the above-cited use the two-stage procedure for empirical evidence by regressing risk metrics on the estimated centrality of the financial institutions in a certain network. In social media such as Twitter or Facebook, networks often serve a role in information transmission or sharing. Individuals in the center of the social network, i.e., the "influencers", can significantly expedite information dissemination (Shao et al., 2018, Bovet and Makse, 2019). Identifying these central individuals can have a wide range of implications from marketing to information censorship.

One challenge of measuring centralities is that we often observe networks with observational error due to the cost of data collection (Lakhina et al., 2003, Banerjee et al., 2013, De Paula, 2017, Breza et al., 2020, Zhu and Yang, 2020). For example, to measure the social connections between people, researchers usually use the friendship on Facebook or Twitter to measure the tie, which is obviously not a perfect measure of the social connection strength. In particular, Banerjee et al. (2013) uses the self-reported friendship in villages of India to measure the social tie between villagers, which is subject to self-reporting and subjective bias. Zhu and Yang (2020) uses the patent citations to measure the knowledge flow between companies, which neglects the communication among workers or executives. The measurement noise of the links could significantly worsen the estimation of the centrality; see Borgatti et al. (2006), Frantz et al. (2009), Wang et al. (2012), Martin and Niemeyer (2019) and the references therein. The inaccurate centrality estimation will further affect the subsequent analysis.

#### 3. A UNIFIED FRAMEWORK

# 3.1. Set-up and notation

We observe a sample of n observations  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$  where  $y_i \in \mathbb{R}$  is the response and  $\mathbf{x}_i \in \mathbb{R}^p$  is the vector of p covariates for the *i*-th observation as in the multivariate regression setting. The data can be represented in matrix form. Let  $\mathbf{y} \in \mathbb{R}^n$  denote the column vector of outcome and  $\mathbf{X} \in \mathbb{R}^{n \times p}$  denote the design matrix with n rows and p columns. We consider the fixed design by treating  $\mathbf{X}$  fixed. In addition, we observe a weighted and directed network  $\mathcal{G} = (V, E)$  representing connections between the observations. Here,  $V = (1, 2, \dots, n)$  are the nodes corresponding to n observations and  $E \subset V \times V$  is the set of directed edges. The directed network  $\mathcal{G}$  can be represented by an asymmetric adjacency matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  where  $a_{ij} = w_{ij}$  if  $(i, j) \in E$  with edge weight  $w_{ij} \in \mathbb{R} \setminus \{0\}$  and  $a_{ij} = 0$  otherwise. We adopt the following notation conventions. Given a matrix  $\mathbf{B} \in \mathbb{R}^{m \times n}$ , define its Frobenius

We adopt the following notation conventions. Given a matrix  $\boldsymbol{B} \in \mathbb{R}^{m \times n}$ , define its Frobenius norm:  $\|\boldsymbol{B}\|_F = \sqrt{\sum_{ij} b_{ij}^2}$ , its  $L_2$  operator norm:  $\|\boldsymbol{B}\|_2 = \sup_{\|\mathbf{x}\|_2=1} \|\boldsymbol{B}\mathbf{x}\|_2$ , and the vectorization operator that converts  $\boldsymbol{B}$  into a column vector of dimension mn as:

$$\operatorname{vec}(\boldsymbol{B}) = (b_{11}, \dots, b_{m1}, b_{12}, \dots, b_{m2}, \dots, b_{1n}, \dots, b_{mn})^{\top}.$$

 $P_B$  denotes the projection matrix that projects onto the column space of B.

# 3.2. A unified framework

The hub and authority centralities are widely used in many fields as shown in Section 2, thanks to their natural implication on network positions. However, literature that investigates

their statistical properties is scarce. The most direct model for them is based on the low-rank decomposition of the network as follows: the adjacency matrix for the observed directed network is generated by some underlying true hub and authority centralities corrupted with noise,

$$\boldsymbol{A} = \boldsymbol{A}_0 + \boldsymbol{E} = d\boldsymbol{u}\boldsymbol{v}^\top + \boldsymbol{E} \tag{2}$$

where the true authority and hub centralities  $u, v \in \mathbb{R}^n$  are the parameters to be estimated,  $A_0$  is the true network, and E is the noise.

Such a low-rank mean plus noise model has been commonly adopted for matrix estimation or matrix denoise (Shabalin and Nobel, 2013, Yang et al., 2016) and matrix completion (Candes and Plan, 2010, Negahban and Wainwright, 2011). Since the hub and authority centralities are the *leading* left and right singular vectors of  $A_0$  respectively, it is natural to consider the noiseless *rank-one* structure for the network  $A_0$ . This framework can be extended to general rank  $r \leq (n-p)/2$ , although the implications of the *non-leading* singular vectors as centralities are unclear.

Note that u and v are only identifiable up to a scalar. Typically, in SVD, people assume u and v have unit length. However, in our framework of network analysis, we assume  $||u||_2 = ||v||_2 = \sqrt{n}$ , in view of the fact that the network can grow and consequently the centralities should roughly be on the same scale no matter how large the network is. Furthermore, we assume the noise matrix has entries with zero mean, i.e.,  $\mathbb{E}[e_{ij}] = 0$  for i, j = 1, 2, ..., n.

Under Model (2), with an extra assumption that all entries of E are i.i.d. with variance  $\sigma_a^2$ , Shabalin and Nobel (2013) has shown that the angle between the leading left singular vector of A and that of  $A_0$ , i.e., the estimated hub centrality and the true hub centrality u, can converge to a nonzero quantity or even asymptotically orthogonal as n goes to infinity, if the signal-tonoise ratio  $d/\sigma_a$  is not large enough. This implies that the naive estimation of the centralities by implementing SVD on the observed network will fail in the presence of large noise.

Fortunately, when the naive estimation is inconsistent, it is still possible to obtain a consistent estimation of the centralities by including additional information from other sources. As we discussed in Sections 1 and 2, the positions of agents (nodes) impact the agents' behaviors and thus shape their outcomes. Since the centralities measure the positions of nodes, researchers often study the relationship between the centralities and a certain response variable of interest so as to investigate the network effect. And this relationship contains an additional source of information for the centralities.

To explore such a relationship, the common practice is to regress the response variable on the *estimated* hub and authority centralities, which are obtained through the SVD of the observed network  $A_0$ . But the generative model is actually prescribed as follows: for the *i*-th observation, the outcome  $y_i$  depends on the *true* hub and authority centralities  $u_i, v_i$  along with the covariates  $x_i \in \mathbb{R}^p$ :

$$y_i = \beta_x^{\top} \boldsymbol{x}_i + \beta_u u_i + \beta_v v_i + \epsilon_i, \qquad (3)$$

where  $\beta_x \in \mathbb{R}^p$  is the vector of regression coefficients and  $\beta_u, \beta_v \in \mathbb{R}$  are the coefficients of the hub and authority centralities respectively. The nuance of using the *estimated* versus the *true* centralities has consequences which we will explicate later. At this stage, no extra assumptions are imposed on the distribution of the regression error, except that we assume  $\mathbb{E}[\epsilon] = 0$  and  $\operatorname{Var}(\epsilon) = \sigma_y^2 I_n$ . We further assume n > p + 2 and  $\mathbf{X}^\top \mathbf{X}$  is invertible.

Putting (2) and the matrix version of (3) together, we propose the following unified framework that encapsulates the two models:

$$\int \boldsymbol{A} = d\boldsymbol{u}\boldsymbol{v}^{\top} + \boldsymbol{E}, \tag{4a}$$

$$(y = X\beta_x + u\beta_u + v\beta_v + \epsilon.$$
 (4b)

With the unified framework (4) and the observed data  $\{A, X, y\}$ , our original two goals of centrality estimation and network regression analysis can be specified as the following three: (i) estimate the true centralities u, v and the true network  $A_0 = duv^{\top}$ ; (ii) estimate the regression coefficients of the predictors  $\beta_x$  and the centralities  $\beta_u, \beta_v$ ; (iii) construct *valid* confidence intervals for the centralities as well as the regression coefficients that account for the randomness in the observed network.

The unified framework unites our estimation goals and provides a theoretical framework to study the behaviors of the two-stage procedure. Furthermore, unifying the two models motivates our supervised network centrality estimation (SuperCENT) methodology, which we will describe formally in the next section. We name it the "supervised centrality estimation" because (X, y) in the regression (4b) can be thought of as the supervisors that offer additional supervision to the centrality estimation. It is expected that if the centralities indeed have strong predictive power (that is, the centrality regression coefficients  $\beta_u, \beta_v$  are large compared with the regression noise level  $\sigma_y$ ), the estimation of the centralities will be better when combining (4a) and (4b) together instead of only considering (4a). With the improve the estimation and inference in the regression model. A similar idea of supervision has also been implemented in matrix decomposition (Li et al., 2016), albeit the absence of response prediction.

**Remark 1.** The unified framework (4) can be extended to a general case with rank r, that is,

$$\int \boldsymbol{A} = \boldsymbol{U}\boldsymbol{D}\boldsymbol{V}^{\top} + \boldsymbol{E},\tag{5a}$$

$$(y = X\beta_x + U\beta_u + V\beta_v + \epsilon,$$
 (5b)

where U and V are orthonormal matrices of size  $n \times r$  and D is a  $r \times r$  diagonal matrix with the singular values as the diagonal entries. Such a natural extension may contribute to the existing literature on regression with network information, because most existing papers consider the centralities as the predictors, which are the leading singular vectors. But potentially, the few leading singular vectors can offer additional predictive power as well.

## 4. METHODOLOGY

Given the unified framework (4), we first elaborate on the widely used two-stage procedure which first estimates the centralities and then estimates and provides inference for the centrality effects in the regression in Section 4.1. We then describe the SuperCENT methodology that simultaneously solves the centrality estimation and network regression in Section 4.2. Section 4.3 is devoted to the prediction problem where new nodes are included together with covariates and a corresponding new network. Section 4.4 describes strategies for tuning parameter selection in SuperCENT.

# 4.1. A naive two-stage procedure

As mentioned in Sections 1 and 2, given the unified framework (4) and the observed data  $\{A, X, y\}$ , a natural and naive procedure is the two-stage estimator, which can serve as a benchmark. To make notations consistent, we now formally introduce the estimator.

In view of (4a), the first stage is to perform SVD on the observed adjacency matrix A and take its leading left and right singular vectors and rescale them to have length  $\sqrt{n}$ , denoted as  $\hat{u}^{ts}$  and  $\hat{v}^{ts}$ , as the estimates for the centralities u and v respectively. The superscript ts stands for the two-stage procedure. In view of (4b), given the estimates  $\hat{u}^{ts}$  and  $\hat{v}^{ts}$ , the second stage performs the ordinary least square (OLS) regression of y on X and  $\hat{u}^{ts}, \hat{v}^{ts}$ , treating

 $\hat{u}^{ts}, \hat{v}^{ts}$  as given covariates. To be specific, the two-stage procedure solves the following two optimization problems sequentially,

$$(\hat{d}^{ts}, \hat{\boldsymbol{u}}^{ts}, \hat{\boldsymbol{v}}^{ts}) := \operatorname*{arg\,min}_{d, \|\boldsymbol{u}\|_{2} = \|\boldsymbol{v}\|_{2} = \sqrt{n}} \|\boldsymbol{A} - d\boldsymbol{u}\boldsymbol{v}^{\top}\|_{F}^{2}, \tag{6a}$$

$$\hat{\boldsymbol{\beta}}^{ts} = ((\hat{\boldsymbol{\beta}}_{x}^{ts})^{\mathsf{T}}, \hat{\beta}_{u}^{ts}, \hat{\beta}_{v}^{ts})^{\mathsf{T}} := \underset{\boldsymbol{\beta}_{x}, \beta_{u}, \beta_{v}}{\arg\min} \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}_{x} - \hat{\boldsymbol{u}}^{ts}\beta_{u} - \hat{\boldsymbol{v}}^{ts}\beta_{v}\|_{2}^{2}.$$
(6b)

Algorithm 1 outlines the two-stage procedure.

# Algorithm 1: The two-stage procedure

**Result**:  $\hat{d}^{ts}$ ,  $\hat{\boldsymbol{u}}^{ts}$ ,  $\hat{\boldsymbol{v}}^{ts}$  and  $\hat{\boldsymbol{\beta}}^{ts}$ . **Input:** the observed network  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , the design matrix  $\mathbf{X} \in \mathbb{R}^{n \times p}$ , the response vector  $\mathbf{y} \in \mathbb{R}^{n}$ . 1.  $(\hat{d}^{ts}, \hat{\boldsymbol{u}}^{ts}, \hat{\boldsymbol{v}}^{ts}) = \arg\min_{d, \|\boldsymbol{u}\|_{2} = \|\boldsymbol{v}\|_{2} = \sqrt{n}} \|\boldsymbol{A} - d\boldsymbol{u}\boldsymbol{v}^{\top}\|_{F}^{2}$ ; 2.  $\widehat{\mathbf{W}} = (\mathbf{X}, \hat{\boldsymbol{u}}^{ts}, \hat{\boldsymbol{v}}^{ts})$ ; 3.  $\hat{\boldsymbol{\beta}}^{ts} = (\widehat{\mathbf{W}^{\top}}\widehat{\mathbf{W}})^{-1}\widehat{\mathbf{W}^{\top}}\mathbf{y}$ .

**Remark 2.** Besides the estimation of the unknown parameters, valid inference is necessary to evaluate the network effect. In numerous empirical studies, researchers usually construct confidence intervals of the regression coefficients from the second stage regression by assuming that  $\hat{u}^{ts}$  and  $\hat{v}^{ts}$  are fixed and noiseless. This assumption simplifies the inferential statement because it follows that  $\operatorname{cov}(\hat{\beta}^{ts}) = \sigma_y^2 (\widehat{W}^\top \widehat{W})^{-1}$  where  $\widehat{W} = (\mathbf{X}, \hat{u}^{ts}, \hat{v}^{ts})$ . However, the observed network A is one realization from  $A_0 + E$  as in Model (4a), which makes its singular vectors  $\hat{u}^{ts}, \hat{v}^{ts}$  random. If one proceeds with inference ignoring the randomness, then the inference loses its justifications and the ensuing validity due to violation of the assumption. We refer to such "ad-hoc" confidence interval as the "two-stage-adhoc" method. To correct for the randomness of the estimated singular vectors  $\hat{u}^{ts}, \hat{v}^{ts}$  and make valid inferences, the asymptotic distribution of the estimator is derived rigorously in Section 5. Remarks 5 and 7 further discuss the theoretical property of the two-stage-adhoc method. The toy experiment in Figure 1 and simulation results in Section 6 show that the two-stage-adhoc method is either conservative with low network noise level or invalid with high network noise level.

# 4.2. SuperCENT methodology

From the two-stage procedure, we observe that the second step of estimation and inference in the regression model depends on the first step of centrality estimation. The more accurate the centrality estimates are, the better we are able to make inference in the regression model. On the other hand, the centralities are incorporated in the regression model as regressors, (X, y)have supervising effect on centrality estimation and can boost the estimation accuracy.

Motivated by the above intuition, we propose to optimize the following objective function to obtain the SuperCENT estimates,

$$(\hat{d}, \hat{\boldsymbol{u}}, \hat{\boldsymbol{v}}, \hat{\boldsymbol{\beta}}_x, \hat{\beta}_u, \hat{\beta}_v) := \underset{\substack{\boldsymbol{\beta}_x, \beta_u, \beta_v \\ d, \|\boldsymbol{\boldsymbol{u}}\| = \|\boldsymbol{\boldsymbol{v}}\| = \sqrt{n}}{\operatorname{arg\,min}} \frac{1}{n} \|\boldsymbol{\boldsymbol{y}} - \boldsymbol{\boldsymbol{X}}\boldsymbol{\beta}_x - \boldsymbol{\boldsymbol{u}}\beta_u - \boldsymbol{\boldsymbol{v}}\beta_v\|_2^2 + \frac{\lambda}{n^2} \|\boldsymbol{\boldsymbol{A}} - d\boldsymbol{\boldsymbol{u}}\boldsymbol{\boldsymbol{v}}^\top\|_F^2.$$
(7)

The above objective function combines the residual sum of squares (6b) and the rank-one approximation error of the observed network A (6a). The connection between the two terms is the centralities. The tradeoff between the two terms can be tuned through a proper selection of the hyper-parameter  $\lambda$ . This idea is somewhat similar to the supervised SVD method in Li et al. (2016), but both the supervising mechanism and the optimization objective function of the SuperCENT are different.

To solve (7), we can use a block gradient descent algorithm by updating  $(\hat{d}, \hat{u}, \hat{v}, \hat{\beta})$  iteratively until convergence, where  $\hat{\beta} = (\hat{\beta}_x^T, \hat{\beta}_u, \hat{\beta}_v)^T$ . Such an iterative algorithm requires an initialization, which can be the first stage of the two-stage procedure, i.e.,  $(\hat{d}^{ts}, \hat{u}^{ts}, \hat{v}^{ts})$  from the SVD of A. The complete algorithm with a given tuning parameter  $\lambda$  is shown in Algorithm 2. We use  $(\hat{d}^{(t)}, \hat{u}^{(t)}, \hat{v}^{(t)}, \hat{\beta}^{(t)})$  to denote the estimations in the *t*-th iteration. The derivation of Algorithm 2 and the algorithm for a symmetric network with the eigenvector centrality, a special case of Model (4), are deferred to the supplement. We will discuss the methods to choose  $\lambda$  in Section 4.4, including cross-validation and others.

Algorithm 2: SuperCENT( $\mathbf{A}, \mathbf{X}, \mathbf{y}, \lambda$ ), an algorithm to solve (7).

**Result**:  $\hat{d}$ ,  $\hat{u}$ ,  $\hat{v}$ , and  $\hat{\beta}$ .

Input: the observed network  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , the design matrix  $\mathbf{X} \in \mathbb{R}^{n \times p}$ , the response vector  $\mathbf{y} \in \mathbb{R}^{n}$ , the tuning penalty parameter  $\lambda$ , the tolerance parameter  $\rho > 0$ , the maximum number of iteration T; Initiate  $(d^{(0)}, \mathbf{u}^{(0)}, \mathbf{v}^{(0)}) = \arg \min_{d, \|\mathbf{u}\|_{2} = \|\mathbf{v}\|_{2} = \sqrt{n}} \|\mathbf{A} - d\mathbf{u}\mathbf{v}^{\top}\|_{F}^{2}$ , t = 1; while  $\max(\|\mathbf{P}_{\mathbf{u}^{(t)}} - \mathbf{P}_{\mathbf{u}^{(t-1)}}\|_{2}, \|\mathbf{P}_{\mathbf{v}^{(t)}} - \mathbf{P}_{\mathbf{v}^{(t-1)}}\|_{2}) > \rho$  and t < T do 1.  $\mathbf{W}^{(t-1)} = (\mathbf{X}, \mathbf{u}^{(t-1)}, \mathbf{v}^{(t-1)});$ 2.  $\beta^{(t)} = (\mathbf{W}^{(t-1)^{\top}}\mathbf{W}^{(t-1)})^{-1}\mathbf{W}^{(t-1)^{\top}}\mathbf{y};$ 3.  $d^{(t)} = \frac{1}{n^{2}}\mathbf{u}^{(t-1)^{\top}}\mathbf{A}\mathbf{v}^{(t-1)};$ 4.  $\mathbf{u}^{(t)} = ((\beta_{u}^{(t)})^{2} + \lambda(d^{(t)})^{2})^{-1} [\beta_{u}^{(t)}(\mathbf{y} - \mathbf{X}\beta_{x}^{(t)} - \mathbf{v}^{(t-1)}\beta_{v}^{(t)}) + \frac{1}{n}\lambda d^{(t)}\mathbf{A}\mathbf{v}^{(t-1)}];$ 5. Normalize  $\mathbf{u}^{(t)}$  such that  $\|\mathbf{u}^{(t)}\|_{2} = \sqrt{n};$ 6.  $\mathbf{v}^{(t)} = ((\beta_{v}^{(t)})^{2} + \lambda(d^{(t)})^{2})^{-1} [\beta_{v}^{(t)}(\mathbf{y} - \mathbf{X}\beta_{x}^{(t)} - \mathbf{u}^{(t)}\beta_{u}^{(t)}) + \frac{1}{n}\lambda d^{(t)}\mathbf{A}^{\top}\mathbf{u}^{(t)}];$ 7. Normalize  $\mathbf{v}^{(t)}$  such that  $\|\mathbf{v}^{(t)}\|_{2} = \sqrt{n};$ 8.  $t \leftarrow t + 1;$ end

Note that although  $\hat{u}$  and  $\hat{v}$  with length  $\sqrt{n}$  are only identifiable up to the sign,  $\hat{u}\hat{\beta}_u$  and  $\hat{v}\hat{\beta}_v$  are uniquely identifiable. One procedure to determine the sign of all the parameters is as follows: first, find the entry that has the largest magnitude in  $\hat{u}$  and  $\hat{v}$ , and make that entry positive so that the sign of either the hub centrality  $\hat{u}$  or the authority centrality  $\hat{v}$  can be fixed; then adjust the sign of the other centrality since the product  $\hat{u}\hat{v}^{\top}$  is identifiable, and finally determine the signs of  $\hat{\beta}_u, \hat{\beta}_v$  accordingly.

# 4.3. Prediction

Once the model is fitted with training data, it can be used for prediction. Suppose there are  $n^*$  new observations, they have not only the new covariates  $X^*$  and the new network among themselves  $A^*$ , but also new edges connecting them with the *n* training observations. The

original network A is augmented to  $A^{all}$  as follows

$$\boldsymbol{A}^{all} = \begin{pmatrix} \boldsymbol{A}_{11} \ \boldsymbol{A}_{12} \\ \boldsymbol{A}_{21} \ \boldsymbol{A}_{22} \end{pmatrix} = \begin{pmatrix} \boldsymbol{A} \ \boldsymbol{A}_{12} \\ \boldsymbol{A}_{21} \ \boldsymbol{A}^* \end{pmatrix},$$
(8)

where  $A^{all}$  is of size  $(n + n^*) \times (n + n^*)$ .

Given the augmented network  $\hat{A}^{all}$  and new covariates  $X^*$ , the task is to make prediction for  $y^*$ . To make use of the regression equation  $\hat{y}^* = X^* \hat{\beta}_x + u^* \hat{\beta}_u + v^* \hat{\beta}_v$ , it is necessary to estimate  $u^*$  and  $v^*$ . Similarly as (4a), we assume  $A^{all} = du^{all}v^{all^{\top}} + E^{all}$  where  $u^{all} = (u^{\top}, u^{*^{\top}})^{\top}$  and  $v^{all} = (v^{\top}, v^{*^{\top}})^{\top}$ . Therefore, we obtain the model

$$\boldsymbol{A}^{all} = \begin{pmatrix} \boldsymbol{A}_{11} \ \boldsymbol{A}_{12} \\ \boldsymbol{A}_{21} \ \boldsymbol{A}_{22} \end{pmatrix} = \begin{pmatrix} \boldsymbol{A} \ \boldsymbol{A}_{12} \\ \boldsymbol{A}_{21} \ \boldsymbol{A}^* \end{pmatrix} = d \begin{pmatrix} \boldsymbol{u} \boldsymbol{v}^\top \ \boldsymbol{u} \boldsymbol{v}^{*\top} \\ \boldsymbol{u}^* \boldsymbol{v}^\top \ \boldsymbol{u}^* \boldsymbol{v}^* \end{pmatrix} + \begin{pmatrix} \boldsymbol{E}_{11} \ \boldsymbol{E}_{12} \\ \boldsymbol{E}_{21} \ \boldsymbol{E}_{22} \end{pmatrix}.$$
(9)

In view of (9), to obtain estimates of  $u^*$  and  $v^*$ , one can either perform SVD of  $A^*$  or SVD of  $A^{all}$  and reserve only the relevant components of the singular vectors. The latter approach is more accurate and is formally described in Algorithm 3.

Algorithm 3: SVD of  $A^{all}$  to obtain  $\hat{u}^*$  and  $\hat{v}^*$ 

**Result**:  $\hat{\boldsymbol{u}}^*$  and  $\hat{\boldsymbol{v}}^*$ . **Input**: The augmented network  $\boldsymbol{A}^{all}$ . 1.  $(\hat{\boldsymbol{u}}^{all}, \hat{\boldsymbol{v}}^{all})$  are the left and right singular vectors of  $\boldsymbol{A}^{all}$ ; 2.  $\hat{\boldsymbol{u}}^{all} = \operatorname{sign}(\hat{\boldsymbol{u}}^{\top} \hat{\boldsymbol{u}}_{1:n}^{all}) \hat{\boldsymbol{u}}^{all}$  and  $\hat{\boldsymbol{v}}^{all} = \operatorname{sign}(\hat{\boldsymbol{v}}^{\top} \hat{\boldsymbol{v}}_{1:n}^{all}) \hat{\boldsymbol{v}}^{all}$ ; 3.  $\hat{\boldsymbol{u}}^* = \hat{\boldsymbol{u}}_{(n+1):(n+n*)}^{all}$  and  $\hat{\boldsymbol{v}}^* = \hat{\boldsymbol{v}}_{(n+1):(n+n*)}^{all}$ ; 4. Rescale  $\hat{\boldsymbol{u}}^* = \frac{\sqrt{n}}{\|\hat{\boldsymbol{u}}_{1:n}^{all}\|_2} \hat{\boldsymbol{u}}^*$  and  $\hat{\boldsymbol{v}}^* = \frac{\sqrt{n}}{\|\hat{\boldsymbol{v}}_{1:n}^{all}\|_2} \hat{\boldsymbol{v}}^*$ .

**Remark 3.** (Sign and scaling issue) Since  $\hat{u}^{all}$  and  $\hat{v}^{all}$  are only identifiable up to sign, we determine their signs as Step 2 of Algorithm 3 such that the angles between the training proportions and the SuperCENT estimates are less than 90 degrees, i.e.,  $\operatorname{sign}(\hat{u}^{\top}\hat{u}_{1:n}^{all}) > 0$  and  $\operatorname{sign}(\hat{v}^{\top}\hat{v}_{1:n}^{all}) > 0$ . In addition,  $\hat{u}^*$  and  $\hat{v}^*$  need to be scaled to match with  $\hat{u}$  and  $\hat{v}$ . Recall that for identifiability,  $\hat{u}$  and  $\hat{v}$  are of norm  $\sqrt{n}$ , and  $\beta_u$  and  $\beta_v$  are of the corresponding scale. In the prediction process, we need to scale  $\hat{u}^*$  and  $\hat{v}^*$  accordingly so that  $\beta_u \hat{u}^* + \beta_v \hat{v}^*$  is on par with  $\beta_u \hat{u} + \beta_v \hat{v}$ . Step 3 of Algorithm 3 is designed for this purpose.

# 4.4. Selection of the tuning parameter $\lambda$

The tuning parameter  $\lambda$  can be selected using the K-fold cross-validation. Given the prediction procedure in Section 4.3, the cross-validation procedure can be easily carried out as follows. For each fold of validation data, we first fit the model using the remaining K - 1 folds with the corresponding induced subnetwork and obtain the estimates for the regression coefficients by implementing Algorithm 2; we then obtain the estimates of the centralities for the validation data by applying Algorithm 3; we last obtain the total prediction error for the validation data by combining the outcomes from the first two steps. The best tuning parameter  $\lambda$  is set to be the minimizer of the total cross-validation error that sums over all folds. Algorithm 4 outlines this procedure in more detail. Another strategy for selecting the tuning parameter is through generalized cross-validation (GCV), which can save computational time. Denote  $\hat{y} = H(\lambda)y$ , then the GCV criterion is

$$GCV(\lambda) = \frac{\|\hat{\boldsymbol{y}} - \boldsymbol{y}\|_2^2/n}{(1 - \operatorname{trace}(\boldsymbol{H}(\lambda))/n)^2}.$$
(10)

The explicit form of  $H(\lambda)$  can be derived from the proof of the theorems. The performance of GCV is left for future investigation.

As a third strategy, Remark 12 in the next section offers an alternative way to choose the tuning parameter based on the theoretical analysis of SuperCENT, which is less time-consuming than cross-validation. However, we recommend using the cross-validation strategy for the best performance based on the simulation results.

<b>Algorithm 4:</b> The cross-validation algorithm for SuperCENT to choose $\lambda$ .
<b>Result</b> : $\lambda_{min}$ .
Input: $A, X, y$ .
for $\lambda$ on a exponentially regular grid <b>do</b>
for each fold do
0. Split the covariates and response into training $\mathbf{X}_{fold,train}$ , $\mathbf{y}_{fold,train}$ and
validation $\mathbf{X}_{fold,val}$ , $\mathbf{y}_{fold,val}$ and denote the the induced sub-network
corresponding to the training data $A_{fold,train}$ ;
1. $(\hat{\boldsymbol{\beta}}_{fold,\lambda}, \hat{\boldsymbol{u}}_{fold,\lambda}, \hat{\boldsymbol{v}}_{fold,\lambda}) \leftarrow \text{SuperCENT}(\mathbf{A}_{fold,train}, \mathbf{X}_{fold,train}, \mathbf{y}_{fold,train}, \lambda)$
by Algorithm 2;
2. $\hat{u}_{fold,val}, \hat{v}_{fold,val} \leftarrow \text{SVD}(A)$ and re-scale by Algorithm 3;
3. $SSE_{fold,\lambda} = \ \mathbf{y}_{fold,val} - (\mathbf{X}_{fold,val}, \hat{\boldsymbol{u}}_{fold,val}, \hat{\boldsymbol{v}}_{fold,val}) \hat{\boldsymbol{\beta}}_{fold,\lambda}\ _{2}^{2};$
end
end
$\lambda_{min} = \min_{\lambda} \sum_{fold} SSE_{fold,\lambda}$

## 5. THEORETICAL PROPERTIES

We investigate the statistical properties of the two-stage procedure in Section 5.1 and SuperCENT in Section 5.2. The two main theorems provide the asymptotic distributions of the estimators under appropriate conditions, which can be used for inference. The corollaries state the convergence rates and the bias of the relevant quantities.

We first introduce some notations and assumptions. Recall that  $P_{\cdot}$  denotes the projection matrix, such as  $P_{u}$ ,  $P_{v}$  and  $P_{X}$ . Define  $\tilde{u} = (I - P_{X})u$ ,  $\tilde{v} = (I - P_{X})v$ , which are the centralities projected onto the orthogonal space of X. Denote  $c = \tilde{u}^{\top} \tilde{u} \tilde{v}^{\top} \tilde{v} - (\tilde{u}^{\top} \tilde{v})^{2}$  and  $C_{\tilde{u}\tilde{v}} = \begin{pmatrix} \tilde{u}^{\top} \tilde{u} \ \tilde{v}^{\top} \tilde{v} \end{pmatrix} = \begin{pmatrix} \tilde{u}^{\top} \\ \tilde{v}^{\top} \end{pmatrix} (\tilde{u} \ \tilde{v}).$ 

Assumption 1. Under the unified framework (4), the noise of the network  $e_{ij}$  independently follows  $N(0, \sigma_a^2)$ , and the noise of the outcome regression  $\epsilon_i$  independently follows  $N(0, \sigma_y^2)$ . Assumption 2. The fixed design matrix  $X \in \mathbb{R}^{n \times p}$  with n > p + 2 and  $X^{\top}X$  is invertible. The dimension p is not diverging.

Assumption 3. The scaled network noise-to-signal ratio  $\kappa = \frac{\sigma_a^2}{d^2n} \rightarrow 0$ .

In Assumption 1, the independence is assumed for simplicity. If the network noise  $e_{ij}$  or the regression noise  $\epsilon_i$  are dependent with known covariance, the theorems and corollaries still hold with slight modifications; if they are dependent with unknown covariance, extra assumptions on the covariance structure need to be made and new methodologies and theories should be developed. Assumption 3 is required for the consistency of the SVD of the observed noisy network A of model (4a), which can be seen from Corollary 1 below.

# 5.1. Theoretical properties for the two-stage procedure

Under the three aforementioned assumptions, the two-stage procedure in Algorithm 1 is consistent, and the asymptotic distribution of its estimators is given in Theorem 1. The convergence rates of the estimators are given in Corollaries 1 and 2. When Assumption 3 is violated, the two-stage procedure is no longer consistent – the centralities estimation is inconsistent in Stage 1, which leads to bias in the regression coefficients estimation in Stage 2 resembling the measurement error problem as shown in Corollary 3, consequently to the detriment of the inference.

Recall that the two-stage estimates from Algorithm 1 are denoted as  $\hat{d}^{ts}$ ,  $\hat{u}^{ts}$ ,  $\hat{v}^{ts}$  and  $\hat{\beta}^{ts} = (\hat{\beta}_u^{ts}, \hat{\beta}_v^{ts}, (\hat{\beta}_x^{ts})^{\top})^{\top}$ . Let  $\hat{A}^{ts} = \hat{d}^{ts} \hat{u}^{ts} \hat{v}^{ts\top}$  be the estimate of  $A_0$ .

THEOREM 1: Under the unified framework (4) and Assumptions 1, 2 and 3, the two-stage estimates converge to the following normal distribution asymptotically,

$$\begin{pmatrix} \hat{\boldsymbol{u}}^{ts} - \boldsymbol{u} \\ \hat{\boldsymbol{v}}^{ts} - \boldsymbol{v} \\ \operatorname{vec} \left( \hat{\boldsymbol{A}}^{ts} - \boldsymbol{A}_0 \right) \\ \hat{\beta}_u^{ts} - \beta_u \\ \hat{\beta}_v^{ts} - \beta_v \\ \hat{\boldsymbol{\beta}}_x^{ts} - \beta_x \end{pmatrix} \xrightarrow{\mathcal{D}} N \Big( \boldsymbol{0}_{(2n+n^2+2+p)\times 1}, \boldsymbol{\Sigma}^{ts} \Big), \tag{11}$$

where  $\Sigma^{ts} = C^{ts} \begin{pmatrix} \sigma_y^2 I_n & \mathbf{0}_{n \times n^2} \\ \mathbf{0}_{n^2 \times n} & \sigma_a^2 I_{n^2} \end{pmatrix} C^{ts^{\top}}$  and  $C^{ts}$  is a function of  $(d, u, v, X, \beta_u, \beta_v)$ , whose specific form is given in the supplement.

Recall the two-stage procedure first estimates the centralities u and v and then plugs the estimated centralities into the regression model. Therefore, the asymptotic distributions for  $\hat{u}^{ts}$ ,  $\hat{v}^{ts}$  and  $\hat{A}^{ts}$  only depend on the noise E from the network model, not the regression noise  $\epsilon$ . This can also be seen in the definition of  $C^{ts}$ , where the three top left blocks are zeros.

**Remark 4.** (Covariance of  $\hat{\boldsymbol{\beta}}^{ts}$ ) One important fact to emphasize is that the covariance of  $\hat{\boldsymbol{\beta}}^{ts}$  is not  $\sigma_y^2(\widehat{\mathbf{W}}^{\top}\widehat{\mathbf{W}})^{-1}$  where  $\widehat{\mathbf{W}} = (\boldsymbol{X}, \hat{\boldsymbol{u}}^{ts}, \hat{\boldsymbol{v}}^{ts})$ , which is the classical results of regression when  $\boldsymbol{X}, \hat{\boldsymbol{u}}^{ts}, \hat{\boldsymbol{v}}^{ts}$  are considered fixed. This makes sense, as in our model, the observed network contains noise, which makes the estimated centralities  $\hat{\boldsymbol{u}}^{ts}, \hat{\boldsymbol{v}}^{ts}$  from the first stage random quantities and invalidates the traditional covariance result. As a matter of fact, the bottom right three blocks of  $C^{ts}$  are not zero, which highlights this phenomenon. Corollary 2 further illustrates this fact and its consequences on the inference.

In what follows, we present the convergence rates of the estimated centralities  $\hat{u}^{ts}$ ,  $\hat{v}^{ts}$ , and the network  $\hat{A}^{ts}$  in Corollary 1, and the convergence rates of the regression coefficients  $\hat{\beta}^{ts}$  and the prediction error in Corollary 2. In Corollary 3, we show the bias of  $\hat{\beta}_u$  and  $\hat{\beta}_v$  when the two-stage is inconsistent.

The convergence rates of the centralities depend on the selection of the loss function. Ideally, since the scales of the centralities are not fully determined, one prefers the loss function  $\|P_{\hat{u}} - P_{u}\|_{2}^{2}$ , which equals the squared sine of the angle between  $\hat{u}$  and u,  $\sin^{2} \angle (\hat{u}, u)$ . However, the exact form of this loss function is not clean mathematically. Instead, we use the scaled Euclidean distance  $\|\hat{u} - \operatorname{sign}(\hat{u}^{\top}u)u\|_{2}^{2}/n = 2 - 2\cos^{2} \angle (\hat{u}, u)$ , which has a cleaner expression and is connected to the squared sine through  $\|P_{\hat{u}} - P_{u}\|_{2}^{2} = (\|\hat{u} - \operatorname{sign}(\hat{u}^{\top}u)u\|_{2}^{2}/n)[1 - (\|\hat{u} - \operatorname{sign}(\hat{u}^{\top}u)u\|_{2}^{2}/n)/4]$ . These two losses are approximately equivalent when the estimator is consistent and the loss goes to zero.

COROLLARY 1: (Rate of  $\hat{u}^{ts}$ ,  $\hat{v}^{ts}$ , and  $\hat{A}^{ts}$ ) Under the unified framework (4) and Assumptions 1, 2 and 3, the two-stage estimators satisfy the following

$$\frac{1}{n}\mathbb{E}\|\hat{\boldsymbol{u}}^{ts} - \boldsymbol{u}\|_{2}^{2} = \frac{1}{n}\mathbb{E}\|\hat{\boldsymbol{v}}^{ts} - \boldsymbol{v}\|_{2}^{2} = \frac{\sigma_{a}^{2}(n-1)}{d^{2}n^{2}}(1+o(1)) = O\left(\frac{\sigma_{a}^{2}}{d^{2}n}\right) = O(\kappa), \quad (12)$$

$$\mathbb{E}\frac{\|\widehat{\boldsymbol{A}}^{ts} - \boldsymbol{A}_0\|_F^2}{\|\boldsymbol{A}_0\|_F^2} = \frac{\sigma_a^2(2n-1)}{d^2n^2}(1+o(1)) = O\left(\frac{\sigma_a^2}{d^2n}\right) = O(\kappa).$$
(13)

According to Corollary 1, the rate for  $\hat{\boldsymbol{u}}^{ts}$ ,  $\hat{\boldsymbol{v}}^{ts}$  in (12) is  $\frac{\sigma_a^2(n-1)}{d^2n^2}$  and the rate for  $\hat{\boldsymbol{A}}^{ts}$  in (13) is  $\frac{\sigma_a^2(2n-1)}{d^2n^2}$ . They have the same order as  $\kappa = \frac{\sigma_a^2}{d^2n}$ , which suggests that the noise-to-signal ratio  $\kappa$  is the critical quantity that determines the consistency of the two-stage procedure.

Note that, for the two-stage procedure to be consistent, one needs  $\kappa \to 0$  as  $n \to \infty$ , which depends on the three parameters  $d, \sigma_a, n$  and their relationships. We first discuss two consistent scenarios: 1) the signal strength d and the noise level  $\sigma_a$  are of constant order while n diverges; 2) the noise level  $\sigma_a$  stays constant, but the signal strength d can decrease when more nodes are collected for the network (because the network edge density might decay with more nodes), in which case as long as  $d^2n \to \infty$ ,  $\kappa$  still goes to zero. For real data, it is possible that the observed network gets noisier with more nodes  $\sigma_a \to \infty$  and the signal strength decays  $d \to 0$ . In this case, it is highly likely that the two-stage procedure will be inconsistent, for example, with a fast diverging noise level; then SuperCENT can improve the performance and remain consistent as discussed below in Section 5.2.

COROLLARY 2: (Rate of  $\hat{\beta}^{ts}$ ) Under the unified framework (4) and Assumptions 1, 2 and 3, the two-stage estimators satisfy

$$\mathbb{E}(\hat{\beta}_{u}^{ts} - \beta_{u})^{2} = \left(\frac{\sigma_{y}^{2}}{c}\tilde{\boldsymbol{v}}^{\top}\tilde{\boldsymbol{v}}\right)$$
(14)

$$+\frac{\sigma_a^2}{c^2}\frac{1}{d^2n}\left[\beta_v^2\tilde{\boldsymbol{v}}^{\top}\tilde{\boldsymbol{v}}\tilde{\boldsymbol{u}}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{v}})\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}^{\top}\tilde{\boldsymbol{v}}+\beta_u^2\tilde{\boldsymbol{u}}^{\top}\tilde{\boldsymbol{v}}\tilde{\boldsymbol{v}}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{u}})\tilde{\boldsymbol{v}}\tilde{\boldsymbol{u}}^{\top}\tilde{\boldsymbol{v}}\right]\right)(1+o(1))$$
(15)

$$=O\left(\frac{\sigma_y^2}{n} + \frac{\sigma_a^2(\beta_u^2 + \beta_v^2)}{d^2n^2}\right),\tag{16}$$

$$\mathbb{E}(\hat{\beta}_v^{ts} - \beta_v)^2 = \left(\frac{\sigma_y^2}{c} \tilde{\boldsymbol{u}}^\top \tilde{\boldsymbol{u}}\right)$$
(17)

$$+\frac{\sigma_{a}^{2}}{c^{2}}\frac{1}{d^{2}n}\left[\beta_{u}^{2}\tilde{\boldsymbol{u}}^{\top}\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{u})\tilde{\boldsymbol{v}}\tilde{\boldsymbol{u}}^{\top}\tilde{\boldsymbol{u}}+\beta_{v}^{2}\tilde{\boldsymbol{v}}^{\top}\tilde{\boldsymbol{u}}\tilde{\boldsymbol{u}}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{v})\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}^{\top}\tilde{\boldsymbol{u}}\right]\right)(1+o(1))$$
(18)

$$=O\left(\frac{\sigma_y^2}{n} + \frac{\sigma_a^2(\beta_u^2 + \beta_v^2)}{d^2n^2}\right),\tag{19}$$

$$Cov\left(\hat{\boldsymbol{\beta}}_{x}^{ts}-\boldsymbol{\beta}_{x}\right)=\sigma_{y}^{2}\left[\left(\boldsymbol{X}^{\top}\boldsymbol{X}\right)^{-1}+\left(\boldsymbol{X}^{\top}\boldsymbol{X}\right)^{-1}\boldsymbol{X}^{\top}\left(\boldsymbol{u}\;\boldsymbol{v}\right)\boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1}\left(\boldsymbol{u}^{\top}_{\boldsymbol{v}}\right)\boldsymbol{X}\left(\boldsymbol{X}^{\top}\boldsymbol{X}\right)^{-1}\right]$$

$$(20)$$

$$+\sigma_a^2 \frac{1}{d^2 n} (\boldsymbol{X}^\top \boldsymbol{X})^{-1} \boldsymbol{X}^\top \bigg[ \beta_u^2 (\boldsymbol{I} - \boldsymbol{P}_u) + \beta_v^2 (\boldsymbol{I} - \boldsymbol{P}_v)$$
(21)

$$+ \left(\boldsymbol{u} \, \boldsymbol{v}\right) \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \beta_{v}^{2} \tilde{\boldsymbol{u}}^{\top} (\boldsymbol{I} - \boldsymbol{P}_{v}) \tilde{\boldsymbol{u}} & 0\\ 0 & \beta_{u}^{2} \tilde{\boldsymbol{v}}^{\top} (\boldsymbol{I} - \boldsymbol{P}_{u}) \tilde{\boldsymbol{v}} \end{pmatrix} \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \boldsymbol{u}^{\top}\\ \boldsymbol{v}^{\top} \end{pmatrix} \end{bmatrix} \boldsymbol{X} (\boldsymbol{X}^{\top} \boldsymbol{X})^{-1} (1 + o(1)) (22)$$

**Remark 5.** (Comments on (14)-(16) for  $\beta_u$ ) For the variance of  $\hat{\beta}_u^{ts}$ , the first term (14) is the typical expression for traditional regression with deterministic predictors. The additional term (15) is caused by the randomness of  $\hat{u}^{ts}$ ,  $\hat{v}^{ts}$ . Note that the second term (15) is non-negative and is zero if  $\sigma_a = 0$ , or  $\tilde{u} \perp \tilde{v}$ . Furthermore, the first term in (16) is of order  $\sigma_y^2/n$  while the second term is of order  $\frac{\sigma_a^2(\beta_u^2 + \beta_v^2)}{d^2n^2} = \kappa \frac{\beta_u^2 + \beta_v^2}{n}$ . So if  $\beta_u$  and  $\beta_v$  are of constant order and  $\kappa \rightarrow 0$ , the second term is of smaller order than the first term.

The above Theorem and Corollaries assume Assumptions 1-3. When Assumption 3 is violated, i.e.,  $\kappa \rightarrow 0$ , the two-stage cannot estimate u and v consistently, and consequently  $\hat{\beta}_u^{ts}$ and  $\hat{\beta}_v^{ts}$  are biased.

COROLLARY 3: (Bias of  $\hat{\beta}_u^{ts}$  and  $\hat{\beta}_v^{ts}$  when the two-stage is inconsistent) Let  $\rho = \operatorname{cor}(\boldsymbol{u}, \boldsymbol{v})$ . If  $\boldsymbol{\beta}_x = \mathbf{0}$  or  $\operatorname{cov}(\boldsymbol{X}, (\boldsymbol{uv})) = \mathbf{0} \in \mathbb{R}^{p \times 2}$ , then under the unified framework (4) and Assumptions 1-2,

$$\operatorname{plim} \hat{\beta}_{u}^{ts} = \frac{(1+\kappa-\rho^{2})\beta_{u}+\kappa\rho\beta_{v}}{(1+\kappa)^{2}-\rho^{2}},$$
(23)

$$\operatorname{plim} \hat{\beta}_{v}^{ts} = \frac{(1+\kappa-\rho^{2})\beta_{v}+\kappa\rho\beta_{u}}{(1+\kappa)^{2}-\rho^{2}}.$$
(24)

**Remark 6.** (Conditions for the bias of  $\hat{\beta}_u^{ts}$  and  $\hat{\beta}_v^{ts}$ ) Corollary 3 assumes either  $\beta_x = 0$ , namely the true regression model only involves two centrality predictors, both of which have measurement errors, or  $\operatorname{cov}(\boldsymbol{X}, (\boldsymbol{uv})) = \mathbf{0} \in \mathbb{R}^{p \times 2}$ , namely the true regression model involves p predictors without measurement error and two centrality predictors with measurement errors, but the noiseless predictors and two centrality predictors are uncorrelated. These assumptions are adopted so that the expression of the bias has intuitive explanations to be followed momentarily. In general, when the p noiseless ones and the two centralities are correlated, the bias persists as long as the signal-to-noise of the network is small when  $\kappa \to 0$ , although the expressions are less comprehensive.

**Remark 7.** (Special cases for the bias of  $\hat{\beta}_u^{ts}$  and  $\hat{\beta}_v^{ts}$ ) There are a few special cases for Corollary 3. (i) When  $\kappa \to 0$ , plim  $\hat{\beta}_u^{ts} = \beta_u$  and plim  $\hat{\beta}_v^{ts} = \beta_v$ . That implies that when  $\hat{u}^{ts}$  and  $\hat{v}^{ts}$  are consistent estimators of u and v respectively, the existence of the estimation error of the centralities does not affect the bias-ness of the second stage regression. This

has correspondence to the errors-in-variables literature: when the relative amount of measurement error is small compared to the total variance of the observed variable, the OLS estimate is unbiased. In other words,  $\frac{1}{1+\kappa}$  can be viewed as the reliability, attenuation factor, or the signal-to-total variance ratio. (ii) When  $\kappa \rightarrow 0$ , but the two true centralities are uncorrelated,  $\rho = 0$ , we have plim  $\hat{\beta}_u^{ts} = \frac{1}{1+\kappa} \beta_u$  and plim  $\hat{\beta}_v = \frac{1}{1+\kappa} \beta_v$ . The OLS estimate is biased towards zero, and the degree of bias depends on  $\frac{1}{1+\kappa}$ . (iii) When  $\kappa \rightarrow 0$  and  $\rho \neq 0$ , if  $|\beta_u| \gg |\beta_v|$ , then plim  $\hat{\beta}_u^{ts} \approx \frac{(1+\kappa-\rho^2)\beta_u}{(1+\kappa)^2-\rho^2}$ , which is equivalent to plim  $\hat{\beta}_u^{ts} - \beta_u \approx -\frac{(1+\kappa)\kappa}{(1+\kappa)^2-\rho^2}\beta_u$ . This implies that  $\hat{\beta}_v^{ts}$  is biased away from zero. (iv) When  $\beta_u$  and  $\beta_v$  have similar size, the directions of the biases depends on the  $\beta_u$ ,  $\beta_v$ ,  $\rho$  and  $\kappa$ . For  $\hat{\beta}_u^{ts}$ , the asymptotic bias is plim  $\hat{\beta}_u^{ts} - \beta_u = \frac{-\kappa[(1+\kappa)\beta_u - \rho\beta_v]}{(1+\kappa)^2 - \rho^2}$ . Since the denominator is always larger than 0 because  $\kappa > 0$  and  $0 < \rho < 1$ , the direction of the bias depends on the  $\beta_u < \frac{\rho}{1+\kappa}\beta_v$ , plim  $\hat{\beta}_u^{ts} - \beta_u > 0$ . Similar conclusions

**Remark 8.** (Review of inference property in the framework of classical measurement error with one predictor) Consider the simplest classical population model with measurement error  $y = u\beta_u + \epsilon$ . For simplicity and clarity, we exclude X and v from the model. We only observe  $\hat{u} = u + \delta_u$  with measurement error  $\delta_u$  instead of the true u. It can be shown that the OLS estimate that regresses y on  $\hat{u}$  satisfies that  $\hat{\beta}_u \to \frac{1}{1+\kappa}\beta_u := \gamma\beta_u$ , where  $\kappa = \frac{\sigma_a^2}{d^2n}$  is our scaled noise-to-signal ratio and  $\gamma$  is the attenuation factor.

Furthermore, with the OLS estimate, the residual sum of squares has limit  $n(\sigma_y^2 + (1 - \gamma)^2 \beta_u^2 + \gamma^2 \beta_u^2 \kappa) = n(\sigma_y^2 + \beta_u^2 \kappa/(1 + \kappa))$ . Hence, the traditional estimate of  $\sigma_y^2$ ,  $\|\hat{\boldsymbol{y}} - \boldsymbol{y}\|_2^2/(n - 1)$ , over-estimates  $\sigma_y^2$ , the larger the  $\kappa$  and  $\beta_u$ , the larger the over-estimation. It can also be shown that the standard error of OLS  $\hat{\beta}_u$  converges to  $\gamma \sigma_y^2 + \gamma (1 - \gamma) \beta_u^2$ . Combined, the *t*-ratio goes to  $\sqrt{\gamma} \frac{\beta_u}{\sqrt{\sigma_y^2 + (1 - \gamma)\beta_u^2}}$ , which is smaller than  $\beta_u/\sigma_y$ . So the traditional inference ignoring the measurement error is conservative. Please refer to Pischke (2007) and Wooldridge (2015) for more details on measurement error.

**Remark 9.** (CI for  $\beta_u$  from the ad-hoc two-stage method) Remark 8 has a few implications on the confidence interval (CI) for  $\beta_u$  under the unified framework (4) when the two-stage estimator is consistent ( $\kappa \rightarrow 0$ ). On one hand, when all of the quantities in (14)-(15) (including  $\sigma_a, \sigma_y, d, n, \beta_u, \beta_v, \tilde{u}, \tilde{v}$ ) are known, one should use both terms to make valid inference. If one uses (14) alone while assuming noiseless  $\hat{u}, \hat{v}$  to construct the CI, i.e. the "two-stageadhoc" method, the inference is invalid unless  $\sigma_a = 0$ , or  $\tilde{u} \perp \tilde{v}$ . But the degree of invalidity is typically small because (15)  $\ll$  (14) when  $\kappa \rightarrow 0$ . On the other hand, when the quantities are unknown and need to be estimated and plugged into these two terms, the stories are different. Note that from the classical measurement error literature,  $\hat{\sigma}_y^2$  from the two-stage over-estimates  $\sigma_y^2$  and the inference based on the first estimated term (14) alone is already conservative. As a consequence, the inference based on the two estimated terms would be even more conservative with unnecessarily large width.

In the simulation study, we consider three relevant methods: two-stage-oracle uses both terms with the true parameters, two-stage uses both terms with the estimated parameters, and the two-stage-adhoc uses only the first term with the estimated parameters. The results are consistent with the above discussion.

can be drawn for  $\hat{\beta}_{w}^{ts}$ .

To further add to Remark 4 on the covariance of  $\hat{\beta}_x^{ts}$ , due to the randomness of  $\hat{u}^{ts}$ ,  $\hat{v}^{ts}$ , the covariance of  $\hat{\beta}_x^{ts}$  in Corollary 2 involves three terms, where (20) is the term involving  $\sigma_y^2$  and would be the traditional covariance of  $\hat{\beta}_x^{ts}$  if  $\hat{u}^{ts}$  and  $\hat{v}^{ts}$  were indeed fixed, the two additional terms (21)-(22) are the extra covariance caused by the randomness of  $\hat{u}^{ts}$  and  $\hat{v}^{ts}$ .

## 5.2. Theoretical properties for SuperCENT

This section states the theoretical properties of our proposed SuperCENT method. Intuitively, we expect SuperCENT to be superior to the two-stage procedure when the signal-to-noise ratio of the regression model is large (so that the supervision information from the regression model is strong enough), especially when the observed network is noisy. Theorem 2 shows the asymptotic distribution of the SuperCENT estimators under the same set of conditions as the two-stage and Corollaries 4 and 5 provide their convergence rates. Comparing SuperCENT with the two-stage, we further explicate the discrepancy between SuperCENT and the two-stage and the conditions such that SuperCENT outperforms the two-stage, particularly when the two-stage is inconsistent. Note that the SuperCENT estimates from Algorithm 2 with a given tuning parameter  $\lambda$  are denoted as  $\hat{d}$ ,  $\hat{u}$ ,  $\hat{v}$ , and  $\hat{\beta} = ((\hat{\beta}_x)^{\top}, \hat{\beta}_u, \hat{\beta}_v)^{\top}$ . Let  $\hat{A} = \hat{d}\hat{u}\hat{v}^{\top}$  be the estimate of  $A_0$ .

THEOREM 2: Under the unified framework (4) and Assumptions 1, 2 and 3, the SuperCENT estimators converge to the following normal distribution asymptotically,

$$\begin{pmatrix} \hat{\boldsymbol{u}} - \boldsymbol{u} \\ \hat{\boldsymbol{v}} - \boldsymbol{v} \\ \operatorname{vec} \left( \hat{\boldsymbol{A}} - \boldsymbol{A}_0 \right) \\ \hat{\beta}_u - \beta_u \\ \hat{\beta}_v - \beta_v \\ \hat{\boldsymbol{\beta}}_x - \boldsymbol{\beta}_x \end{pmatrix} \xrightarrow{\mathcal{D}} N \Big( \boldsymbol{0}_{(2n+n^2+2+p)\times 1}, \boldsymbol{\Sigma} \Big), \tag{25}$$

where  $\boldsymbol{\Sigma} = \boldsymbol{C} \begin{pmatrix} \sigma_y^2 \boldsymbol{I}_n & \boldsymbol{0}_{n \times n^2} \\ \boldsymbol{0}_{n^2 \times n} & \sigma_a^2 \boldsymbol{I}_{n^2} \end{pmatrix} \boldsymbol{C}^{\top}$  and  $\boldsymbol{C}$  is a function of  $(d, \boldsymbol{u}, \boldsymbol{v}, \boldsymbol{X}, \beta_u, \beta_v, \lambda)$ , whose specific form is given in the supplement.

Corollaries 4 and 5 provide the convergence rates of the SuperCENT estimators. To explicate the difference between the two-stage and SuperCENT, let

$$\delta_{ts,sc} = \frac{1}{(\lambda d^2 + \beta_u^2 + \beta_v^2)^2} \Big[ \frac{2\lambda d^2 + \beta_u^2 + \beta_v^2}{d^2 n} \sigma_a^2 - \sigma_y^2 \Big].$$
(26)

COROLLARY 4: (*Rate of*  $\hat{u}$ ,  $\hat{v}$  and  $\hat{A}$ ). Under the unified framework (4) and Assumptions 1, 2 and 3, the SuperCENT estimators satisfy the following,

$$\mathbb{E}\|\hat{\boldsymbol{u}} - \boldsymbol{u}\|_{2}^{2}/n = \left(\frac{\sigma_{a}^{2}(n-1)}{d^{2}n^{2}} - \frac{n-p-2}{n}\beta_{u}^{2}\delta_{ts,sc}\right)(1+o(1))$$
(27)

$$= O\left(\frac{\sigma_a^2}{d^2n} - \beta_u^2 \delta_{ts,sc}\right),\tag{28}$$

$$\mathbb{E}\|\hat{\boldsymbol{v}} - \boldsymbol{v}\|_{2}^{2}/n = \left(\frac{\sigma_{a}^{2}(n-1)}{d^{2}n^{2}} - \frac{n-p-2}{n}\beta_{v}^{2}\delta_{ts,sc}\right)(1+o(1))$$
(29)

$$= O\left(\frac{\sigma_a^2}{d^2n} - \beta_v^2 \delta_{ts,sc}\right),\tag{30}$$

$$\mathbb{E}\frac{\left\|\widehat{\boldsymbol{A}} - \boldsymbol{A}_{0}\right\|_{F}^{2}}{\left\|\boldsymbol{A}_{0}\right\|_{F}^{2}} = \left(\frac{\sigma_{a}^{2}(2n-1)}{d^{2}n^{2}} - \frac{n-p-2}{n}(\beta_{u}^{2} + \beta_{v}^{2})\delta_{ts,sc}\right)(1+o(1))$$
(31)

$$= O\left(\frac{\sigma_a^2}{d^2n} - \left(\beta_u^2 + \beta_v^2\right)\delta_{ts,sc}\right).$$
(32)

**Remark 10.** (The role of  $\delta_{ts,sc}$ ) Let us consider the estimation of u. Similar messages can be obtained for v and  $A_0$ . Comparing the rate of  $\hat{u}$  and  $\hat{u}^{ts}$  in (28) and (12),  $\mathbb{E} \| \hat{u}^{ts} - u \|_2^2 / n - \mathbb{E} \| \hat{u} - u \|_2^2 / n = \beta_u^2 \delta_{ts,sc}$ , where  $\delta_{ts,sc}$  is defined in (26). As one can see,  $\delta_{ts,sc}$  is the crucial quantity that measures the discrepancy between the two-stage and SuperCENT estimators for u. When  $\delta_{ts,sc} > 0$ , SuperCENT outperforms two-stage and vice versa. The positiveness  $\delta_{ts,sc} > 0$  requires  $d^2 n \sigma_y^2 < (2\lambda d^2 + \beta_u^2 + \beta_v^2) \sigma_a^2$ . On one hand, this inequality can be satisfied when  $\lambda$  is not too small and is satisfied when  $\lambda$  takes the optimal value  $\lambda_0 = n\sigma_u^2/\sigma_a^2$  given in the remark below. It implies that when the tuning parameter is properly selected, SuperCENT performs better than the two-stage. On the other hand, this inequality is more likely to hold when the signal of the regression  $\beta_u, \beta_v$  is large, or the noise of the regression  $\sigma_y$  is small, or the signal of the network d is small, or the noise of the network  $\sigma_a$  is large. This exactly corresponds to our intuition: when the signal-to-noise ratio in the regression model is high, we gain information from the regression model to assist centrality estimation; and the advantage is more pronounced when the signal-to-noise ratio in the network is low, which is exactly when the two-stage behaves poorly.

**Remark 11.** (Optimal  $\lambda$ ) SuperCENT achieves the best performance with the following  $\lambda$ value

$$\lambda_0 = \frac{n\sigma_y^2}{\sigma_a^2}.\tag{33}$$

SuperCENT with this  $\lambda_0$  also obtains the most improvement over the two-stage procedure. It can be directly derived by minimizing (28), (30), or (32). Plugging the optimal value  $\lambda_0$  into the SuperCENT objective function (7) leads to

$$(\hat{d}, \hat{\boldsymbol{u}}, \hat{\boldsymbol{v}}, \hat{\boldsymbol{\beta}}_x, \hat{\boldsymbol{\beta}}_u, \hat{\boldsymbol{\beta}}_v) = \underset{\substack{\boldsymbol{\beta}_x, \boldsymbol{\beta}_u, \boldsymbol{\beta}_v\\d, \|\boldsymbol{u}\|_2 = \|\boldsymbol{v}\|_2 = \sqrt{n}}{\operatorname{arg\,min}} \frac{\|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}_x - \boldsymbol{u}\boldsymbol{\beta}_u - \boldsymbol{v}\boldsymbol{\beta}_v\|_2^2}{\sigma_y^2} + \frac{\|\boldsymbol{A} - d\boldsymbol{u}\boldsymbol{v}^\top\|_F^2}{\sigma_a^2}, (34)$$

which is -2 times log likelihood when the errors  $\epsilon$  and E are normally distributed. Nevertheless, the objective function is just the scaled residual sum of squares of the regression and the scaled rank-one approximation error of the observed network, which does not require the normality assumption.

**Remark 12.** (SuperCENT- $\hat{\lambda}_0$  and SuperCENT- $\hat{\lambda}_{cv}$ ) The benefit of the optimal value  $\lambda_0$  is twofold: 1) to benchmark the cross-validation procedure in Algorithm 4; 2) to avoid the timeconsuming cross-validation by supplying a candidate for the tuning parameter  $\lambda$ . We can obtain a crude estimate of  $\lambda_0$  by plugging in the two-stage estimates of  $\sigma_y^2$  and  $\sigma_a^2$ . To be specific, after we obtain  $\hat{A}^{ts}$  and  $\hat{y}^{ts} = X\hat{\beta}_x^{ts} + \hat{u}^{ts}\hat{\beta}_u^{ts} + \hat{v}^{ts}\hat{\beta}_v^{ts}$  from the two-stage procedure

.....

by Algorithm 1, we estimate  $(\hat{\sigma}_y^{ts})^2 = \frac{1}{n-p-2}(\hat{y}^{ts} - y)^2$  and  $(\hat{\sigma}_a^{ts})^2 = \frac{1}{n^2} \|\hat{A}^{ts} - A_0\|_F^2$ . We then plug in and obtain  $\hat{\lambda}_0 = n(\hat{\sigma}_y^{ts})^2/(\hat{\sigma}_a^{ts})^2$ . We refer to the SuperCENT with given  $\hat{\lambda}_0$  as SuperCENT- $\hat{\lambda}_0$ , whose empirical performance is given in the simulation study. Furthermore,  $\hat{\lambda}_0$  can be used as a guide to lay out the cross-validation grid points in Algorithm 4, to obtain  $\hat{\lambda}_{cv}$  and SuperCENT- $\hat{\lambda}_{cv}$ .

**Remark 13.** (Comparison of the estimation of  $u, v, A_0$  when two-stage is *inconsistent*) For the two-stage procedure,  $\hat{u}^{ts}, \hat{v}^{ts}, \hat{A}^{ts}$  is consistent if and only if  $\kappa = \frac{\sigma_a^2}{d^2n} \rightarrow 0$ , which implies the network signal-to-noise ratio has to be large enough for the two-stage to be consistent. When  $\kappa = O(1)$ , the two-stage procedure is inconsistent. Can the SuperCENT estimates remain consistent under this regime?

The answer is positive. Plugging in the optimal  $\lambda_0$ , the rate of  $\mathbb{E} \| \hat{\boldsymbol{u}} - \boldsymbol{u} \|_2^2 / n$  in (28) becomes

$$\kappa \frac{1 + \kappa \frac{\beta_v^2}{\sigma_y^2}}{1 + \kappa \left(\frac{\beta_u^2}{\sigma_y^2} + \frac{\beta_v^2}{\sigma_y^2}\right)},\tag{35}$$

which is obviously smaller than  $\kappa$ , the rate of  $\mathbb{E} \| \hat{\boldsymbol{u}}^{ts} - \boldsymbol{u} \|_2^2 / n$  in (12). We want the above rate converges to 0 when  $\kappa = O(1)$ . Given (35), the convergence of  $\hat{\boldsymbol{u}}$  boils down to the signal-tonoise ratio of  $\boldsymbol{u}$  and  $\boldsymbol{v}$  in the network regression model, i.e.,  $\frac{\beta_u^2}{\sigma_y^2}$  and  $\frac{\beta_v^2}{\sigma_y^2}$ . One sufficient condition for convergence is then  $\frac{\beta_u^2}{\sigma_y^2} = O(n^c)$ , c > 0 and  $\frac{\beta_v^2}{\sigma_y^2} = O(1)$ , meaning the signal for  $\boldsymbol{u}$  has to be stronger than both the signal for  $\boldsymbol{v}$  and the noise  $\sigma_u$  to guarantee convergence of  $\hat{\boldsymbol{u}}$ .

stronger than both the signal for v and the noise  $\sigma_y$  to guarantee convergence of  $\hat{u}$ . If we want to guarantee the convergence of  $\hat{v}$  under this regime, one sufficient condition is  $\frac{\beta_v^2}{\sigma_y^2} = O(n^c)$ , c > 0 and  $\frac{\beta_u^2}{\sigma_y^2} = O(1)$ . This conflicts with the requirement of the convergence of  $\hat{u}$ . Fortunately, the rates of both  $\hat{u}$  and  $\hat{v}$  are smaller than those of  $\hat{u}^{ts}$  and  $\hat{v}^{ts}$ , so SuperCENT always improves the estimation: when  $\beta_u^2/\beta_v^2 = o(1)$  or  $\beta_v^2/\beta_u^2 = o(1)$ , one of  $\hat{u}$  and  $\hat{v}$  will be consistent. We will demonstrate this phenomenon in the simulation.

Lastly, for the estimation of  $A_0$  with  $\lambda_0$ , the rate of  $\mathbb{E} \| \hat{A} - A_0 \|_F^2 / \|A_0\|_F^2$  in (32) becomes

$$\kappa \frac{1}{1 + \kappa \left(\frac{\beta_u^2}{\sigma_y^2} + \frac{\beta_v^2}{\sigma_y^2}\right)}$$

which is much smaller than  $\kappa$ , the rate of  $\mathbb{E} \| \hat{A}^{ts} - A_0 \|_F^2 / \| A_0 \|_F^2$  in (13). Better yet, to ensure  $\hat{A}$  is consistent, we only require either  $\frac{\beta_u^2}{\sigma_y^2} = O(n^c)$ ,  $c > \text{or } \frac{\beta_v^2}{\sigma_y^2} = O(n^c)$ , c > 0. This means that, as long as one of SuperCENT  $\hat{u}$  or  $\hat{v}$  is consistent, SuperCENT  $\hat{A}$  is consistent as well, while two-stage  $\hat{A}^{ts}$  is only consistent when both  $\hat{u}^{ts}$  and  $\hat{v}^{ts}$  are consistent.

When the two-stage estimator is consistent, the supervision effect of SuperCENT only takes place for the estimation of u, v, A and the in-sample prediction, but not for  $\beta_u, \beta_v, \beta_x$ , as shown in Corollary 5.

COROLLARY 5: (Rate of  $\hat{\beta}^{ts}$ ) Under the unified framework (4) and Assumptions 1, 2 and 3, the SuperCENT estimators satisfy the following,

$$\begin{split} \mathbb{E}(\hat{\beta}_u - \beta_u)^2 &= \mathbb{E}(\hat{\beta}_u^{ts} - \beta_u)^2 = O\left(\frac{\sigma_y^2}{n} + \frac{\sigma_a^2(\beta_u^2 + \beta_v^2)}{d^2 n^2}\right),\\ \mathbb{E}(\hat{\beta}_v - \beta_v)^2 &= \mathbb{E}(\hat{\beta}_v^{ts} - \beta_v)^2 = O\left(\frac{\sigma_y^2}{n} + \frac{\sigma_a^2(\beta_u^2 + \beta_v^2)}{d^2 n^2}\right),\\ Cov\left(\hat{\boldsymbol{\beta}}_x - \boldsymbol{\beta}_x\right) &= Cov\left(\hat{\boldsymbol{\beta}}_x^{ts} - \boldsymbol{\beta}_x\right). \end{split}$$

**Remark 14.** (Comparison of the estimation of  $\beta_u, \beta_v$ ) When  $\kappa \to 0$  and the two-stage is consistent, from the perspective of regression coefficient estimation, SuperCENT and the two-stage are similar. However, when  $\kappa \to 0$  and the two-stage regression coefficient estimation is biased as shown in Corollary 3. While for SuperCENT, if the conditions in Remark 13 hold, the SuperCENT coefficient estimates are still consistent and satisfy Corollary 5.

## 6. SIMULATION

In this section, we investigate the empirical performances, including the *estimation accuracy* and *inference property* of the two-stage and SuperCENT estimators under various simulation setups. We describe our simulation setup and overview the simulation results in Section 6.1 and show the simulation results for the inconsistent regime of the two-stage procedure in Section 6.2. The simulation results for the consistent regime of the two-stage procedure is deferred to the supplement.

## 6.1. Simulation setup and results overview

We generate the network as  $\mathbf{A} = d\mathbf{u}\mathbf{v}^{\top} + \mathbf{E} \in \mathbb{R}^{n \times n}$ , where  $\mathbf{u} \in \mathbb{R}^{n \times 1}$  and  $\mathbf{v} \in \mathbb{R}^{n \times 1}$  are vectors of the hub and authority centralities and all the entries of  $\mathbf{E}$  follow  $N(0, \sigma_a^2)$  independently. The elements of  $\mathbf{u}$  are first generated from i.i.d. N(0, 1) and  $\mathbf{v} = 0.5\mathbf{u} + \epsilon_{\mathbf{v}}$  where  $\epsilon_{\mathbf{v}}$  are generated from i.i.d. N(0, 1).  $\mathbf{u}$  and  $\mathbf{v}$  are then re-scaled to have norm  $\sqrt{n}$ . For the regression model,  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta}_x + \mathbf{u}\beta_u + \mathbf{v}\beta_v + \epsilon$ , we set p = 3, the coefficients for the covariates  $\boldsymbol{\beta}_x = (1, 3, 5)^{\top}$ , the covariate matrix  $\mathbf{X}$  consists of a column of 1's and p - 1 columns whose entries follow N(0, 1) independently, and  $\boldsymbol{\epsilon} \sim N(0, \sigma_y^2 \mathbf{I}_n)$ .

Since only the signal-to-noise of the network  $d/\sigma_a$  and the signal-to-noise of the regression  $(\beta_u/\sigma_y, \beta_v/\sigma_y)$  matter to the properties of our estimator and inference, we fix  $n = 2^8$ , d = 1, and  $\beta_v = 1$  to study the effect of  $\sigma_a, \sigma_y$ , and  $\beta_u$ . To study the effect of the signal-to-noise of the regression, we vary  $\sigma_y \in 2^{-4,-2,0}$  and  $\beta_u \in 2^{0,2,4}$  so that  $\frac{\beta_u^2}{\sigma_y^2} = O(n^c)$ ,  $c \ge 0$  and  $\frac{\beta_v^2}{\sigma_y^2} = O(1)$ . Now that the signal-to-noise of the network is solely controlled by  $\sigma_a$ , we vary  $\sigma_a$  to differentiate the regimes when the two-stage estimator is consistent with small  $\sigma_a$  and inconsistent with large  $\sigma_a$ . Specifically, for the *consistent regime of the two-stage*, i.e., when the network noise-to-signal ratio  $\kappa = \frac{\sigma_a^2}{d^2n} \to 0$ , we vary  $\sigma_a \in 2^{-4,-2}$  to keep  $\kappa < 2^{-12}$ . For the *inconsistent regime of the two-stage*, i.e.,  $\kappa = O(1)$ , we vary  $\sigma_a \in 2^{0,2}$  so that  $\kappa \in 2^{-8,-4} = O(1)$ .

For each setting, we compare several two-stage-based and SuperCENT-based procedures in terms of estimation accuracy and inference property as follows.

*Estimation accuracy* For the estimation accuracy, we compare the following procedures:

- 1. **Two-stage**: the two-stage procedure as in Algorithm 1;
- 2. **SuperCENT**- $\lambda_0$ : SuperCENT Algorithm 2 with the optimal  $\lambda_0 = n\sigma_y^2/\sigma_a^2$  as in Remark 11, where the true  $\sigma_y, \sigma_a$  are used. In other words, SuperCENT- $\lambda_0$  is not implementable for real data and it only serves as the benchmark;
- 3. **SuperCENT**- $\hat{\lambda}_0$ : SuperCENT with estimated tuning parameter  $\hat{\lambda}_0 = n(\hat{\sigma}_y^{ts})^2/(\hat{\sigma}_a^{ts})^2$ based on Remark 12, where  $(\hat{\sigma}_y^{ts})^2 = \frac{1}{n-p-2} \|\hat{y}^{ts} - y\|_2^2$  and  $(\hat{\sigma}_a^{ts})^2 = \frac{1}{n^2} \|\hat{A}^{ts} - A_0\|_F^2$ are estimated from the two-stage procedure. This way, SuperCENT- $\hat{\lambda}_0$  is implementable for real data;
- 4. SuperCENT- $\hat{\lambda}_{cv}$ : SuperCENT with tuning parameter  $\hat{\lambda}_{cv}$  chosen by 10-fold cross-validation as in Algorithm 4.

In the SuperCENT Algorithm 2, the tolerance parameter for the stopping criterion is set to be  $\rho = 10^{-4}$ . The following five performance metrics are used to compare these four procedures: the first two are from the perspective of network and centralities, and the last three are from the network regression.

- 1. The loss for estimating  $\boldsymbol{u}$  and  $\boldsymbol{v}$ ,  $l(\hat{\boldsymbol{u}}, \boldsymbol{u}) = \|\boldsymbol{P}_{\hat{\boldsymbol{u}}} \boldsymbol{P}_{\boldsymbol{u}}\|_{2}^{2}$  and  $l(\hat{\boldsymbol{v}}, \boldsymbol{v}) = \|\boldsymbol{P}_{\hat{\boldsymbol{v}}} \boldsymbol{P}_{\boldsymbol{v}}\|_{2}^{2}$ ;
- 2. The loss for estimating  $A_0$ ,  $l(\hat{A}, A_0) = \|\hat{A} A_0\|_F^2 / \|A_0\|_F^2$ ;
- 3. The normalized squared error loss for estimating  $\beta_u$  and  $\beta_v$ ,  $l(\hat{\beta}_u, \beta_u) = (\hat{\beta}_u \beta_u)^2 / \beta_u^2$ and  $l(\hat{\beta}_v, \beta_v) = (\hat{\beta}_v - \beta_v)^2 / \beta_v^2$ ;
- 4. The estimation bias for regression coefficients  $\hat{\beta}_u \beta_u$  and  $\hat{\beta}_v \beta_v$ .

Inference property For the inference property, let  $z_{1-\alpha/2}$  denote the  $(1-\alpha/2)$ -quantile of the standard normal distribution and we consider the following procedures to construct the confidence intervals (CIs) for the regression coefficient,  $CI_{\beta_u}$  and  $CI_{\beta_v}$ :

- 1. **Two-stage-adhoc**:  $\hat{\beta}^{ts} \pm z_{1-\alpha/2} \hat{\sigma}^{OLS}(\hat{\beta}^{ts})$ , where  $\hat{\beta}^{ts}$  is the two-stage estimate of  $\beta$  and  $\hat{\sigma}^{OLS}(\hat{\beta}^{ts})$  is the standard error from OLS, assuming  $\hat{\boldsymbol{u}}^{ts}, \hat{\boldsymbol{v}}^{ts}$  are fixed predictors;
- 2. **Two-stage-oracle**:  $\hat{\beta}^{ts} \pm z_{1-\alpha/2}\sigma(\hat{\beta}^{ts})$ , where  $\sigma(\hat{\beta}^{ts})$  is the standard error of  $\hat{\beta}^{ts}$ , whose mathematical expressions are given in (14)-(15) or (17)-(18) and the true parameters are plugged into those expressions;
- 3. **Two-stage**:  $\hat{\beta}^{ts} \pm z_{1-\alpha/2}\hat{\sigma}(\hat{\beta}^{ts})$ , where  $\hat{\sigma}(\hat{\beta}^{ts})$  is the standard error of  $\hat{\beta}^{ts}$  by plugging all the two-stage estimators into (14)-(15) or (17)-(18).
- 4. SuperCENT- $\lambda_0$ -oracle;  $\hat{\beta}^{\lambda_0} \pm z_{1-\alpha/2}\sigma(\hat{\beta}^{\lambda_0})$ , where  $\hat{\beta}^{\lambda_0}$  is the estimate of  $\beta$  by SuperCENT- $\lambda_0$  and  $\sigma(\hat{\beta}^{\lambda_0})$  follows (14)-(15) or (17)-(18), with the true parameters plugged in;
- 5. SuperCENT- $\hat{\lambda}_{cv}$ :  $\hat{\beta}^{\hat{\lambda}_{cv}} \pm z_{1-\alpha/2}\hat{\sigma}(\hat{\beta}^{\hat{\lambda}_{cv}})$ , where  $\hat{\beta}^{\hat{\lambda}_{cv}}$  is the estimate of  $\beta$  by SuperCENT- $\hat{\lambda}_{cv}$  and  $\hat{\sigma}(\hat{\beta}^{\hat{\lambda}_{cv}})$  is obtained by plugging the SuperCENT- $\hat{\lambda}_{cv}$  estimates into (14)-(15) or (17)-(18).

For these five methods, we compare the empirical coverage probability (CP) and the average width of the confidence intervals. The experiments are repeated 500 times. In parallel, the same five procedures can be used to study the inference regarding the true network  $A_0$ . The only difference is that the inference is made for each entry of  $A_0$ , which we denote as  $CI_{a_{ij}}$ ; and the CP and the average width reported are the average over all the entries  $a_{ij}$ .

To give an overview of the simulation results, Table I summarizes the comparison of the twostage method and SuperCENT in the consistent and the inconsistent regimes of the two-stage from the perspectives of both estimation and inference. SuperCENT universally outperforms the two-stage in terms of centrality estimation, regression coefficients estimation, and infer-

#### TABLE I

COMPARISON OF THE TWO-STAGE AND SUPERCENT IN THE CONSISTENT AND INCONSISTENT REGIMES OF THE TWO-STAGE WHEN  $\beta_u^2/\sigma_y^2 = O(n^c)$ , c > 0 and  $\beta_v^2/\sigma_y^2 = O(1)$ . IF  $\beta_v^2/\sigma_y^2 = O(n^c)$ , c > 0 and  $\beta_u^2/\sigma_y^2 = O(1)$ , THE RESULTS FOR u and v,  $\beta_u$  and  $\beta_v$  will be switched. In the estimation panel,  $\checkmark$  INDICATES THAT THE ESTIMATION IS CONSISTENT,  $\checkmark$  INDICATES THAT THE ESTIMATION IS INCONSISTENT. IN THE INFERENCE PANEL,  $\checkmark$  INDICATES THAT THE EMPIRICAL COVERAGE OF CONFIDENCE INTERVAL IS NO LESS THAN THE NOMINAL LEVEL,  $\checkmark$  INDICATES THAT THE CONFIDENCE INTERVAL FAILS TO REACH THE NOMINAL LEVEL. IN EACH ROW, SUPERCENT IS UNDERLINED WHENEVER IT OUTPERFORMS THE TWO-STAGE.

	$\kappa \to 0$ : two-s	tage consistent	$\kappa \rightarrow 0$ : two-stage inconsistent			
	Two-stage	SuperCENT	Two-stage	SuperCENT		
Estimation						
$\boldsymbol{u}$	1	Improved	×	1		
$\boldsymbol{v}$	1	Slightly Improved	×	✗ (Slightly Improved)		
$oldsymbol{A}_0$	1	Improved	×	<u> </u>		
$\beta_u$	1	$\checkmark$	✗ (Biased)	$\checkmark$		
$\beta_v$	<ul> <li>✓</li> </ul>	1	X (Biased)	X (Biased)		
Inference						
$CI_{\beta_u}$	✓ (Conservative)	✓ (Shorter)	×	$\checkmark$		
$CI_{\beta_v}$	1	$\checkmark$	×	×		
$CI_{a_{ij}}$	✓	$\checkmark$ (Shorter)	×	$\checkmark$		

ence. In what follows, we focus on the inconsistent regime of the two-stage and defer the consistent regime of the two-stage to the supplement.

## 6.2. Simulation results for the inconsistent regime of the two-stage procedure

In the inconsistent regime of two-stage where  $\kappa = \frac{\sigma_a^2}{d^2n} = O(1)$ , the two-stage procedure cannot consistently estimate  $\boldsymbol{u}$  and  $\boldsymbol{v}$ , i.e.,  $\hat{\boldsymbol{u}}^{ts}$  and  $\hat{\boldsymbol{v}}^{ts}$  are inconsistent. Recall that we fix  $\beta_v = 1$  and vary  $\sigma_y \in 2^{-4,-2,0}$ ,  $\beta_u \in 2^{0,2,4}$ , and  $\sigma_a \in 2^{0,2}$ , so that the scaled noise-to-signal of the network  $\kappa \in 2^{-8,-4} = O(1)$ , and the signal-to-noise of the regression  $\frac{\beta_z^2}{\sigma_y^2} = O(n^c)$ , c > 0 and  $\frac{\beta_v^2}{\sigma_y^2} = O(1)$ . In such a range, we expect the two-stage to be inaccurate for the estimation of  $\boldsymbol{u}, \boldsymbol{v}, \boldsymbol{A}_0$  and SuperCENT to be much more accurate for  $\boldsymbol{u}, \boldsymbol{A}_0$ , as mentioned in Remark 13. In addition, the two-stage estimates of  $\beta_u, \beta_v$  are biased as shown in Corollary 3 while SuperCENT estimates remain consistent as in Remark 14. As to the inference property, the two-stage-based confidence intervals are expected to be under-covered and wider than necessary while SuperCENT confidence intervals are valid and narrower.

Estimation accuracy Figure 3 shows the boxplot of the logarithm of  $l(\hat{u}, u)$  across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$  with d = 1 and  $\beta_v = 1$ . The rows correspond to  $\log_2(\sigma_a)$  and the columns correspond to  $\log_2(\beta_u)$ . For each panel, the x-axis is  $\log_2(\sigma_y)$  and the y-axis is  $\log_{10}(l(\hat{u}, u))$ . The super-imposed red symbols show the theoretical rates of  $\hat{u}^{ts}$  in Corollary 1 and that of  $\hat{u}$  in Corollary 4. The two-stage estimator performs the same no matter how large  $\sigma_y$  and  $\beta_u$  are, and it has a smaller error with smaller  $\sigma_a$ . The performance of SuperCENT is better with smaller  $\sigma_a$ ,  $\sigma_y$  or larger  $\beta_u$ . As expected, the three SuperCENT-based methods estimate u much more accurately than the two-stage procedure. In particular, the supervision effect of (X, y) is more pronounced when the noise of the outcome regression,  $\sigma_y$ , is small, or when the signal of



FIGURE 3.—Inconsistent regime of two-stage: The boxplot of  $\log_{10}(l(\hat{u}, u))$  for the four estimators across different  $\sigma_a, \sigma_y$  and  $\beta_u$  with fixed d = 1 and  $\beta_v = 1$  where  $l(\hat{u}, u) = \| P_{\hat{u}} - P_u \|_2^2$ . The super-imposed red symbols show the theoretical rates of the two-stage in Corollary 1 and SuperCENT in Corollary 4.



FIGURE 4.—Inconsistent regime of two-stage: The boxplot of  $\log_{10}(l(\hat{v}, v))$  for the four estimators across different  $\sigma_a, \sigma_y$  and  $\beta_u$  with d = 1 and  $\beta_v = 1$ . where  $l(\hat{v}, v) = \|P_{\hat{v}} - P_v\|_2^2$ . The super-imposed red symbols show the theoretical rates of the two-stage in Corollary 1 and SuperCENT in Corollary 4.

the outcome regression,  $\beta_u$ , is large, or when the network noise-to-signal,  $\sigma_a/d = \sigma_a$  is large. This validates Remarks 10 and 13 on the theoretical comparison of the estimators. Comparing the three SuperCENT-based methods, the benchmark SuperCENT- $\lambda_0$  is always the best, SuperCENT- $\hat{\lambda}_{cv}$  and SuperCENT- $\hat{\lambda}_0$  are sometimes worse than SuperCENT- $\lambda_0$ , but still better than the two-stage. SuperCENT- $\hat{\lambda}_0$  is typically comparable to or worse than SuperCENT- $\hat{\lambda}_{cv}$ , because SuperCENT- $\hat{\lambda}_0$  fails to locate the optimal  $\lambda_0$  due to inaccurate estimate of  $\sigma_a$  and  $\sigma_y$  from the two-stage procedure.



FIGURE 5.—Inconsistent regime of two-stage: The boxplot of  $\log_{10}(l(\hat{A}, A_0))$  for four estimators across different  $\sigma_a, \sigma_y$  and  $\beta_u$  with d = 1 and  $\beta_v = 1$  where  $l(\hat{A}, A_0) = \|\hat{A} - A_0\|_F^2 / \|A_0\|_F^2$ . The super-imposed red symbols show the theoretical rates of the two-stage in Corollary 1 and SuperCENT in Corollary 4.

For the estimation of v shown in Figure 4, the improvement of SuperCENT over two-stage is not as large as that of the estimation of u when  $\beta_u \in 2^{2,4}$ , because  $\frac{\beta_u^2}{\sigma_y^2} = O(n^c)$ , c > 0and  $\frac{\beta_v^2}{\sigma_y^2} = O(1)$ . But the improvement is still quite significant when  $\beta_u = 2^0 \approx \beta_v = 1$ . It is worth noting that the supervised effect to  $\hat{v}$  shrinks as  $\beta_u$  increases, leading to a different trend comparing Figures 3 and 4. This phenomena aligns with Remark 13 where we discuss the estimation of u and v when the two-stage is inconsistent. Specifically, the roles of u and v are not exchangeable, because here we have  $\beta_v \leq \beta_u$  by fixing  $\beta_v = 1$  and varying  $\beta_u \in 2^{0,2,4}$ . On the other hand, when  $\beta_v \gg \beta_u$  we should expect the improvement in estimating v to increase.

The conclusion for the estimation of  $A_0$  is similar to that of u as shown in Figure 5. With the improvement from estimating u and v, SuperCENT- $\hat{\lambda}_{cv}$  estimates  $A_0$  more accurately across all the settings. As claimed in Remark 13, the convergence of  $\hat{A}$  in this regime only requires  $\frac{\beta u^2}{\sigma_y^2} = O(n^c)$ , c > 0 or  $\frac{\beta v^2}{\sigma_y^2} = O(n^c)$ , c > 0. Therefore, with  $\beta_v = 1$ ,  $\beta_u > 1$ ,  $\hat{A}$  converges and  $l(\hat{A}, A_0) < l(\hat{A}^{ts}, A_0)$ . Comparing Figures 3, 4 and 5 altogether, when  $\beta_u = 2^0$ , SuperCENT improves the estimation of both u and v significantly; when  $\beta_u \in 2^{2,4}$ , SuperCENT improves the estimation of  $A_0$  a lot for all the ranges of  $\beta_u$ .

The attention is next turned to the regression coefficient  $\beta_u$ . Based on Corollary 3 on the issues of measurement error, the two-stage coefficient estimates tend to have bias under the inconsistent regime for the two-stage and the directions of the bias depend on the size of  $\beta_u$ ,  $\beta_v$ ,  $\kappa$ , and the correlation  $\rho$  between u and v. Figure 6 shows the estimation bias  $\hat{\beta}_u - \beta_u$ . With large  $\sigma_a$  or large  $\beta_u$ , the two-stage estimates suffer from sever attenuation bias, while SuperCENT can alleviate the bias. The attenuation bias by the two-stage can be explained by Remark 7 as follows. In this regime where  $\kappa = 2^{-8,-4} \rightarrow 0$ , u and v are correlated with  $\rho = \frac{2^{-1}}{\sqrt{1.25}}$ , and  $\beta_u \in 2^{2,4} > 1 \cdot 2^{-1.2} = \beta_v \frac{\rho}{(1+\kappa)}$ , then plim  $\hat{\beta}_u^{ts} - \beta_u < 0$ . Hence,  $\hat{\beta}_u^{ts}$  has an attenuation bias and the bias becomes larger as  $\beta_u$  increases. On the other hand, this also



FIGURE 6.—Inconsistent regime of two-stage: The boxplot of the bias  $\hat{\beta}_u - \beta_u$  across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$  with fixed d = 1 and  $\beta_v = 1$ . The dashed lines show  $\hat{\beta}_u - \beta_u = 0$ . The super-imposed red points show the median of the bias.



FIGURE 7.—Inconsistent regime of two-stage: The boxplot of  $\log_{10}(l(\hat{\beta}_u, \beta_u))$  for four estimators across different  $\sigma_a, \sigma_y$  and  $\beta_u$  with d = 1 and  $\beta_v = 1$  where  $l(\hat{\beta}_u, \beta_u) = (\hat{\beta}_u - \beta_u)^2 / \beta_u^2$ . The super-imposed red points show the median of  $\log_{10}(l(\hat{\beta}_u, \beta_u))$ .

implies that the two-stage estimation of  $\beta_v$  is biased away from zero and the bias is also larger as  $\beta_u$  increases as shown in the supplement. The improvement of SuperCENT over the twostage is relatively small and sometimes negligible for  $\beta_v$ .

In terms of the squared error loss of the estimation of  $\beta_u$  as shown in Figure 7, since  $\hat{\beta}_u^{ts}$  suffers from an attenuation bias with large  $\sigma_a$  and  $\beta_u$  and SuperCENT can alleviate the bias, SuperCENT improves over two-stage in the mean squared error  $l(\hat{\beta}_u, \beta_u)$ , corroborating Remark 14.



O SuperCENT-λ₀-oracle • SuperCENT-λ̂<sub>cv</sub> △ Two-stage-oracle ▼ Two-stage-adhoc ▲ Two-stage

FIGURE 8.—Inconsistent regime of two-stage: Empirical coverage of  $CI_{\beta_u}$  across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$  with d = 1 and  $\beta_v = 1$ . SuperCENT variants are labelled as circles ( $\circ \bullet$ ) and the two-stage variants are labelled as triangles ( $\Delta \checkmark \blacktriangle$ ). The hollow ones are for oracles and the solid ones are for non-oracles. The dashed lines show the nominal confidence level 0.95.



FIGURE 9.—Inconsistent regime of two-stage: The  $\log_{10}$  of the width of  $CI_{\beta_u}$  across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$  with d = 1 and  $\beta_v = 1$ . SuperCENT variants are labelled as circles ( $\circ \bullet$ ) and the two-stage variants are labelled as triangles ( $\Delta \checkmark \blacktriangle$ ). The hollow ones are for oracles and the solid ones are for non-oracles.

Inference property We now switch gears to the inference property. The bias in the estimation of  $\beta_u$  by the two-stage further affects its confidence interval. Figures 8 and 9 show the empirical coverage and the average width of the 95% confidence interval for  $\beta_u$  respectively. For the empirical coverage, when  $\beta_u$  is small (leftmost column), all the methods are close to the nominal level. When  $\beta_u$  increases and  $\sigma_a$  remains small (top right two panels), all the methods remain valid except for two-stage-oracle, but different methods remain valid for different rea-



FIGURE 10.—Inconsistent regime of two-stage:  $log_{10}$  of the average width of  $CI_{a_{ij}}$  across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$  with d = 1 and  $\beta_v = 1$ . SuperCENT variants are labelled as circles ( $\circ \bullet$ ) and the two-stage variants are labelled as triangles ( $\Delta \checkmark \blacktriangle$ ). The hollow ones are for oracles and the solid ones are for non-oracles.

sons. The two SuperCENT-based methods remain valid because there is no estimation bias and the estimation of the standard error is accurate. Two-stage and two-stage-adhoc remain valid mainly because they over-estimate  $\sigma_u^2$ , and this conservative-ness covers up the issue of bias. Two-stage-oracle uses the true  $\sigma_u^2$  and the issue of bias uncovers itself, consequently invalidating the inference. When  $\beta_u$  increases and  $\sigma_a$  gets large as well (bottom right two panels), the over-estimation of  $\sigma_y^2$  can no longer conceal the issue of bias and all two-stage related methods are not valid anymore. Again, the SuperCENT can mitigate the bias and the coverage probability is closer to the nominal level.

As for the width of  $CI_{\beta_u}$ , Figure 9 shows that the confidence intervals by the SuperCENTbased methods have better coverage and are narrower than those by the two-stage methods. The improvement in the width is more significant with larger  $\beta_u, \sigma_a$ .

For the confidence interval of the authority centrality coefficient,  $CI_{\beta_n}$ , the improvement of SuperCENT over two-stage is relatively small in terms of both the coverage probability and width as shown in the supplement, a similar phenomena as of the estimation accuracy of  $\beta_{v}$ .

Finally, we investigate the average coverage and the average width of confidence intervals for all the entries  $a_{ij}$  of  $A_0$  respectively. The average coverage probability of all the methods, Average<sub>*ij*</sub>( $CP(CIa_{ij})$ ), achieves the nominal level of 95% as shown in the supplement. The coverage tends to be slightly below the nominal coverage as  $\sigma_a$  increases, because the estimation becomes worse and the theorem only holds up to 1 + o(1). SuperCENT- $\lambda_{cv}$  is the closest to the nominal coverage in all the settings compared to the others. Figure 10 shows the  $\log_{10}$ of the average width of the CIs, Average<sub>ij</sub> (Width(CIa<sub>ij</sub>)). SuperCENT- $\lambda_0$ -oracle provides the shortest width among the four methods, followed by SuperCENT- $\lambda_{cv}$ . The widths of the confidence intervals of both SuperCENT-based methods are shorter than those of the two-stage methods. Again, the improvement of SuperCENT over the two-stage increases as  $\sigma_a$  and  $\beta_u$ increase or  $\sigma_y$  decreases.

### 7. CASE STUDY: GLOBAL TRADE NETWORK AND CURRENCY RISK PREMIUM

We consider a real case study with a triplet of  $\{A, X, y\}$ , where A is the country-level trade network, y is the currency risk premium, and X is GDP share, whose detailed information and construction will be given shortly. In this case study, we demonstrate that SuperCENT can provide more accurate estimation of the centrality, which is closely related to currency risk premium, and hence has a profound and lucrative implication on portfolio management. We further show that the SuperCENT method outperforms the two-stage methods in the inference of regression coefficients, providing less biased estimates and narrower CIs, and thus strengthens a related economic theory.

In the literature of international finance, economists have been studying currency risk premium extensively and are puzzled by its driving forces. The currency risk premium is formally defined as the excess return from holding foreign currency compared to holding the US dollar. Specifically, for an investor going long in a country/region i, the log risk premium "rx" at time t + 1 is

$$\mathbf{r}\mathbf{x}_{i,t+1} := r_{it} - r_t - \Delta q_{i,t+1},\tag{36}$$

where  $r_{it}$  is log interest rate of country/region *i*,  $r_t$  is the log interest rate of the U.S. and  $\Delta q_{i,t+1}$  is the appreciation of U.S. dollar.

In this case study, we investigate how the global trade network drives the currency risk premium and build a regression model that regresses the currency risk premium on the centrality from the international trade network. Such a predictive relationship is motivated by Richmond (2019), which developed a general equilibrium with international trade between countries and showed that countries' positions in the trade network can explain the difference in currency premiums across countries. Specifically, he showed an economic theory that countries that are central in the trade network exhibit lower currency risk premiums. This has two implications: (i) the regression coefficients for the centralities should be negative; (ii) international investors can obtain profit through taking a long-short strategy for foreign exchange – take a long position in currencies of countries with low centralities and a short position in currencies of countries with high centralities. Therefore, if the centralities can be estimated accurately, one can yield a significant investment return based on the strategy.

To verify the economic theory and construct a profitable portfolio, we first need to compute the currency risk premium and construct the global trade network following Richmond (2019), because these data are not directly available. We then apply the developed methodologies and theories to compare the two-stage and SuperCENT. We focus on the period between 1999 and 2013<sup>2</sup>.

To compute the currency risk premium, we obtain the interest rates and the exchange rates from DataStream. The currency risk premium can be calculated by plugging the interest rates and exchange rates into the definition of risk premium (36). Only 25 countries/regions have exchange rates available during the period of interest. We further exclude the region of Europe as it is not comparable to the others in the trade network, resulting in 24 countries/regions in the end<sup>3</sup>. We use a 5-year moving average of the currency risk premium. Specifically, when considering year t, all the relevant quantities are the averages from year t - 4 to year t. Figure 11 shows the time series plot of the rank of the 5-year moving average of risk premium from 2003 to 2012<sup>4</sup> for the 24 countries/regions. In each year, we rank the 24 countries/regions from

<sup>&</sup>lt;sup>2</sup>Euro was first adopted in 1999. Exchange rate for Malaysia (MYS) is not available from 1999 to 2004. The bilateral trade data is only available till 2013.

<sup>&</sup>lt;sup>3</sup>The list of country acronyms is provided in the supplement.

<sup>&</sup>lt;sup>4</sup>We leave the last available year 2013 for the validation purpose.



FIGURE 11.—Time series of the risk premium ranking in descending order from 2003 to 2012. The vertical dashed line indicates 2008, the year of the financial crisis.

1 to 24 from the largest risk premium to the smallest, i.e., in descending order of the risk premium.

Richmond (2019) defined the trade linkage as the trade amount normalized by the pair-wise total GDP, which represents the relative trade (export/import) intensity between two countries. Specifically, the trade linkage between two countries is computed as

$$a_{ijt} = \frac{S_{ijt}}{GDP_{it} + GDP_{jt}},\tag{37}$$

where  $S_{ijt}$  is the dollar value of goods and commodities exported from country *i* to country *j* at time *t*, and  $GDP_{it}$  is the GDP of country *i* at time *t* in U.S. dollar. The bilateral trade data come from the correlates of war project (COW) (Barbieri et al., 2009) and the International Monetary Fund (IMF) Direction of Trade Statistics<sup>5</sup>. Current U.S. dollar GDP (using 2015 as the base year) data are from the World Bank's World Development Indicators<sup>6</sup>. Same as the currency risk premium, we also use the 5-year moving average in the following analysis. In the supplement, we show a circular plot to visualize the average trade volume from 2003 to 2012.

As neither the two-stage nor SuperCENT is applicable for panel data, we will repeat the analysis for each year from 2003 to 2012. Besides the network and the response variable, we also include the predictor of GDP share, which is defined as the percentage of country/region GDP among the world GDP, where the world GDP is the total GDP of all available countries

<sup>&</sup>lt;sup>5</sup>https://data.imf.org/?sk=9D6028D4-F14A-464C-A2F2-59B2CD424B85

<sup>&</sup>lt;sup>6</sup>https://databank.worldbank.org/source/world-development-indicators

in the sample for that year. In summary, the models are, for each t,

$$a_{ijt} = d \operatorname{Hub}_{it} \times \operatorname{Authority}_{it} + e_{ijt}, \tag{38}$$

$$\mathbf{r}\mathbf{x}_{it} = \alpha + \beta_{ut} \cdot \mathbf{Hub}_{it} + \beta_{vt} \cdot \mathbf{Authority}_{it} + \beta_{xt} \cdot \mathbf{GDP} \text{ share}_{it} + \epsilon_{it}.$$
(39)

In Sections 5 and 6, we have demonstrated that the two-stage procedure is problematic under large network noise. In this case study, the observational error of the network comes from two sources: GDPs and the trade volumes, because each entry of the observed network  $a_{ijt}$  is defined as (37), i.e., the trade tie  $S_{ijt}$  between country *i* and country *j* normalized by their GDPs. GDPs and the trade volumes are often measured with errors. For the GDP, its accounting has been a challenge in macroeconomics (Landefeld et al., 2008). For the trade volume, the measurement errors are mostly due to (i) underground or illegal import and export; (ii) not including service trade; (iii) trade cost like transportation or taxes (Lipsey, 2009). Consequently, the observed trade network *A* can be very noisy and the two-stage method can perform badly.

On the other hand, SuperCENT can significantly improve over the two-stage when the network noise is large. In what follows, we focus on SuperCENT- $\hat{\lambda}_{cv}$  using 10-fold crossvalidation. We will refer SuperCENT- $\hat{\lambda}_{cv}$  to SuperCENT for simplicity and use the superscript sc for the SuperCENT- $\hat{\lambda}_{cv}$  estimates. Figure 12 shows the time series plots of the ranking of the hub centrality estimated by two-stage and SuperCENT for the 24 countries/regions, together with the ranking of the currency risk premium. Figure 13 is for the authority centrality. We rank the centrality in ascending order and the risk premium in descending order. Based on the negative relationship between centralities and risk premium established in Richmond (2019), the closer the trends of rankings between centralities and risk premium are, the better the centralities capture the time variation in the risk premium. In general, the centrality estimated by the two-stage procedure is relatively more stable over time compared to SuperCENT, since SuperCENT incorporates information of both the GDP share and currency risk premium, which is more volatile than the trade network itself. Asian trade hubs such as Hong Kong (HKG) and Singapore (SGP) are the most central; while countries like South Africa (ZAF) and New Zealand (NZL) are peripheral. Comparing with the ranking of risk premium, the time variation is not reflected in the centrality estimated by the two-stage procedure, while it can be captured by SuperCENT. In 2008, the year of the financial crisis, the SuperCENT centralities fluctuate together with risk premium while the two-stage centralities almost remained unchanged.

To emphasize the importance of an accurate centrality estimation for portfolio management, we examine whether a long-short strategy based on our estimated centrality can significantly boost investment performance. Specifically, Richmond (2019) showed in theory that a country with low centrality exhibits a higher expected currency premium than the ones with high centralities. We have shown in theory and simulations that the two-stage centrality estimation can be off from or even orthogonal to the truth when the signal-to-noise ratio of the network is low. Thus, if we can estimate the centralities more accurately than the two-stage procedure, we expect to obtain a higher return through longing countries with low centralities and shorting countries with high centralities.

To illustrate this insight, for each method, we take a long position on the currencies with the lowest 3 centralities (bottom 10%) and a short position on the currencies with the highest 3 centralities (top 10%). That is, we obtain a long-short return based on the estimated centrality of the period between year t - 4 and t. Figure 14 shows the year t + 1 return based on this strategy. The return based on the centrality estimated by SuperCENT is much higher than that of the two-stage procedure. Table II shows the 10-year average return based on this strategy with the top and bottom 3, 4, and 5 currencies, respectively. The 10-year average return based



FIGURE 12.—Time series of hub centrality ranking in ascending order from 2003 to 2012. The ranking of risk premium is in descending order, same as Figure 11. If the trend of centralities is close to the trend of risk premium, then the centralities capture the time variation of risk premium, based on the negative relationship between the two as claimed in Richmond (2019). The vertical dashed line indicates 2008, the year of the financial crisis.

on the SuperCENT centralities increased more than twice from that of the two-stage procedure. Thus, an accurate estimate of the centrality can significantly boost the average portfolio return.

We further demonstrate the superiority of SuperCENT in inference. Again since our method is not directly applicable to longitudinal data, we take the 10-year average of trade volume and GDP to construct a 10-year trade network and GDP share. Similarly, we take the 10-year average of risk premium as the response.

To better understand the behavior of the two-stage and SuperCENT estimators, it is crucial to know which regime the trade network belongs to. However, the true noise-to-signal ratio  $\kappa$  of the trade network is unknown, so we estimate it using SuperCENT:  $\hat{\kappa}^{sc} = 0.154 \approx 2^{-3}$ , which falls in the inconsistent regime of the two-stage. Note that for the simulation study, when  $\kappa = 2^{-8}$ , two-stage already shows inconsistency.

To further comprehend the behavior of SuperCENT and gauge how much improvement SuperCENT can potentially achieve in the inconsistent regime, we estimate the signal-to-noise ratio of the regression:  $(\hat{\beta}_u^{sc}/\hat{\sigma}_y^{sc})^2 = 7.6 \times 10^6 \approx 2^{23}$  and  $(\hat{\beta}_v^{sc}/\hat{\sigma}_y^{sc})^2 = 1.8 \times 10^5 \approx 2^{17}$ . Compared with the simulation settings in the inconsistent regime where  $\kappa = 2^{-4}$ ,  $\beta_u^2/\sigma_y^2 \leq 2^{16}$  and  $\beta_v^2/\sigma_y^2 \leq 2^8$ , we should expect SuperCENT to improve enormously over two-stage for both the estimation and inference of  $\beta_u$ , thanks to a large  $(\hat{\beta}_u^{sc}/\hat{\sigma}_y^{sc})^2$  and  $\hat{\beta}_u^{sc} \gg \hat{\beta}_v^{sc}$  under a relatively large  $\hat{\kappa}^{sc}$ .



FIGURE 13.—Time series of authority centrality ranking in ascending order from 2003 to 2012. The ranking of risk premium is in descending order, same as Figure 11. Similar to the hub centrality, if the trend of centralities is close to the trend of risk premium, then the centralities capture the time variation of risk premium, based on the negative relationship between the two as claimed in Richmond (2019). The vertical dashed line indicates 2008, the year of the financial crisis.

The three columns of Table III show the coefficient estimation, standard error and significant level for the two-stage-adhoc, two-stage, and SuperCENT, respectively. For the hub centrality  $\beta_u$ , (i) the estimate from two-stage-adhoc and two-stage is -0.0011, while the estimate from SuperCENT is -0.0021, which demonstrates the severe bias of two-stage in the inconsistent regime and the bias is towards zero because  $|\hat{\beta}_u^{sc}| = 0.0021 \gg 0.0003 = |\hat{\beta}_v^{sc}|^7$ ; (ii) the standard errors from two-stage-adhoc and two-stage are close to 0.0006, much larger than 0.0002 from SuperCENT, which reinforces the problem of overestimation of  $\sigma_y^2$  in two-stage-adhoc and two-stage; (iii) the above two facts combined make the confidence intervals by two-stage-adhoc; and two-stage unnecessarily wide, yet still invalid: consequently the hub centrality  $\beta_u$  is barely significant at level 0.1 using two-stage and is insignificant using two-stage-adhoc; (iv) the two facts in (i) and (ii) also lead to a valid but narrower confidence interval for SuperCENT, making the hub centrality a significant factor at level 0.01 for the currency risk premium; (v) conclusions drawn from the two-stage-adhoc and two-stage methods contradict the theory in Richmond (2019), while SuperCENT supports the theory.

<sup>&</sup>lt;sup>7</sup>Specifically,  $\hat{\beta}_{u}^{sc} = -0.0021 < -0.0003 \times \frac{0.673}{1+0.154} \approx -0.0002 = \hat{\beta}_{v}^{sc} \frac{\rho^{sc}}{1+\kappa^{sc}}$ , then plim  $\hat{\beta}_{u}^{ts} - \beta_{u} > 0$  as in Remark 7, and therefore the two-stage estimate is biased towards zero.



FIGURE 14.—Time series of the next-year return from 2004 to 2013 based on the long-short strategy that takes a long position on the currencies with the lowest 3 centralities and a short position on the currencies with the highest 3 centralities estimated from 2003 to 2012 respectively.

TABLE II

The 10-year average return							
	Top/Bottom 3		Top/Bottom 4		Top/Bottom 5		
	Hub	Authority	Hub	Authority	Hub	Authority	
SuperCENT Two-stage Relative difference	0.0031 0.0003 1 136%	0.0021 -0.0014 253%	0.0036 0.0008 338%	0.0019 -0.001 285%	0.0033 0.001 237%	0.0014 -0.0006 320%	

Let us consider other regression coefficients. For  $\beta_v$ , the estimate from two-stage-adhoc and two-stage is -0.0005, while the estimate from SuperCENT is -0.0003, implying a bias away from zero of two-stage in the inconsistent regime due to  $|\hat{\beta}_v^{sc}| \ll |\hat{\beta}_u^{sc}|$ . SuperCENT still improves its estimation and confidence interval, even though the improvement is not as large as  $\beta_u$  due to  $(\hat{\beta}_v^{sc})^2 \ll (\hat{\beta}_u^{sc})^2$  and the nonexchangeable roles of u and v. For  $\beta_x$ , the estimates from two-stage and SuperCENT are comparable, as there is no attenuation bias, but the widths of the confidence intervals from two-stage-adhoc and two-stage are much larger than that from SuperCENT, again as a result of the over-estimation of  $\sigma_y^2$ . Hence,  $\beta_x$  is barely significant when using two-stage-adhoc and two-stage, while very significant by SuperCENT.

## 8. CONCLUSION AND DISCUSSION

Motivated by the rising use of centralities in empirical literature, we highlighted and examined the centrality estimation and inference problems on a noisy network and presented the impact of the estimated centralities on the subsequent network regression. We showed that the commonly used two-stage procedure could yield inaccurate centrality estimates, biased regression coefficient estimates, and invalid inference, especially when the noise-to-signal ratio of the network is large. We proposed SuperCENT which incorporates the network-generation model with the network regression model to simultaneously estimate the centralities and the effects of centralities on the outcome. The additional information from the regression model supervises

### TABLE III

THE SUMMARY TABLE OF THE OUTCOME REGRESSION COMPARING THREE METHODS IN TERMS OF
COEFFICIENT ESTIMATION, STANDARD ERROR (IN PARENTHESIS) AND THE SIGNIFICANT LEVEL (BY
ASTERISKS).

	Dependent variable: Risk premium			
	Two-stage-adhoc	Two-stage	SuperCENT- $\hat{\lambda}_{cv}$	
GDP share $\beta_r$	-0.0159*	-0.0159*	-0.0162***	
·	(0.0083)	(0.0083)	(0.0037)	
Hub $\beta_u$	-0.0011	-0.0011*	-0.0021***	
, _	(0.0006)	(0.0006)	(0.0002)	
Authority $\beta_v$	-0.0005	-0.0005	-0.0003	
• / -	(0.0006)	(0.0006)	(0.0003)	
Note:		*p<0.1; **p	0<0.05; ***p<0.01	

centrality estimations and thus boosts the accuracy of the estimations. The better centrality estimations, in turn, benefits the network regression. We further derived the convergence rate and distribution of the SuperCENT estimator and provided valid confidence intervals for all the parameters of interest. The theoretical results are corroborated with extensive simulations and a real case study.

The SuperCENT model can be extended in multiple directions. One can consider a generalized linear model for the outcome model and extend SuperCENT to generalized SuperCENT. In the case when only a subset of covariates and outcomes are observed, semi-supervised SuperCENT can be developed. In the case when the network is partially observed, we can perform matrix completion with supervision. SuperCENT can be extended to a longitudinal model with additional assumptions by using techniques from Tensor decomposition as well as functional data analysis to obtain centralities that are smooth over time. In the case of ultra highdimensional problems, sparsity can be imposed on centralities due to the existence of abundant peripheral nodes.
#### SUPERCENT

#### REFERENCES

- Abel, A. B. (2017). Classical measurement error with several regressors. Technical report, Tech Rep. 2017. Tech Rep) Working Paper.
- Ahern, K. R. (2013). Network centrality and the cross section of stock returns. Available at SSRN 2197370.
- Allen, F., Cai, J., Gu, X., Qian, J., Zhao, L., and Zhu, W. (2019). Ownership network and firm growth: What do five million companies tell about chinese economy. *Available at SSRN 3465126*.
- Banerjee, A., Chandrasekhar, A. G., Duflo, E., and Jackson, M. O. (2013). The diffusion of microfinance. *Science*, 341(6144).
- Banerjee, A., Chandrasekhar, A. G., Duflo, E., and Jackson, M. O. (2019). Using gossips to spread information: Theory and evidence from two randomized controlled trials. *The Review of Economic Studies*, 86(6):2453–2490.
- Barbieri, K., Keshk, O. M., and Pollins, B. M. (2009). Trading data: Evaluating our assumptions and coding rules. Conflict Management and Peace Science, 26(5):471–491.
- Battaglini, M., Patacchini, E., and Rainone, E. (2021). Endogenous social interactions with unobserved networks. *The Review of Economic Studies*.
- Benzi, M., Estrada, E., and Klymko, C. (2013). Ranking hubs and authorities using matrix functions. *Linear Algebra and its Applications*, 438(5):2447–2474.
- Binkiewicz, N., Vogelstein, J. T., and Rohe, K. (2017). Covariate-assisted spectral clustering. *Biometrika*, 104(2):361–377.
- Borgatti, S. P., Carley, K. M., and Krackhardt, D. (2006). On the robustness of centrality measures under conditions of imperfect data. *Social networks*, 28(2):124–136.
- Bovet, A. and Makse, H. A. (2019). Influence of fake news in twitter during the 2016 us presidential election. *Nature communications*, 10(1):1–14.
- Bramoullé, Y., Djebbari, H., and Fortin, B. (2009). Identification of peer effects through social networks. *Journal of econometrics*, 150(1):41–55.
- Breza, E. and Chandrasekhar, A. G. (2019). Social networks, reputation, and commitment: evidence from a savings monitors experiment. *Econometrica*, 87(1):175–216.
- Breza, E., Chandrasekhar, A. G., McCormick, T. H., and Pan, M. (2020). Using aggregated relational data to feasibly identify network structure without network data. *American Economic Review*, 110(8):2454–84.
- Butts, C. T. (2003). Network inference, error, and informant (in) accuracy: a bayesian approach. *social networks*, 25(2):103–140.
- Candes, E. J. and Plan, Y. (2010). Matrix completion with noise. Proceedings of the IEEE, 98(6):925–936.
- Chandrasekhar, A. and Lewis, R. (2011). Econometrics of sampled networks. Unpublished manuscript, MIT.[422].
- Chen, M., Fernández-Val, I., and Weidner, M. (2021). Nonlinear factor models for network and panel data. *Journal* of *Econometrics*, 220(2):296–324.
- De Paula, A. (2017). Econometrics of network models. In *Advances in Economics and Econometrics: Theory and Applications: Eleventh World Congress*, volume 1, pages 268–323. Cambridge University Press Cambridge.
- De Paula, Á., Rasul, I., and Souza, P. (2019). Identifying network ties from panel data: theory and an application to tax competition. *arXiv preprint arXiv:1910.07452*.
- Elliott, M. and Golub, B. (2019). A network approach to public goods. *Journal of Political Economy*, 127(2):730–776.
- Elliott, M., Golub, B., and Jackson, M. O. (2014). Financial networks and contagion. *American Economic Review*, 104(10):3115–53.
- Fogli, A. and Veldkamp, L. (2021). Germs, social networks, and growth. *The Review of Economic Studies*, 88(3):1074–1100.
- Frantz, T. L., Cataldo, M., and Carley, K. M. (2009). Robustness of centrality measures under uncertainty: Examining the role of network topology. *Computational and Mathematical Organization Theory*, 15(4):303–328.
- Garber, S. and Klepper, S. (1980). Extending the classical normal errors-in-variables model. *Econometrica: Journal of the Econometric Society*, pages 1541–1546.
- Glasserman, P. and Young, H. P. (2016). Contagion in financial networks. *Journal of Economic Literature*, 54(3):779– 831.
- Gofman, M. (2017). Efficiency and stability of a financial architecture with too-interconnected-to-fail institutions. *Journal of Financial Economics*, 124(1):113–146.
- Graham, B. and De Paula, A. (2020). The Econometric Analysis of Network Data. Academic Press.
- Graham, B. S. (2017). An econometric model of network formation with degree heterogeneity. *Econometrica*, 85(4):1033–1063.
- Griliches, Z. (1986). Economic data issues. Handbook of econometrics, 3:1465-1514.
- Handcock, M. S. and Gile, K. J. (2010). Modeling social networks from sampled data. *The Annals of Applied Statistics*, 4(1):5.

- Hochberg, Y. V., Ljungqvist, A., and Lu, Y. (2007). Whom you know matters: Venture capital networks and investment performance. *The Journal of Finance*, 62(1):251–301.
- Hsieh, C.-S. and Lee, L. F. (2016). A social interactions model with endogenous friendship formation and selectivity. *Journal of Applied Econometrics*, 31(2):301–319.
- Jackson, M. O. (2010). Social and economic networks. Princeton university press.
- Jackson, M. O., Rogers, B. W., and Zenou, Y. (2017). The economic consequences of social-network structure. *Journal of Economic Literature*, 55(1):49–95.
- Jochmans, K. and Weidner, M. (2019). Fixed-effect regressions on network data. Econometrica, 87(5):1543–1560.
- Kleinberg, J. M. (1999). Authoritative sources in a hyperlinked environment. *Journal of the ACM (JACM)*, 46(5):604–632.
- Kolaczyk, E. D. (2010). Statistical analysis of network data: Methods and models.
- Lakhina, A., Byers, J. W., Crovella, M., and Xie, P. (2003). Sampling biases in ptopology measurements. In *IEEE INFOCOM 2003. Twenty-second Annual Joint Conference of the IEEE Computer and Communications Societies (IEEE Cat. No. 03CH37428)*, volume 1, pages 332–341. IEEE.
- Landefeld, J. S., Seskin, E. P., and Fraumeni, B. M. (2008). Taking the pulse of the economy: Measuring gdp. Journal of Economic Perspectives, 22(2):193–216.
- Le, C. M., Levin, K., and Levina, E. (2018). Estimating a network from multiple noisy realizations. *Electronic Journal of Statistics*, 12(2):4697–4740.
- Le, C. M. and Li, T. (2020). Linear regression and its inference on noisy network-linked data. *arXiv preprint* arXiv:2007.00803.
- Lee, L.-F. (2007). Identification and estimation of econometric models with group interactions, contextual factors and fixed effects. *Journal of Econometrics*, 140(2):333–374.
- Lee, L.-f., Liu, X., and Lin, X. (2010). Specification and estimation of social interaction models with network structures. *The Econometrics Journal*, 13(2):145–176.
- Li, G., Yang, D., Nobel, A. B., and Shen, H. (2016). Supervised singular value decomposition and its asymptotic properties. *Journal of Multivariate Analysis*, 146:7–17.
- Li, T., Levina, E., Zhu, J., et al. (2019). Prediction models for network-linked data. *The Annals of Applied Statistics*, 13(1):132–164.
- Lipsey, R. E. (2009). 1. Measuring International Trade in Services. University of Chicago Press.
- Liu, E. (2019). Industrial policies in production networks. The Quarterly Journal of Economics, 134(4):1883–1948.
- Ma, Z., Ma, Z., and Yuan, H. (2020). Universal latent space model fitting for large networks with edge covariates. J. Mach. Learn. Res., 21:4–1.
- Manski, C. F. (1993). Identification of endogenous social effects: The reflection problem. *The review of economic studies*, 60(3):531–542.
- Martin, C. and Niemeyer, P. (2019). Influence of measurement errors on networks: Estimating the robustness of centrality measures. *Network Science*, 7(2):180–195.
- Negahban, S. and Wainwright, M. J. (2011). Estimation of (near) low-rank matrices with noise and high-dimensional scaling. *The Annals of Statistics*, pages 1069–1097.
- Newman, M. E. (2018). Estimating network structure from unreliable measurements. *Physical Review E*, 98(6):062321.
- Newman, M. E. and Clauset, A. (2016). Structure and inference in annotated networks. *Nature communications*, 7(1):1–11.
- Pischke, S. (2007). Lecture notes on measurement error. London School of Economics, London.
- Richmond, R. J. (2019). Trade network centrality and currency risk premia. *The Journal of Finance*, 74(3):1315–1361.
- Rohe, K. (2019). A critical threshold for design effects in network sampling. The Annals of Statistics, 47(1):556–582.
- Shabalin, A. A. and Nobel, A. B. (2013). Reconstruction of a low-rank matrix in the presence of gaussian noise. *Journal of Multivariate Analysis*, 118:67–76.
- Shao, C., Ciampaglia, G. L., Varol, O., Yang, K.-C., Flammini, A., and Menczer, F. (2018). The spread of lowcredibility content by social bots. *Nature communications*, 9(1):1–9.
- Shen, D., Shen, H., and Marron, J. (2016). A general framework for consistency of principal component analysis. *The Journal of Machine Learning Research*, 17(1):5218–5251.
- Van Loan, C. F. and Golub, G. (1996). Matrix computations (johns hopkins studies in mathematical sciences).
- Vohra, R., Xing, Y., and Zhu, W. (2020). The network effects of agency conflicts. Available at SSRN 3556298.
- Wang, D. J., Shi, X., McFarland, D. A., and Leskovec, J. (2012). Measurement error in network data: A reclassification. *Social Networks*, 34(4):396–409.
- Wooldridge, J. M. (2015). Introductory econometrics: A modern approach. Cengage learning.
- Yan, T., Jiang, B., Fienberg, S. E., and Leng, C. (2019). Statistical inference in a directed network model with covariates. *Journal of the American Statistical Association*, 114(526):857–868.

- Yang, D., Ma, Z., and Buja, A. (2016). Rate optimal denoising of simultaneously sparse and low rank matrices. *The Journal of Machine Learning Research*, 17(1):3163–3189.
- Zhang, Y., Levina, E., and Zhu, J. (2016). Community detection in networks with node features. *Electronic Journal* of *Statistics*, 10(2):3153–3178.
- Zhu, W. and Yang, Y. (2020). Networks and business cycles. Available at SSRN.
- Zhu, X., Pan, R., Li, G., Liu, Y., and Wang, H. (2017). Network vector autoregression. *The Annals of Statistics*, 45(3):1096–1123.

# SUPPLEMENT TO "NETWORK REGRESSION AND SUPERVISED CENTRALITY ESTIMATION"

JUNHUI CAI University of Pennsylvania

DAN YANG The University of Hong Kong

> WU ZHU Tsinghua University

HAIPENG SHEN The University of Hong Kong

LINDA ZHAO University of Pennsylvania

This supplementary material contains more details and proofs that are deferred from the main text and is organized as follows. Section S1 provides the derivation of Algorithm 2. Section S2 presents the algorithm of SuperCENT for an undirected network with the eigenvector centrality. Section S3 gives the explicit mathematical expressions of the covariance matrices in Theorems 1 and 2. Section S4 shows more simulation results. Section S5 provides additional information for the case study. Finally, the proofs of the theoretical results are in Section S6.

## S1. DERIVATION OF ALGORITHM 2

In the paper, the proposed SuperCENT obtains estimates by optimizing the following objective function which combines the network model and the network regression model in the unified framework (4),

$$(\hat{d}, \hat{\boldsymbol{u}}, \hat{\boldsymbol{v}}, \hat{\boldsymbol{\beta}}_x, \hat{\beta}_u, \hat{\beta}_v) := \underset{\substack{\boldsymbol{\beta}_x, \beta_u, \beta_v\\d, \|\boldsymbol{u}\| = \|\boldsymbol{v}\| = \sqrt{n}}{\arg\min} \frac{1}{n} \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}_x - \boldsymbol{u}\beta_u - \boldsymbol{v}\beta_v\|_2^2 + \frac{\lambda}{n^2} \|\boldsymbol{A} - d\boldsymbol{u}\boldsymbol{v}^\top\|_F^2.$$
(S1)

To solve (S1), we use a block gradient descent algorithm by updating  $(\hat{d}, \hat{u}, \hat{v}, \hat{\beta})$  where  $\hat{\beta} = (\hat{\beta}_x^{\top}, \hat{\beta}_u, \hat{\beta}_v)^{\top}$ . The derivation of each step in each iteration of Algorithm 2 is described below. Denote  $\boldsymbol{W} = (\boldsymbol{X}, \boldsymbol{u}, \boldsymbol{v}), \beta = (\beta_x, \beta_u, \beta_v)$ , and  $\mathcal{L} := \frac{1}{n} \|\boldsymbol{y} - \boldsymbol{X}\beta_x - \boldsymbol{u}\beta_u - \boldsymbol{v}\beta_v\|_2^2 + \frac{\lambda}{n^2} \|\boldsymbol{A} - d\boldsymbol{u}\boldsymbol{v}^{\top}\|_F^2$ . Given  $\lambda$ , we minimize the objection function (S1) by setting the partial derivatives of all the parameters as zero. The partial derivatives are as follows.

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\beta}} = -\frac{2}{n} \boldsymbol{W}^{\mathsf{T}} (\boldsymbol{y} - \boldsymbol{W} \boldsymbol{\beta}), \tag{S2}$$

Dan Yang: dyanghku@hku.hk

Wu Zhu: zhuwu@sem.tsinghua.edu.cn

Haipeng Shen: haipeng@hku.hk

Junhui Cai: junhui@wharton.upenn.edu

Linda Zhao: lzhao@wharton.upenn.edu

$$\frac{\partial \mathcal{L}}{\partial d} = 2\lambda d - 2\lambda \boldsymbol{u}^{\mathsf{T}} \mathbf{A} \boldsymbol{v} / n^2, \tag{S3}$$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{u}} = -\frac{2}{n}\beta_u(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}_x - \boldsymbol{u}\beta_u - \boldsymbol{v}\beta_v) + \frac{2}{n}\lambda d^2\boldsymbol{u} - \frac{2}{n^2}\lambda d\mathbf{A}\boldsymbol{v},$$
(S4)

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{v}} = -\frac{2}{n}\beta_{\boldsymbol{v}}(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}_{x} - \boldsymbol{u}\beta_{u} - \boldsymbol{v}\beta_{v}) + \frac{2}{n}\lambda d^{2}\boldsymbol{v} - \frac{2}{n^{2}}\lambda d\mathbf{A}^{\top}\boldsymbol{u}.$$
(S5)

Setting the partial derivatives above as zeros yields the estimates

$$\widehat{\boldsymbol{\beta}} = (\widehat{\boldsymbol{W}}^{\top} \widehat{\boldsymbol{W}})^{-1} \widehat{\boldsymbol{W}}^{\top} \boldsymbol{y}, \tag{S6}$$

$$\hat{d} = \frac{1}{n^2} \hat{\boldsymbol{u}}^{\mathsf{T}} \boldsymbol{A} \hat{\boldsymbol{v}},\tag{S7}$$

$$\hat{\boldsymbol{u}} = \left(\hat{\beta}_u^2 + \lambda \hat{d}^2\right)^{-1} \left[\hat{\beta}_u(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_x - \hat{\boldsymbol{v}}\hat{\beta}_v) + \frac{1}{n}\lambda \hat{d}\mathbf{A}\hat{\boldsymbol{v}}\right],\tag{S8}$$

$$\hat{\boldsymbol{v}} = \left(\hat{\beta}_v^2 + \lambda \hat{d}^2\right)^{-1} \left[\hat{\beta}_v(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_x - \hat{\boldsymbol{u}}\hat{\beta}_u) + \frac{1}{n}\lambda \hat{d}\mathbf{A}^{\mathsf{T}}\hat{\boldsymbol{u}}\right],\tag{S9}$$

with constraints

$$\hat{\boldsymbol{u}}^{\top}\hat{\boldsymbol{u}} = n \quad \text{and} \quad \hat{\boldsymbol{v}}^{\top}\hat{\boldsymbol{v}} = n \tag{S10}$$

where  $\widehat{W} = (X, \hat{u}, \hat{v})$ . Denote  $(\hat{d}^{(t)}, \hat{u}^{(t)}, \hat{v}^{(t)}, \hat{\beta}^{(t)})$  as the estimations from the *t*-th iteration. Combining (S8)-(S10) and substituting the corresponding estimates from the previous updates, we obtain each update step in each iteration.

## S2. SUPERCENT FOR AN UNDIRECTED NETWORK

When the network is undirected with the eigenvector centrality, it can be represented by a symmetric matrix A. Denote u as the eigenvector centrality. The objective function of Super-CENT estimation is a special case of (S1), i.e.,

$$(\hat{d}, \hat{\boldsymbol{u}}, \hat{\boldsymbol{\beta}}_x, \hat{\beta}_u) := \underset{\substack{\boldsymbol{\beta}_x, \beta_u\\d, \|\boldsymbol{u}\|_2 = \sqrt{n}}}{\operatorname{arg\,min}} \frac{1}{n} \| \mathbf{y} - \mathbf{X} \boldsymbol{\beta}_x - \boldsymbol{u} \beta_u \|_2^2 + \frac{\lambda}{n^2} \| \mathbf{A} - d\boldsymbol{u} \boldsymbol{u}^\top \|_F^2.$$
(S11)

To solve (S11), we adopt a similar strategy as Algorithm 2 – a block gradient descent algorithm by updating  $(\hat{d}, \hat{u}, \hat{\beta})$  iteratively until convergence, where  $\hat{\beta} = (\hat{\beta}_x^{\top}, \hat{\beta}_u)^{\top}$ . The initialization can be obtained from the eigen decomposition of A. Given a tuning parameter  $\lambda$ , Algorithm S1 describes the algorithm for a symmetric matrix A. Similarly, the tuning parameters  $\lambda$  can be chosen using cross-validation and others as described in Section 4.4.

Algorithm S1: SuperCENT( $\mathbf{A}, \mathbf{X}, \mathbf{y}, \lambda$ ) to solve (S11) for a symmetric  $\mathbf{A}$ .

**Result**:  $\hat{d}$ ,  $\hat{u}$  and  $\hat{\beta}$  **Input**: the observed network  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , the design matrix  $\mathbf{X} \in \mathbb{R}^{n \times p}$ , the response vector  $\mathbf{y} \in \mathbb{R}^{n}$ , the tuning penalty parameter  $\lambda$ , the tolerance parameter  $\rho > 0$ , the maximum number of iteration T; Initiate  $(d^{(0)}, \boldsymbol{u}^{(0)}) = \arg \min_{d, \parallel} \boldsymbol{u}_{\parallel_2 = \sqrt{n}} \parallel A - d\boldsymbol{u} \boldsymbol{u}^{\top} \parallel_F^2, t = 1$ ; while  $\parallel \boldsymbol{P}_{\boldsymbol{u}^{(t)}} - \boldsymbol{P}_{\boldsymbol{u}^{(t-1)}} \parallel_2 > \rho$  and t < T do 1.  $\mathbf{W}^{(t-1)} = (\mathbf{X}, \boldsymbol{u}^{(t-1)})$ 2.  $\beta^{(t)} = (\mathbf{W}^{(t-1)\top} \mathbf{W}^{(t-1)})^{-1} \mathbf{W}^{(t-1)\top} \mathbf{Y}$ 3.  $d^{(t)} = \boldsymbol{u}^{(t-1)\top} \mathbf{A} \boldsymbol{u}^{(t-1)} / n^2$ 4.  $\boldsymbol{u}^{(t)} = ((\beta_u^{(t)})^2 \mathbf{I} + \lambda (d^{(t)})^2 \mathbf{I} - 2\lambda d^{(t)} \mathbf{A} / n)^{-1} \beta_u^{(t)} (\mathbf{y} - \mathbf{X} \boldsymbol{\beta}_x^{(t)})$ 5. Normalize  $\boldsymbol{u}^{(t)}$  such that  $\parallel \boldsymbol{u}^{(t)} \parallel_2 = \sqrt{n}$ 6.  $t \leftarrow t + 1$ end

The derivation of Algorithm 1 is similar to Section S1. Denote W = (X, u),  $\beta = (\beta_x, \beta_u)$ , and  $\mathcal{L}_{sym} := \frac{1}{n} \|\mathbf{y} - \mathbf{X}\beta_x - u\beta_u\|_2^2 + \frac{\lambda}{n^2} \|\mathbf{A} - duu^{\top}\|_F^2$  where the subscript *sym* denotes the objective function for a symmetric matrix  $\mathbf{A}$ . Given  $\lambda$ , we minimize the objection function (S1) by setting the partial derivatives with all the parameters as zero. The partial derivatives are as follows.

$$\frac{\partial \mathcal{L}_{sym}}{\partial \boldsymbol{\beta}} = -\frac{2}{n} \boldsymbol{W}^{\mathsf{T}} (\boldsymbol{y} - \boldsymbol{W} \boldsymbol{\beta}), \qquad (S12)$$

$$\frac{\partial \mathcal{L}_{sym}}{\partial d} = 2\lambda d - 2\lambda \boldsymbol{u}^{\mathsf{T}} \mathbf{A} \boldsymbol{u} / n^2, \tag{S13}$$

$$\frac{\partial \mathcal{L}_{sym}}{\partial \boldsymbol{u}} = -\frac{2}{n}\beta_u(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}_x - \boldsymbol{u}\beta_u) + \frac{2}{n}\lambda d^2\boldsymbol{u} - \frac{4}{n^2}\lambda d\mathbf{A}\boldsymbol{u}.$$
 (S14)

Setting the partial derivatives above as zero yields the estimates

$$\widehat{\boldsymbol{\beta}} = (\widehat{\boldsymbol{W}}^{\top} \widehat{\boldsymbol{W}})^{-1} \widehat{\boldsymbol{W}}^{\top} \boldsymbol{y}, \qquad (S15)$$

$$\hat{d} = \frac{1}{n^2} \hat{\boldsymbol{u}}^{\mathsf{T}} \boldsymbol{A} \hat{\boldsymbol{u}}, \tag{S16}$$

$$\hat{\boldsymbol{u}} = \left(\hat{\beta}_u^2 \boldsymbol{I} + \lambda \hat{d}^2 \boldsymbol{I} - \frac{2}{n} \lambda \hat{d} \boldsymbol{A}\right)^{-1} \hat{\beta}_u (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}_x),$$
(S17)

with a constraint

$$\hat{\boldsymbol{u}}^{\top}\hat{\boldsymbol{u}}=\boldsymbol{n},\tag{S18}$$

where  $\widehat{W} = (X, \widehat{u})$ . Similarly, denote  $(\widehat{d}^{(t)}, \widehat{u}^{(t)}, \widehat{\beta}^{(t)})$  as the estimations from the *t*-th iteration. Taking together (S15)-(S18) and substituting corresponding estimates from the previous update, we obtain each update step in each iteration.

### S3. EXPLICIT MATHEMATICAL EXPRESSIONS OF COVARIANCES IN THEOREMS 1 AND 2

Let K be the  $n^2 \times n^2$  commutation matrix such that  $vec(E^{\top}) = K vec(E)$  and  $\otimes$  denote the Kronecker product. Recall  $\tilde{u} = (I - P_X)u$ ,  $\tilde{v} = (I - P_X)v$  and  $C_{\tilde{u}\tilde{v}} = \begin{pmatrix} \tilde{u}^\top \tilde{u} \ \tilde{u}^\top \tilde{v} \\ \tilde{u}^\top \tilde{v} \ \tilde{v}^\top \tilde{v} \end{pmatrix}$ .

# S3.1. Specific form of $\Sigma^{ts}$ and $C^{ts}$ in Theorem 1

The asymptotic variance of the two-stage estimator in (11)  $\Sigma^{ts} = C^{ts} \begin{pmatrix} \sigma_y^2 I_n & \mathbf{0}_{n \times n^2} \\ \mathbf{0}_{n^2 \times n} & \sigma_a^2 I_{n^2} \end{pmatrix} C^{ts^{\top}}.$ 

Denote  $C^{ts} = \begin{pmatrix} \mathbf{0}_{n \times n} & \mathbf{C}_{12} \\ \mathbf{0}_{n \times n} & \mathbf{C}_{22}^{ts} \\ \mathbf{0}_{n^2 \times n} & \mathbf{C}_{32}^{ts} \\ \mathbf{C}_{41}^{ts} & \mathbf{C}_{42}^{ts} \\ \mathbf{C}_{51}^{ts} & \mathbf{C}_{52}^{ts} \\ \mathbf{C}_{51}^{ts} & \mathbf{C}_{52}^{ts} \end{pmatrix}$  whose specific forms are as follows. The matrices related to  $\hat{u}^{ts}$ 

$$\begin{pmatrix} \boldsymbol{C}_{12}^{ts} \\ \boldsymbol{C}_{22}^{ts} \end{pmatrix} = (dn)^{-1} \begin{pmatrix} \boldsymbol{v}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P} \boldsymbol{u}) \\ (\boldsymbol{u}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P} \boldsymbol{v})) \boldsymbol{K} \end{pmatrix},$$
(S19)

the matrix related to  $\hat{A}^{ts}$  is

$$\boldsymbol{C}_{32}^{ts} = (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \otimes \boldsymbol{P}_{\boldsymbol{u}} + \boldsymbol{P}_{\boldsymbol{v}} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) + \boldsymbol{P}_{\boldsymbol{v}} \otimes \boldsymbol{P}_{\boldsymbol{u}},$$
(S20)

the matrices related to  $\hat{\beta}_{u}^{ts}$  and  $\hat{\beta}_{v}^{ts}$  are

$$\begin{pmatrix} \boldsymbol{C}_{41}^{ts} \, \boldsymbol{C}_{42}^{ts} \\ \boldsymbol{C}_{51}^{ts} \, \boldsymbol{C}_{52}^{ts} \end{pmatrix} = \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} \begin{pmatrix} -\beta_u \boldsymbol{I}_n & -\beta_v \boldsymbol{I}_n & \boldsymbol{I}_n \end{pmatrix} \begin{pmatrix} \boldsymbol{C}_{11}^{ts} \, \boldsymbol{C}_{12}^{ts} \\ \boldsymbol{C}_{21}^{ts} \, \boldsymbol{C}_{22}^{ts} \\ \boldsymbol{I}_n \, \boldsymbol{0}_{n \times n^2} \end{pmatrix}, \qquad (S21)$$

and the matrices related to  $\hat{\boldsymbol{\beta}}_x^{ts}$  are

$$(\boldsymbol{C}_{61}^{ts} \boldsymbol{C}_{62}^{ts}) = (\boldsymbol{X}^{\top} \boldsymbol{X})^{-1} \boldsymbol{X}^{\top} (-\beta_{u} \boldsymbol{I}_{n} - \beta_{v} \boldsymbol{I}_{n} - \boldsymbol{u} - \boldsymbol{v} \boldsymbol{I}_{n}) \begin{pmatrix} \boldsymbol{C}_{11}^{ts} & \boldsymbol{C}_{12}^{ts} \\ \boldsymbol{C}_{21}^{ts} & \boldsymbol{C}_{22}^{ts} \\ \boldsymbol{C}_{41}^{ts} & \boldsymbol{C}_{42}^{ts} \\ \boldsymbol{C}_{51}^{ts} & \boldsymbol{C}_{52}^{ts} \\ \boldsymbol{I}_{n} & \boldsymbol{0}_{n \times n^{2}} \end{pmatrix}.$$
(S22)

## S3.2. Specific form of $\Sigma$ and C in Theorem 2

The asymptotic variance of the SuperCENT estimator in (25)  $\boldsymbol{\Sigma} = \boldsymbol{C} \begin{pmatrix} \sigma_y^2 \boldsymbol{I}_n & \boldsymbol{0}_{n \times n^2} \\ \boldsymbol{0}_{n^2 \times n} & \sigma_c^2 \boldsymbol{I}_{n^2} \end{pmatrix} \boldsymbol{C}^{\top}.$ 

Denote 
$$C = \begin{pmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \\ C_{31} & C_{32} \\ C_{41} & C_{42} \\ C_{51} & C_{52} \\ C_{61} & C_{62} \end{pmatrix}$$
 whose specific forms are given as follows. The matrices related

to  $\hat{u}$  and  $\hat{v}$  are

$$\begin{pmatrix} \boldsymbol{C}_{11} \ \boldsymbol{C}_{12} \\ \boldsymbol{C}_{21} \ \boldsymbol{C}_{22} \end{pmatrix} = (\lambda d^2)^{-1} \begin{pmatrix} \boldsymbol{I}_{2n} - (\lambda d^2 + \beta_u^2 + \beta_v^2)^{-1} \begin{pmatrix} \beta_u^2 & \beta_u \beta_v \\ \beta_u \beta_v & \beta_v^2 \end{pmatrix} \otimes \begin{pmatrix} \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \end{pmatrix} \end{pmatrix} \\ \begin{pmatrix} \beta_u \begin{pmatrix} \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \end{pmatrix} & \lambda dn^{-1} \boldsymbol{v}^\top \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) \\ \beta_v \begin{pmatrix} \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \end{pmatrix} & \lambda dn^{-1} (\boldsymbol{u}^\top \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}})) \boldsymbol{K} \end{pmatrix}$$
(S23)

where  $P_{(Xuv)}$  is the projection matrix that projects onto the column space of (X, u, v), the matrices related to  $\hat{A}$  are

$$(\boldsymbol{C}_{31} \boldsymbol{C}_{32}) = \left( d(\boldsymbol{v} \otimes \boldsymbol{I}_n) \ d\boldsymbol{K}(\boldsymbol{u} \otimes \boldsymbol{I}_n) \ \boldsymbol{P}_{\boldsymbol{v}} \otimes \boldsymbol{P}_{\boldsymbol{u}} \right) \begin{pmatrix} \boldsymbol{C}_{11} \ \boldsymbol{C}_{12} \\ \boldsymbol{C}_{21} \ \boldsymbol{C}_{22} \\ \boldsymbol{0}_{n^2 \times n} \ \boldsymbol{I}_{n^2} \end{pmatrix}, \quad (S24)$$

the matrices related to  $\hat{\beta}_u$  and  $\hat{\beta}_v$  are

$$\begin{pmatrix} \boldsymbol{C}_{41} \ \boldsymbol{C}_{42} \\ \boldsymbol{C}_{51} \ \boldsymbol{C}_{52} \end{pmatrix} = \boldsymbol{C}_{\boldsymbol{\tilde{u}}\boldsymbol{\tilde{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} (-\beta_u \boldsymbol{I}_n \ -\beta_v \boldsymbol{I}_n \ \boldsymbol{I}_n) \begin{pmatrix} \boldsymbol{C}_{11} \ \boldsymbol{C}_{12} \\ \boldsymbol{C}_{21} \ \boldsymbol{C}_{22} \\ \boldsymbol{I}_n \ \boldsymbol{0}_{n \times n^2} \end{pmatrix}, \quad (S25)$$

and the matrices related to  $\hat{\boldsymbol{\beta}}_{x}$  are

$$(\boldsymbol{C}_{61} \, \boldsymbol{C}_{62}) = (\boldsymbol{X}^{\top} \boldsymbol{X})^{-1} \boldsymbol{X}^{\top} (-\beta_u \boldsymbol{I}_n - \beta_v \boldsymbol{I}_n - \boldsymbol{u} - \boldsymbol{v} \, \boldsymbol{I}_n) \begin{pmatrix} \boldsymbol{C}_{11} & \boldsymbol{C}_{12} \\ \boldsymbol{C}_{21} & \boldsymbol{C}_{22} \\ \boldsymbol{C}_{41} & \boldsymbol{C}_{42} \\ \boldsymbol{C}_{51} & \boldsymbol{C}_{52} \\ \boldsymbol{I}_n & \boldsymbol{0}_{n \times n^2} \end{pmatrix}.$$
(S26)



FIGURE S1.—Consistent regime of two-stage: The boxplot of  $\log_{10}(l(\hat{u}, u))$  for the four estimators across different  $\sigma_a, \sigma_y$  and  $\beta_u$  with fixed d = 1 and  $\beta_v = 1$  where  $l(\hat{u}, u) = \|P_{\hat{u}} - P_{u}\|_2^2$ . The super-imposed red symbols show the theoretical rates of the two-stage in Corollary 1 and SuperCENT in Corollary 4.

#### S4. MORE ON SIMULATIONS

In this section, we show more simulation results that are deferred from Section 6. Section S4.1 is for the consistent regime of two-stage and Section S4.2 is for the inconsistent regime of two-stage.

### S4.1. Consistent regime of two-stage

We present the consistent regime of the two-stage procedure, i.e., when the network noiseto-signal ratio  $\kappa = \frac{\sigma_a^2}{d^2n} \rightarrow 0$ . For the consistency of the two-stage procedure, we keep  $\kappa < 2^{-12}$ (recall d = 1 and  $n = 2^8$ ) by varying  $\sigma_a \in 2^{-4,-2}$ . The configuration is exactly the same as the consistent regime except for  $\sigma_a$ . We expect that SuperCENT improves over the two-stage in terms of both estimation and inference.

Figure S1 shows the boxplot of the logarithm of  $l(\hat{u}, u)$  across different  $\sigma_a, \sigma_y$  and  $\beta_u$  with d = 1 and  $\beta_v = 1$ . The panel structure remains the same: the rows correspond to  $\log_2(\sigma_a)$  and the columns correspond to  $\log_2(\beta_u)$ ; for each panel, the x-axis is  $\log_2(\sigma_y)$  and the y-axis is  $\log_{10}(l(\hat{u}, u))$ . The super-imposed red symbols show the theoretical rates of  $\hat{u}^{ts}$  in Corollary 1 and that of  $\hat{u}$  in Corollary 4. Some messages are consistent between Figure S1 and Figure 3: i) SuperCENT- $\lambda_0$  performs the best, Two-stage is always the worst, SuperCENT- $\hat{\lambda}_{cv}$  and SuperCENT- $\hat{\lambda}_0$  are in-between; ii) The two-stage estimator performs the same no matter how large  $\sigma_y$  and  $\beta_u$  are, and it has smaller error with smaller  $\sigma_a$ ; ii) The improvement of SuperCENT over two-stage is stronger with larger  $\beta_u, \sigma_a$  and smaller  $\sigma_y$ , which aligns with Remark 10. For the comparison of the three SuperCENT- $\hat{\lambda}_0$  are sometimes worse than SuperCENT- $\lambda_0$ , but still better than two-stage.

Figure S2 shows the boxplot of logarithm of  $l(\hat{v}, v)$  across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$ . Other than a noticeable improvement when  $\beta_u \in 2^{0,2}$  and  $\sigma_y = 2^{-4}$ , the improvement of SuperCENT over two-stage is negligible in other settings when  $\beta_u > \beta_v$  and the network signal-to-noise



FIGURE S2.—Consistent regime of two-stage: The boxplot of  $\log_{10}(l(\hat{v}, v))$  for four estimators across different  $\sigma_a, \sigma_y$  and  $\beta_u$  with fixed d = 1 and  $\beta_v = 1$  where  $l(\hat{v}, v) = \|P_{\hat{v}} - P_v\|_2^2$ . The super-imposed red points show the theoretical rates of the two-stage and SuperCENT algorithms.



FIGURE S3.—Consistent regime of two-stage: The boxplot of  $\log_{10}(l(\mathbf{A}, \mathbf{A}_0))$  for four estimators across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$  with fixed d = 1 and  $\beta_v = 1$  where  $l(\hat{\mathbf{A}}, \mathbf{A}_0) = \|\hat{\mathbf{A}} - \mathbf{A}_0\|_F^2 / \|\mathbf{A}_0\|_F^2$ . The super-imposed red points show the theoretical rates of the two-stage and SuperCENT algorithms.

ratio is relatively large. Comparing Figures S1 and S2, the supervised effect to  $\hat{u}$  increases as  $\beta_u$  increases while the supervised effect to  $\hat{v}$  shrinks as  $\beta_u$  increases. This phenomena aligns with the nonexchangeable roles u and v as in Remark 13 – noted that here we have  $\beta_v \leq \beta_u$  by fixing  $\beta_v = 1$  and varying  $\beta_u \in 2^{0,2,4}$ . Similarly, when  $\beta_v \gg \beta_u$  we should expect the improvement in estimating v to increase.

Figure S3 shows the boxplot of logarithm of  $l(\hat{A}, A_0)$  across different  $\sigma_a, \sigma_y$  and  $\beta_u$ . With the improvement in estimating u and v by SuperCENT as shown in Figures S1 and S2, it is not surprising to see the improvement in estimating  $A_0$ . The trend of improvement is similar to that of u in Figure S1. Comparing Figures S1,S2 and S3 altogether, when  $\beta_u = 2^0$  and  $\sigma_y = 2^{-4}$ ,



FIGURE S4.—Consistent regime of two-stage: The boxplot of the bias  $\hat{\beta}_u - \beta_u$  across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$  with fixed d = 1 and  $\beta_v = 1$ . The dashed lines show  $\hat{\beta}_u - \beta_u = 0$ . The super-imposed red points show the median of the bias.



FIGURE S5.—Consistent regime of two-stage: The boxplot of the bias  $\hat{\beta}_v - \beta_v$  across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$  with fixed d = 1 and  $\beta_v = 1$ . The dashed lines show  $\hat{\beta}_v - \beta_v = 0$ . The super-imposed red points show the median of the bias.

SuperCENT improves both the estimation of u and v; when  $\beta_u \in 2^{2,4}$ , SuperCENT improves the estimation of u a lot, and therefore, SuperCENT improves the estimation of  $A_0$  for most ranges of  $\beta_u$ , especially when  $\sigma_y$  is small. The comparison of the three SuperCENT-based methods is similar, with the benchmark SuperCENT- $\lambda_0$  always being the best, SuperCENT- $\hat{\lambda}_{cv}$ and SuperCENT- $\hat{\lambda}_0$  being worse than SuperCENT- $\lambda_0$  sometimes but still not worse than twostage.

Figures S4 and S5 show the boxplots of the bias in estimating  $\beta_u$  and  $\beta_v$ ,  $\hat{\beta}_u - \beta_u$  and  $\hat{\beta}_v - \beta_v$ , respectively across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$ . The performance of all methods are



FIGURE S6.—Consistent regime of two-stage: The boxplot of  $\log_{10}(l(\hat{\beta}_u, \beta_u))$  for four estimators across different  $\sigma_a, \sigma_y$  and  $\beta_u$  with fixed d = 1 and  $\beta_v = 1$  where  $l(\hat{\beta}_u, \beta_u) = (\hat{\beta}_u - \beta_u)^2 / \beta_u^2$ . The super-imposed red points show the theoretical rates of the two-stage in Corollary 2 and SuperCENT algorithms in Corollary 5.



FIGURE S7.—Consistent regime of two-stage: The boxplot of  $\log_{10}(l(\hat{\beta}_v, \beta_v))$  for four estimators across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$  with fixed d = 1 and  $\beta_v = 1$  where  $l(\hat{\beta}_v, \beta_v) = (\hat{\beta}_v - \beta_v)^2 / \beta_v^2$ . The super-imposed red points show the theoretical rates of the two-stage in Corollary 2 and SuperCENT algorithms in Corollary 5.

similar with negligible bias. This is aligned with Remark 7 where we claim that there exists no asymptotic bias when two-stage is consistent.

Figures S6 and S7 show the boxplots of  $l(\hat{\beta}_u, \beta_u)$  and  $l(\hat{\beta}_v, \beta_v)$  respectively across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$ . The performance of all methods are similar in estimating  $\beta_u$  and  $\beta_v$ . This echoes with Corollary 5 where we prove that despite better estimations of u, v by SuperCENT, the rates for estimating  $\beta_u$  and  $\beta_v$  are surprisingly the same for SuperCENT and two-stage.

We now turn our attention to the inference property. Figures S8 and S9 show the empirical coverage and  $\log_{10}$  of the average width, respectively, of the 95% confidence interval for  $\beta_u$  from the five methods. In terms of the empirical coverage, both SuperCENT- $\lambda_0$ -oracle and



O SuperCENT-λ₀-oracle • SuperCENT-λ̂<sub>cv</sub> △ Two-stage-oracle ▼ Two-stage-adhoc ▲ Two-stage

FIGURE S8.—Consistent regime of two-stage: The empirical coverage of  $CI_{\beta_u}$  across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$  with d = 1 and  $\beta_v = 1$ . SuperCENT variants are labelled as circles ( $\circ \bullet$ ) and the two-stage variants are labelled as triangles ( $\Delta \checkmark \blacktriangle$ ). The hollow ones are for oracles and the solid ones are for non-oracles. The dashed lines show the nominal confidence level 0.95.

two-stage-oracle are very close to the nominal level across different settings. Their widths are the same because the same rate of  $\hat{\beta}_u^{ts}$  and  $\hat{\beta}_u$  are shown in Corollary 5 and the true parameters are plugged in. For the non-oracle methods, the two two-stage related methods and SuperCENT- $\hat{\lambda}_{cv}$  are either close to nominal or tend to be conservative. When all of them are conservative, SuperCENT- $\hat{\lambda}_{cv}$  always has the smallest width among the three. Again, the relative narrower width of SuperCENT- $\hat{\lambda}_{cv}$  is more obvious with large  $\sigma_a$ ,  $\beta_u$  or small  $\sigma_y$ . Furthermore, two-stage tends to have wider CI than two-stage-adhoc and is more conservative, which supports Remark 5.

Figures S10 and S11 show the empirical coverage and  $\log_{10}$  of the average width, respectively, of 95% confidence interval for  $\beta_v$  from five methods. In terms of empirical coverage, both SuperCENT- $\lambda_0$ -oracle and two-stage-oracle are very close to the nominal level across different settings. Their widths are the same due to the same rate of  $\hat{\beta}_v^{ts}$  and  $\hat{\beta}_v$  in Corollary 5 with the true parameters plugged in. For the non-oracle methods, two-stage and SuperCENT- $\hat{\lambda}_{cv}$  are mostly close to the nominal level although tend to be conservative when  $\beta_u = 2^4$  and  $\sigma_y = 2^{-4}$ ; while two-stage-adhoc is always close to the nominal. It is worth noting that as  $\beta_u$  increases, the first term of (18) dominates (17) because  $\kappa = \frac{\sigma_a^2}{d^2n} = O(1)$  and  $\beta_u^2 \gg \sigma_y^2$ . Therefore, even with  $\sigma_y^2$  being overestimated, two-stage-adhoc is on par with two-stage-oracle because two-stage-adhoc is more conservative than the two-stage-oracle. The phenomena is more pronounced when  $\sigma_a$  increases (bottom right two panels). In terms of the width, when two-stage and SuperCENT- $\hat{\lambda}_{cv}$  are conservative, their widths are always wider than two-stage-adhoc with two-stage being the widest. Furthermore, the fact of two-stage having wider CI and being more conservative than two-stage having wider CI and being more conservative than two-stage having wider CI and being more conservative than two-stage-adhoc supports Remark 5.

Finally, we investigate the average coverage and the average width of confidence intervals for all the entries  $a_{ij}$  of  $A_0$  respectively. Figure S12 shows that the average coverage probability of all the methods, Average<sub>ij</sub> (CP(CIa<sub>ij</sub>)), achieves the nominal level of 95%. Figure S13 shows



FIGURE S9.—Consistent regime of two-stage:  $\log_{10}$  of the width of  $CI_{\beta_u}$  across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$  with d = 1 and  $\beta_v = 1$ . SuperCENT variants are labelled as circles ( $\circ \bullet$ ) and the two-stage variants are labelled as triangles ( $\Delta \checkmark \blacktriangle$ ). The hollow ones are for oracles and the solid ones are for non-oracles.



FIGURE S10.—Consistent regime of two-stage: Empirical coverage of  $\beta_v$  across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$  with d = 1 and  $\beta_v = 1$ . SuperCENT variants are labelled as circles ( $\circ \bullet$ ) and the two-stage variants are labelled as triangles ( $\Delta \checkmark \blacktriangle$ ). The hollow ones are for oracles and the solid ones are for non-oracles. The dashed lines show the nominal confidence level 0.95.

the  $\log_{10}$  of the average width of the CIs,  $\operatorname{Average}_{ij}(\operatorname{Width}(\operatorname{CI}a_{ij}))$ . SuperCENT- $\lambda_0$ -oracle provides the shortest width among the four methods, followed by SuperCENT- $\hat{\lambda}_{cv}$ , which in turn dominates the two two-stage related methods. Again, the improvement of SuperCENT over two-stage increases as  $\sigma_a$  and  $\beta_u$  increase or  $\sigma_y$  decreases.



FIGURE S11.—Consistent regime of two-stage: Width of  $CI_{\beta_v}$  across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$  with d = 1 and  $\beta_v = 1$ . SuperCENT variants are labelled as circles ( $\circ \bullet$ ) and the two-stage variants are labelled as triangles ( $\Delta \checkmark \blacktriangle$ ). The hollow ones are for oracles and the solid ones are for non-oracles.



FIGURE S12.—Consistent regime of two-stage: The average empirical coverage of  $CI_{a_{ij}}$  across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$  with d = 1 and  $\beta_v = 1$ . SuperCENT variants are labelled as circles ( $\circ \bullet$ ) and the two-stage variants are labelled as triangles ( $\Delta \bullet$ ). The dashed lines show the nominal confidence level 0.95.

### S4.2. Inconsistent regime of two-stage

In the inconsistent regime of two-stage procedure,  $\kappa = \frac{\sigma_a^2}{d^2n} = O(1)$ ,  $\hat{u}^{ts}$  and  $\hat{v}^{ts}$  are inconsistent. Section 6.2 shows that SuperCENT enormously improves over two-stage in terms of estimation of u, v and  $\beta_u$  as well as the inference of  $\beta_u$  and  $A_0$ . In this section, we demonstrate the behaviors of the SuperCENT-based and the two-stage-based estimators of  $\beta_v$  and their corresponding confidence intervals. We also show the plot of the average coverage of the confidence interval of  $A_0$ , which is deferred from the main text.

Figure S14 shows the bias in estimating  $\beta_v$ , i.e.,  $\hat{\beta}_v - \beta_v$ . We observe a bias away from zero in estimating  $\beta_v$  when  $\beta_u$  is large, which is different from the attenuation bias in  $\hat{\beta}_u$ . As in



FIGURE S13.—Consistent regime of two-stage:  $\log_{10}$  of the average width of  $CI_{a_{ij}}$  across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$  with d = 1 and  $\beta_v = 1$ . SuperCENT variants are labelled as circles ( $\circ \bullet$ ) and the two-stage variants are labelled as triangles ( $\Delta \blacktriangle$ ).



FIGURE S14.—Inconsistent regime of two-stage: The boxplot of the bias  $\hat{\beta}_v - \beta_v$  across different  $\sigma_a, \sigma_y$  and  $\beta_u$  with fixed d = 1 and  $\beta_v = 1$ . The dashed lines show  $\hat{\beta}_v - \beta_v = 0$ . The super-imposed red points show the median of the bias.

Remark 7, the bias is away from zero because  $\kappa = 2^{-8,-4} \rightarrow 0$ , u and v are correlated with  $\rho = \frac{2^{-1}}{\sqrt{1.25}}$ , and  $\beta_v = 1 = 2^0 < \{2^2 \cdot 2^{-1.2}, 2^4 \cdot 2^{-1.2}\} = \beta_u \frac{\rho}{(1+\kappa)}$ , leading to plim  $\hat{\beta}_v^{ts} - \beta_v > 0$ . In addition, the bias is larger as  $\beta_u$  increases. SuperCENT in this case can still alleviate the bias but the improvement is not as large as that of  $\beta_u$  since the improvement in estimation of v when  $\beta_u$  is large is limited as shown in Figure 4, thereby slightly improving over the two-stage from the perspective of the estimation bias in Figure S14 as well as the squared error loss in Figure S15.



FIGURE S15.—Inconsistent regime of two-stage: The boxplot of  $\log_{10}(l(\hat{\beta}_v, \beta_v))$  across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$  with d = 1 and  $\beta_v = 1$  where  $l(\hat{\beta}_v, \beta_v) = (\hat{\beta}_v - \beta_v)^2 / \beta_v^2$ . The super-imposed red points show the median of  $\log_{10}(l(\hat{\beta}_v, \beta_v))$ .

Similar to  $\beta_u$ , the bias in the estimation of  $\beta_v$  further affects its confidence interval. Figures S16 and S17 show the empirical coverage and  $\log_{10}$  of the average width, respectively, of the 95% confidence interval for  $\beta_v$ . For the empirical coverage, when  $\sigma_a$  remains small (top panels), all methods remain close to the nominal level, again with different reasons for different methods as discussed for the coverage of  $CI_{\beta_u}$ . When  $\sigma_a$  remains small, the bias in  $\hat{\beta}_v$  is relatively small, leaving a relatively small impact on the coverage. As  $\beta_u$  increases, the oracles tend to under cover since the bias increases. Two-stage is still conservative because  $\sigma_u^2$ is overestimated, covering up the issue of bias. Similar to the phenomena we observed in the consistent regime, the two-stage-adhoc is on par with the two-stage-oracle even with  $\sigma_u^2$  being overestimated and thus below the nominal level. When  $\sigma_a$  increases and the bias is not too severe with  $\beta_u = 2^2$  (the mid-bottom panel), two-stage is still close to the nominal level while two-stage-oracle and two-stage-adhoc are no longer valid. SuperCENT does not improve much from two-stage as the improvement of estimating  $\beta_v$  is limited. The trade-off between  $\beta_u$  and  $\beta_v$  for SuperCENT is desirable – SuperCENT provides valid and shorter intervals for both  $\beta_u$ and  $\beta_v$  if  $\beta_u$  and  $\beta_v$  are similar; if  $\beta_u$  and  $\beta_v$  differ a lot, SuperCENT provides a valid and shorter interval for the larger effect which is more of one's interest.

As for the width of the CI for  $\beta_v$ , Figure S17 shows that when the SuperCENT methods reach the nominal level, the widths are shorter than two-stage.

Finally, Figure S18 shows the average coverage of the confidence intervals of  $A_0$ . All methods remain at the nominal level but tends to be slightly below the nominal coverage as  $\sigma_a$  increases. SuperCENT- $\hat{\lambda}_{cv}$  is the closest to the nominal coverage in all the settings compared to the others. At the same time, SuperCENT- $\hat{\lambda}_{cv}$  retains its advantage in widths compared to the two-stage as shown in 10.

#### S5. ADDITIONAL INFORMATION FOR THE CASE STUDY IN SECTION 7

Table SI provides the country acronyms and full names. Figure S19 shows the average trade volume from 2003 to 2012 among the 24 countries/regions. The arrows reflect the trade directions and the widths represent the volume.



FIGURE S16.—Inconsistent regime of two-stage: Empirical coverage of  $CI_{\beta_v}$  across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$  with d = 1 and  $\beta_v = 1$ . SuperCENT variants are labelled as circles ( $\circ \bullet$ ) and the two-stage variants are labelled as triangles ( $\Delta \checkmark \blacktriangle$ ). The hollow ones are for oracles and the solid ones are for non-oracles. The dashed lines show the nominal confidence level 0.95.



FIGURE S17.—Inconsistent regime of two-stage: Width of  $CI_{\beta_v}$  across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$  with d = 1 and  $\beta_v = 1$ . SuperCENT variants are labelled as circles ( $\circ \bullet$ ) and the two-stage variants are labelled as triangles ( $\Delta \checkmark \blacktriangle$ ). The hollow ones are for oracles and the solid ones are for non-oracles.

#### S6. PROOF OF THE THEORETICAL RESULTS IN SECTION 5

In this section, we collect the proofs of Theorem 1, Corollaries 1, 2, 3 in Section S6.1 and Theorem 2 and Corollaries 4, 5 in Section S6.2.

We begin by providing some basic properties of the Kronecker product and the commutation matrix. The Kronecker product of  $M = (m_{ij}) \in \mathbb{R}^{m \times n}$  and  $N = (n_{ij}) \in \mathbb{R}^{p \times q}$ , denoted by



FIGURE S18.—Inconsistent regime of two-stage: The average empirical coverage  $\operatorname{Average}_{ij}(\operatorname{CP}(\operatorname{CI}_{a_{ij}}))$  across different  $\sigma_a$ ,  $\sigma_y$  and  $\beta_u$  with d = 1 and  $\beta_v = 1$ . SuperCENT variants are labelled as circles  $(\circ \bullet)$  and the two-stage variants are labelled as triangles  $(\land \bullet)$ . The hollow ones are for oracles and the solid ones are for non-oracles. The dashed lines show the nominal confidence level 0.95.

Code	Country	Code	Country	Code	Country
AUS	Australia	JPN	Japan	SGP	Singapore
CAN	Canada	KOR	Korea	SWE	Sweden
CHE	Switzerland	KWT	Kuwait	THA	Thailand
CZE	Czech Republic	MEX	Mexico	ZAF	South Africa
DNK	Denmark	MYS	Malaysia		
GBR	United Kingdom	NOR	Norway		
HKG	Hong Kong	NZL	New Zealand		
HUN	Hungary	PHL	Philippines		
IDN	Indonesia	POL	Poland		
IND	India	SAU	Saudi Arabia		
TABLE SI					

LIST OF COUNTRY ACRONYMS.

 $M \otimes N$ , is defined as

$$\boldsymbol{M} \otimes \boldsymbol{N} = \begin{pmatrix} m_{11} \boldsymbol{N} \cdots m_{1n} \boldsymbol{N} \\ \vdots & \dots & \vdots \\ m_{m1} \boldsymbol{N} \cdots & m_{mn} \boldsymbol{N} \end{pmatrix} \in \mathbb{R}^{mp \times nq}.$$
 (S27)

Denote  $K_{mn} \in \{0,1\}^{mn \times mn}$  as the commutation matrix such that

$$\operatorname{vec}(\boldsymbol{M}^{\top}) = \boldsymbol{K}_{mn} \operatorname{vec}(\boldsymbol{M}).$$
(S28)

We list the following facts about the Kronecker product and the commutation matrix, which are used in the section without specific references. Proofs of these facts can be found in Magnus and Neudecker (1979).

Let  $M \in \mathbb{R}^{m \times n}$ ,  $N \in \mathbb{R}^{p \times q}$ ,  $P \in \mathbb{R}^{n \times t}$ ,  $Q \in \mathbb{R}^{q \times s}$ , and  $Z \in \mathbb{R}^{n \times p}$ . (i)  $(M \otimes N)^{\top} = M^{\top} \otimes N^{\top}$ .



FIGURE \$19.—The average trade volume from 2003 to 2012 among the 24 countries/regions. Each country is in different color.

- (ii)  $(\boldsymbol{M} \otimes \boldsymbol{N})(\boldsymbol{P} \otimes \boldsymbol{Q}) = (\boldsymbol{M}\boldsymbol{P}) \otimes (\boldsymbol{N}\boldsymbol{Q}).$
- (iii)  $\operatorname{vec}(MZN) = (N^{\top} \otimes M) \operatorname{vec}(Z)$ . (iv)  $\operatorname{vec}(MP) = (I \otimes M) \operatorname{vec}(P) = (P^{\top} \otimes I) \operatorname{vec}(M)$ .

- (iv)  $\mathbf{K}_{mn}^{\top} = \mathbf{K}_{nm}$ . (v)  $\mathbf{K}_{mn}^{\top} = \mathbf{K}_{nm}$ . (vi)  $\mathbf{K}_{mn}^{\top} \mathbf{K}_{mn} = \mathbf{K}_{mn}^{\top} \mathbf{K}_{mn} = \mathbf{I}$ . (vii)  $\mathbf{K}_{mp}(\mathbf{M} \otimes \mathbf{N}) \mathbf{K}_{qn} = (\mathbf{N} \otimes \mathbf{M})$ . Equivalently,  $\mathbf{K}_{mp}(\mathbf{M} \otimes \mathbf{N}) = (\mathbf{N} \otimes \mathbf{M}) \mathbf{K}_{qn}$ .
- (viii)  $(M \otimes N) K_{nq}(P \otimes Q) = ((MP) \otimes (NQ)) K_{st} = K_{mp}((NQ) \otimes (MP)).$ (ix)  $\operatorname{tr}(K_{mn}(M \otimes N)) = \operatorname{tr}(MN)$  where  $M, N \in \mathbb{R}^{m \times n}$ .

#### S6.1. Proof of the theoretical results of two-stage

#### S6.1.1. Proof of Theorem 1

PROOF: The naive two-stage procedure first estimates the centralities u and v by the leading left and right singular vectors, rescaled to have norm  $\sqrt{n}$  and denoted as  $\hat{u}^{ts}$ ,  $\hat{v}^{ts}$ , from the SVD on the observed adjacency matrix A, and then performs ordinary least square (OLS) regression of y on X and  $\hat{u}^{ts}, \hat{v}^{ts},$ treating  $\hat{u}^{ts}, \hat{v}^{ts}$  as given covariates. It is, therefore, equivalent to solve the following two optimization problems sequentially,

$$(\hat{d}^{ts}, \hat{\boldsymbol{u}}^{ts}, \hat{\boldsymbol{v}}^{ts}) := \underset{d, \|\boldsymbol{u}\| = \|\boldsymbol{v}\| = \sqrt{n}}{\arg\min} \|\boldsymbol{A} - d\boldsymbol{u}\boldsymbol{v}^{\top}\|_{F}^{2},$$
(S29a)

$$\hat{\boldsymbol{\beta}}^{ts} = ((\hat{\boldsymbol{\beta}}_x^{ts})^\top, \hat{\boldsymbol{\beta}}_u^{ts}, \hat{\boldsymbol{\beta}}_v^{ts})^\top := \underset{\boldsymbol{\beta}_x, \beta_u, \beta_v}{\arg\min} \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}_x - \hat{\boldsymbol{u}}^{ts}\beta_u - \hat{\boldsymbol{v}}^{ts}\beta_v\|_2^2.$$
(S29b)

Denote  $\mathcal{L}_1 := \|\boldsymbol{A} - d\boldsymbol{u}\boldsymbol{v}^\top\|_F^2$  and  $\mathcal{L}_2 := \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}_x - \hat{\boldsymbol{u}}^{ts}\boldsymbol{\beta}_u - \hat{\boldsymbol{v}}^{ts}\boldsymbol{\beta}_v\|_2^2$ . In the following, we refer to solving (S29a) as the first stage and (S29b) as the second stage.

(1) *First stage* We minimize the objection function (S29a) of the first stage by setting the partial derivatives with all the parameters as zero. The partial derivatives are as follows.

$$\frac{\partial \mathcal{L}_1}{\partial d} = 2dn^2 - 2\boldsymbol{u}^\top \mathbf{A}\boldsymbol{v},\tag{S30}$$

$$\frac{\partial \mathcal{L}_1}{\partial \boldsymbol{u}} = 2d^2 n \boldsymbol{u} - 2d \mathbf{A} \boldsymbol{v}, \tag{S31}$$

$$\frac{\partial \mathcal{L}_1}{\partial \boldsymbol{v}} = 2d^2 n \boldsymbol{v} - 2d\mathbf{A}^\top \boldsymbol{u}.$$
(S32)

By setting the partial derivatives above as zero, we have the estimates

$$\hat{d}^{ts} = \hat{\boldsymbol{u}}^{ts\top} \boldsymbol{A} \hat{\boldsymbol{v}}^{ts} / n^2, \qquad (S33)$$

$$\hat{d}^{ts}\hat{\boldsymbol{u}}^{ts} = \mathbf{A}\hat{\boldsymbol{v}}^{ts}/n, \tag{S34}$$

$$\hat{d}^{ts}\hat{\boldsymbol{v}}^{ts} = \mathbf{A}^{\top}\hat{\boldsymbol{u}}^{ts}/n \tag{S35}$$

with constraints

$$\hat{\boldsymbol{u}}^{ts\top}\hat{\boldsymbol{u}}^{ts} = n \quad \text{and} \quad \hat{\boldsymbol{v}}^{ts\top}\hat{\boldsymbol{v}}^{ts} = n.$$
 (S36)

Together they lead to the first order expansion

$$(\boldsymbol{u} + \delta_{\boldsymbol{u}}^{ts})^{\top} (\boldsymbol{u} + \delta_{\boldsymbol{u}}^{ts}) \approx n,$$
(S37)

$$(\boldsymbol{v} + \delta_{\boldsymbol{v}}^{ts})^{\top} (\boldsymbol{v} + \delta_{\boldsymbol{v}}^{ts}) \approx n,$$
 (S38)

$$d + \delta_d^{ts} \approx (\boldsymbol{u} + \delta_{\boldsymbol{u}}^{ts})^\top (d\boldsymbol{u}\boldsymbol{v}^\top + \boldsymbol{E})(\boldsymbol{v} + \delta_{\boldsymbol{v}}^{ts})/n^2,$$
(S39)

$$(d + \delta_d^{ts})(\boldsymbol{u} + \delta_{\boldsymbol{u}}^{ts}) - (d\boldsymbol{u}\boldsymbol{v}^\top + \boldsymbol{E})(\boldsymbol{v} + \delta_{\boldsymbol{v}}^{ts})/n \approx 0,$$
(S40)

$$(d + \delta_d^{ts})(\boldsymbol{v} + \delta_{\boldsymbol{v}}^{ts}) - (d\boldsymbol{v}\boldsymbol{u}^\top + \boldsymbol{E}^\top)(\boldsymbol{u} + \delta_{\boldsymbol{u}}^{ts})/n \approx 0.$$
(S41)

After simplification and dropping the second order terms, (S37) and (S38) become

$$\boldsymbol{u}^{\top} \delta_{\boldsymbol{u}}^{ts} / n \approx 0 \quad \text{and} \quad \boldsymbol{v}^{\top} \delta_{\boldsymbol{v}}^{ts} / n \approx 0.$$
 (S42)

Further simplifying the rest and using (S42), we have

$$\delta_d^{ts} \approx d\boldsymbol{u}^\top \delta_{\boldsymbol{u}}^{ts} + d\boldsymbol{v}^\top \delta_{\boldsymbol{v}}^{ts} + \boldsymbol{u}^\top \boldsymbol{E} \boldsymbol{v} / n^2 \approx \boldsymbol{u}^\top \boldsymbol{E} \boldsymbol{v} / n^2,$$
(S43)

$$nd\delta_{\boldsymbol{u}}^{ts} \approx (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}})\boldsymbol{E}\boldsymbol{v},\tag{S44}$$

$$nd\delta_{\boldsymbol{v}}^{ts} \approx (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}})\boldsymbol{E}^{\top}\boldsymbol{u}.$$
(S45)

Combining (S44) and (S45), we obtain

$$\begin{pmatrix} \delta_{\boldsymbol{u}}^{ts} \\ \delta_{\boldsymbol{v}}^{ts} \end{pmatrix} \approx (dn)^{-1} \begin{pmatrix} \boldsymbol{v}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) \\ (\boldsymbol{u}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}})) K \end{pmatrix} \operatorname{vec}(\boldsymbol{E}) \stackrel{def}{=} \begin{pmatrix} \boldsymbol{C}_{12}^{ts} \\ \boldsymbol{C}_{22}^{ts} \end{pmatrix} \operatorname{vec}(\boldsymbol{E}).$$
(S46)

Let  $\hat{A}^{ts} = \hat{d}^{ts} \hat{u}^{ts} \hat{v}^{ts \top}$ . Then the first order expansion leads to

$$\operatorname{vec}\left(\delta_{\boldsymbol{A}}^{ts}\right) = \operatorname{vec}\left(\hat{d}^{ts}\hat{\boldsymbol{u}}^{ts}\hat{\boldsymbol{v}}^{ts\top}\right) - \operatorname{vec}\left(d\boldsymbol{u}\boldsymbol{v}^{\top}\right)$$
(S47)

$$= \operatorname{vec}(d + \delta_d^{ts})(\boldsymbol{u} + \delta_{\boldsymbol{u}}^{ts})(\boldsymbol{v} + \delta_{\boldsymbol{v}}^{ts})^{\top} - \operatorname{vec}(d\boldsymbol{u}\boldsymbol{v}^{\top})$$
(S48)

$$= d \operatorname{vec}(\boldsymbol{u} \delta_{\boldsymbol{v}}^{ts^{\top}}) + d \operatorname{vec}(\delta_{\boldsymbol{u}}^{ts} \boldsymbol{v}^{\top}) + \delta_{d}^{ts} \operatorname{vec}(\boldsymbol{u} \boldsymbol{v}^{\top})$$
(S49)

$$= d\boldsymbol{K}\operatorname{vec}(\delta_{\boldsymbol{v}}^{ts}\boldsymbol{u}^{\top}) + d\operatorname{vec}(\delta_{\boldsymbol{u}}^{ts}\boldsymbol{v}^{\top}) + \delta_{d}^{ts}\operatorname{vec}(\boldsymbol{u}\boldsymbol{v}^{\top})$$
(S50)

$$= d\boldsymbol{K}(\boldsymbol{u} \otimes \boldsymbol{I})\delta_{\boldsymbol{v}}^{ts} + d(\boldsymbol{v} \otimes \boldsymbol{I})\delta_{\boldsymbol{u}}^{ts} + \delta_{d}^{ts}\operatorname{vec}(\boldsymbol{u}\boldsymbol{v}^{\top})$$
(S51)

$$= n^{-1} \boldsymbol{K}(\boldsymbol{u} \otimes \boldsymbol{I}) \left( \boldsymbol{u}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \right) \boldsymbol{K} \operatorname{vec}(\boldsymbol{E})$$
(S52)

$$+n^{-1}(\boldsymbol{v}\otimes\boldsymbol{I})\left(\boldsymbol{v}^{\top}\otimes(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{u}})\right)\operatorname{vec}(\boldsymbol{E})+n^{-2}\boldsymbol{u}^{\top}\boldsymbol{E}\boldsymbol{v}\operatorname{vec}(\boldsymbol{u}\boldsymbol{v}^{\top})$$
(S53)

$$= [(I - Pv) \otimes Pu + Pv \otimes (I - Pu) + Pv \otimes Pu] \operatorname{vec}(E)$$
(S54)

$$\stackrel{def}{=} \boldsymbol{C}_{32} \operatorname{vec}(\boldsymbol{E}). \tag{S55}$$

Putting together (S46) and (S55)

$$\begin{pmatrix} \hat{\boldsymbol{u}}^{ts} - \boldsymbol{u} \\ \hat{\boldsymbol{v}}^{ts} - \boldsymbol{v} \\ \operatorname{vec} \left( \hat{\boldsymbol{A}}^{ts} - \boldsymbol{A}_0 \right) \end{pmatrix} \approx \begin{pmatrix} \boldsymbol{C}_{12}^{ts} \\ \boldsymbol{C}_{22}^{ts} \\ \boldsymbol{C}_{32}^{ts} \end{pmatrix} \operatorname{vec}(\boldsymbol{E}).$$
(S56)

This finishes the proof of the first stage.

(2) Second stage. For the second stage, we plug in  $\hat{u}^{ts}$  and  $\hat{v}^{ts}$  from the first stage into the objection function (S29b) of the second stage. We then minimize the objection function by setting the partial derivatives with the regression coefficients as zero,

$$\frac{\partial \mathcal{L}_2}{\partial \beta_u} = -\frac{2}{n} \hat{\boldsymbol{u}}^{ts \top} (\boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta}_x - \hat{\boldsymbol{u}}^{ts} \beta_u - \hat{\boldsymbol{v}}^{ts} \beta_v) = 0, \qquad (S57)$$

$$\frac{\partial \mathcal{L}_2}{\partial \beta_v} = -\frac{2}{n} \hat{\boldsymbol{v}}^{ts\top} (\boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta}_x - \hat{\boldsymbol{u}}^{ts} \beta_u - \hat{\boldsymbol{v}}^{ts} \beta_v) = 0, \qquad (S58)$$

$$\frac{\partial \mathcal{L}_2}{\partial \boldsymbol{\beta}_x} = -\frac{2}{n} \boldsymbol{X}^\top (\boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta}_x - \hat{\boldsymbol{u}}^{ts} \boldsymbol{\beta}_u - \hat{\boldsymbol{v}}^{ts} \boldsymbol{\beta}_v) = 0.$$
(S59)

After some algebra, we obtain the estimates as

$$\hat{\beta}_{u}^{ts} = \frac{1}{n} \hat{\boldsymbol{u}}^{ts\top} (\boldsymbol{y} - \boldsymbol{X} \hat{\boldsymbol{\beta}}_{x}^{ts} - \hat{\boldsymbol{v}}^{ts} \hat{\beta}_{v}^{ts}), \qquad (S60)$$

$$\hat{\beta}_{v}^{ts} = \frac{1}{n} \hat{\boldsymbol{v}}^{ts\top} (\boldsymbol{y} - \boldsymbol{X} \hat{\boldsymbol{\beta}}_{x}^{ts} - \hat{\boldsymbol{u}}^{ts} \hat{\boldsymbol{\beta}}_{u}^{ts}),$$
(S61)

$$\hat{\boldsymbol{\beta}}_{x}^{ts} = (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}(\boldsymbol{y} - \hat{\boldsymbol{u}}^{ts}\hat{\boldsymbol{\beta}}_{u}^{ts} - \hat{\boldsymbol{v}}^{ts}\hat{\boldsymbol{\beta}}_{v}^{ts}),$$
(S62)

which leads to the first order expansion

$$\beta_{u} + \delta_{\beta_{u}}^{ts} \approx \frac{1}{n} (\boldsymbol{u} + \delta_{\boldsymbol{u}}^{ts})^{\top} (\boldsymbol{X}\boldsymbol{\beta}_{x} + \boldsymbol{u}\beta_{u} + \boldsymbol{v}\beta_{v} + \boldsymbol{\epsilon} - \boldsymbol{X}\boldsymbol{\beta}_{x} - \boldsymbol{X}\delta_{\boldsymbol{\beta}_{x}}^{ts} - \boldsymbol{v}\beta_{v} - \boldsymbol{v}\delta_{\beta_{v}}^{ts} - \delta_{\boldsymbol{v}}^{ts}\beta_{v}) (S63)$$

$$\beta_{v} + \delta_{\beta_{v}}^{ts} \approx \frac{1}{n} (\boldsymbol{v} + \delta_{\boldsymbol{v}}^{ts})^{\top} (\boldsymbol{X} \boldsymbol{\beta}_{x} + \boldsymbol{u} \beta_{u} + \boldsymbol{v} \beta_{v} + \boldsymbol{\epsilon} - \boldsymbol{X} \boldsymbol{\beta}_{x} - \boldsymbol{X} \delta_{\boldsymbol{\beta}_{x}}^{ts} - \boldsymbol{u} \beta_{u} - \boldsymbol{u} \delta_{\beta_{u}}^{ts} - \delta_{\boldsymbol{u}}^{ts} \beta_{u}) (S64)$$
  
$$\boldsymbol{\beta}_{x} + \delta_{\boldsymbol{\beta}}^{ts} \approx (\boldsymbol{X}^{\top} \boldsymbol{X})^{-1} \boldsymbol{X}^{\top} (\boldsymbol{X} \boldsymbol{\beta}_{x} + \boldsymbol{\epsilon} - \boldsymbol{u} \delta_{\beta_{u}}^{ts} - \delta_{\boldsymbol{u}}^{ts} \beta_{u} - \boldsymbol{v} \delta_{\beta_{v}}^{ts} - \delta_{\boldsymbol{v}}^{ts} \beta_{v}).$$
(S65)

$$\beta_x + \delta_x \sim (\mathbf{A} + \mathbf{A}) + \mathbf{A} (\mathbf{A} \beta_x + \mathbf{C} + \mathbf{u} \delta_{\beta_u} + \delta_u \beta_u + \delta_{\beta_v} + \delta_v \beta_v).$$

After simplification and dropping the second order terms, we have

$$\boldsymbol{u}^{\mathsf{T}}\boldsymbol{\delta}_{\boldsymbol{v}}^{ts}\boldsymbol{\beta}_{v} + \boldsymbol{u}^{\mathsf{T}}\boldsymbol{X}\boldsymbol{\delta}_{\boldsymbol{\beta}_{x}}^{ts} + n\boldsymbol{\delta}_{\boldsymbol{\beta}_{u}}^{ts} + \boldsymbol{u}^{\mathsf{T}}\boldsymbol{v}\boldsymbol{\delta}_{\boldsymbol{\beta}_{v}}^{ts} \approx \boldsymbol{u}^{\mathsf{T}}\boldsymbol{\epsilon},$$
(S66)

$$\boldsymbol{v}^{\top} \delta_{\boldsymbol{u}}^{ts} \beta_{\boldsymbol{u}} + \boldsymbol{v}^{\top} \boldsymbol{X} \delta_{\boldsymbol{\beta}_{x}}^{ts} + \boldsymbol{u}^{\top} \boldsymbol{v} \delta_{\beta_{\boldsymbol{u}}}^{ts} + n \delta_{\beta_{\boldsymbol{v}}}^{ts} \approx \boldsymbol{v}^{\top} \boldsymbol{\epsilon},$$
(S67)

$$\delta_{\boldsymbol{\beta}_{x}}^{ts} \approx (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}(\boldsymbol{\epsilon} - \boldsymbol{u}\delta_{\boldsymbol{\beta}_{u}}^{ts} - \delta_{\boldsymbol{u}}^{ts}\boldsymbol{\beta}_{u} - \boldsymbol{v}\delta_{\boldsymbol{\beta}_{v}}^{ts} - \delta_{\boldsymbol{v}}^{ts}\boldsymbol{\beta}_{v}).$$
(S68)

Plugging (S68) into (S66, S67) and using (S42), we have

$$\boldsymbol{u}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})\delta_{\boldsymbol{u}}^{ts}\beta_{u}+\boldsymbol{u}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})\delta_{\boldsymbol{v}}^{ts}\beta_{v}+\boldsymbol{u}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})\boldsymbol{u}\delta_{\beta_{u}}^{ts}+\boldsymbol{u}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})\boldsymbol{v}\delta_{\beta_{v}}^{ts}$$
(S69)

$$\approx \boldsymbol{u}^{\top} (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{X}}) \boldsymbol{\epsilon}, \tag{S70}$$

$$\boldsymbol{v}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})\delta_{\boldsymbol{u}}^{ts}\beta_{u}+\boldsymbol{v}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})\delta_{\boldsymbol{v}}^{ts}\beta_{v}+\boldsymbol{v}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})\boldsymbol{u}\delta_{\beta_{u}}^{ts}+\boldsymbol{v}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})\boldsymbol{v}\delta_{\beta_{v}}^{ts}$$
(S71)

$$\approx \boldsymbol{v}^{\top} (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{X}}) \boldsymbol{\epsilon}. \tag{S72}$$

Denote

$$\tilde{\boldsymbol{u}} = (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{X}})\boldsymbol{u} \quad \text{and} \quad \tilde{\boldsymbol{v}} = (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{X}})\boldsymbol{v}$$
(S73)

and

$$\boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}} = \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top}\tilde{\boldsymbol{u}} \ \tilde{\boldsymbol{u}}^{\top}\tilde{\boldsymbol{v}} \\ \tilde{\boldsymbol{u}}^{\top}\tilde{\boldsymbol{v}} \ \tilde{\boldsymbol{v}}^{\top}\tilde{\boldsymbol{v}} \end{pmatrix}.$$
(S74)

Then (S70)-(S72) can be written as

$$\boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}\begin{pmatrix}\delta_{\beta_{u}}^{ts}\\\delta_{\beta_{v}}^{ts}\end{pmatrix}\approx\begin{pmatrix}\tilde{\boldsymbol{u}}^{\top}\\\tilde{\boldsymbol{v}}^{\top}\end{pmatrix}(-\beta_{u}\boldsymbol{I}_{n} - \beta_{v}\boldsymbol{I}_{n} \boldsymbol{I}_{n})\begin{pmatrix}\delta_{\boldsymbol{u}}^{ts}\\\delta_{\boldsymbol{v}}^{ts}\\\boldsymbol{\epsilon}\end{pmatrix}.$$
(S75)

Plugging (S46) into (S75) and solving for  $\begin{pmatrix} \delta_{\beta_u}^{ts} \\ \delta_{\beta_v}^{ts} \end{pmatrix}$ , we obtain

$$\begin{pmatrix} \delta_{\beta_{u}}^{ts} \\ \delta_{\beta_{v}}^{ts} \end{pmatrix} \approx \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} (-\beta_{u}\boldsymbol{I}_{n} - \beta_{v}\boldsymbol{I}_{n} \boldsymbol{I}_{n}) \begin{pmatrix} \boldsymbol{0}_{n\times n} & \boldsymbol{C}_{12}^{ts} \\ \boldsymbol{0}_{n\times n} & \boldsymbol{C}_{22}^{ts} \\ \boldsymbol{I}_{n} & \boldsymbol{0}_{n\times n^{2}} \end{pmatrix} \begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix}$$
(S76)

$$\stackrel{def}{=} \begin{pmatrix} C_{41}^{ts} C_{42}^{ts} \\ C_{51}^{ts} C_{52}^{ts} \end{pmatrix} \begin{pmatrix} \epsilon \\ \operatorname{vec}(E) \end{pmatrix}$$
(S77)

where explicitly

$$\begin{pmatrix} C_{41}^{ts} \\ C_{51}^{ts} \end{pmatrix} = C_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix}$$
(S78)

and

$$\begin{pmatrix} \boldsymbol{C}_{42}^{ts} \\ \boldsymbol{C}_{52}^{ts} \end{pmatrix} = (dn)^{-1} \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \beta_u \boldsymbol{v}^\top \otimes (\tilde{\boldsymbol{u}}^\top (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}})) \\ \beta_v \boldsymbol{u}^\top \otimes (\tilde{\boldsymbol{v}}^\top (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}})) \end{pmatrix} \boldsymbol{K} \end{pmatrix}.$$
(S79)

For  $\delta^{ts}_{\ensuremath{\boldsymbol{\beta}}_x}$ , plugging (S46) and (S77) into (S68),

$$\delta_{\boldsymbol{\beta}_{x}}^{ts} \approx (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}(\boldsymbol{\epsilon} - \boldsymbol{u}\delta_{\boldsymbol{\beta}_{u}}^{ts} - \delta_{\boldsymbol{u}}^{ts}\boldsymbol{\beta}_{u} - \boldsymbol{v}\delta_{\boldsymbol{\beta}_{v}}^{ts} - \delta_{\boldsymbol{v}}^{ts}\boldsymbol{\beta}_{v})$$
(S80)

$$= (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top} (-\beta_{u}\boldsymbol{I}_{n} - \beta_{v}\boldsymbol{I}_{n} - \boldsymbol{u} - \boldsymbol{v} \boldsymbol{I}_{n}) \begin{pmatrix} \boldsymbol{\delta}_{u}^{u} \\ \boldsymbol{\delta}_{v}^{ts} \\ \boldsymbol{\delta}_{\beta_{u}}^{ts} \\ \boldsymbol{\delta}_{\beta_{v}}^{ts} \\ \boldsymbol{\epsilon} \end{pmatrix}$$
(S81)

$$= (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top} (-\beta_{u}\boldsymbol{I}_{n} - \beta_{v}\boldsymbol{I}_{n} - \boldsymbol{u} - \boldsymbol{v} \boldsymbol{I}_{n}) \begin{pmatrix} \boldsymbol{0}_{n \times n} & \boldsymbol{C}_{12}^{ts} \\ \boldsymbol{0}_{n \times n} & \boldsymbol{C}_{22}^{ts} \\ \boldsymbol{C}_{41}^{ts} & \boldsymbol{C}_{52}^{ts} \\ \boldsymbol{C}_{51}^{ts} & \boldsymbol{C}_{52}^{ts} \\ \boldsymbol{I}_{n} & \boldsymbol{0}_{n \times n^{2}} \end{pmatrix} \begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix}$$
(S82)

$$\stackrel{def}{=} \left( \boldsymbol{C}_{61}^{ts} \, \boldsymbol{C}_{62}^{ts} \right) \begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix}.$$
(S83)

Finally, putting (S46), (S55), (S77) and (S83) together, we have

$$\begin{pmatrix} \hat{\boldsymbol{u}}^{ts} - \boldsymbol{u} \\ \hat{\boldsymbol{v}}^{ts} - \boldsymbol{v} \\ \operatorname{vec} \left( \hat{\boldsymbol{A}}^{ts} - \boldsymbol{A}_0 \right) \\ \hat{\boldsymbol{\beta}}_{u}^{ts} - \boldsymbol{\beta}_{u} \\ \hat{\boldsymbol{\beta}}_{v}^{ts} - \boldsymbol{\beta}_{v} \\ \hat{\boldsymbol{\beta}}_{x}^{ts} - \boldsymbol{\beta}_{x} \end{pmatrix} = \begin{pmatrix} \delta_{\boldsymbol{u}}^{ts} \\ \delta_{\boldsymbol{v}}^{ts} \\ \operatorname{vec} \left( \delta_{\boldsymbol{A}}^{ts} \right) \\ \operatorname{vec} \left( \delta_{\boldsymbol{A}}^{ts} \right) \\ \delta_{\boldsymbol{\beta}_{u}}^{ts} \\ \delta_{\boldsymbol{\beta}_{v}}^{ts} \\ \delta_{\boldsymbol{\beta}_{v}}^{ts} \\ \delta_{\boldsymbol{\beta}_{x}}^{ts} \end{pmatrix} \approx \begin{pmatrix} \boldsymbol{0}_{n \times n} & \boldsymbol{C}_{12}^{ts} \\ \boldsymbol{0}_{n \times n} & \boldsymbol{C}_{22}^{ts} \\ \boldsymbol{0}_{n^{2} \times n} & \boldsymbol{C}_{32}^{ts} \\ \boldsymbol{C}_{13}^{ts} & \boldsymbol{C}_{43}^{ts} \\ \boldsymbol{C}_{51}^{ts} & \boldsymbol{C}_{52}^{ts} \\ \boldsymbol{C}_{51}^{ts} & \boldsymbol{C}_{52}^{ts} \\ \boldsymbol{C}_{61}^{ts} & \boldsymbol{C}_{62}^{ts} \end{pmatrix} \begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix}^{def} = \boldsymbol{C}^{ts} \begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix}. \quad (S84)$$

Recall that we assume

$$\begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix} \sim N\left( \mathbf{0}_{(n+n^2)\times 1}, \begin{pmatrix} \sigma_y^2 \boldsymbol{I}_n \ \mathbf{0}_{n\times n^2} \\ \mathbf{0}_{n^2 \times n} \ \sigma_a^2 \boldsymbol{I}_{n^2} \end{pmatrix} \right).$$
(S85)

Therefore, the two-stage estimators converge to the following normal distribution asymptotically,

$$\begin{pmatrix} \hat{\boldsymbol{u}}^{ts} - \boldsymbol{u} \\ \hat{\boldsymbol{v}}^{ts} - \boldsymbol{v} \\ \operatorname{vec} \left( \hat{\boldsymbol{A}}^{ts} - \boldsymbol{A}_0 \right) \\ \hat{\boldsymbol{\beta}}_{u}^{ts} - \boldsymbol{\beta}_{u} \\ \hat{\boldsymbol{\beta}}_{x}^{ts} - \boldsymbol{\beta}_{v} \\ \hat{\boldsymbol{\beta}}_{x}^{ts} - \boldsymbol{\beta}_{x} \end{pmatrix} \xrightarrow{\mathcal{D}} N\left( \boldsymbol{0}_{(2n+n^{2}+2+p)\times 1}, \boldsymbol{C}^{ts} \left( \begin{array}{c} \sigma_{y}^{2} \boldsymbol{I}_{n} & \boldsymbol{0}_{n\times n^{2}} \\ \boldsymbol{0}_{n^{2}\times n} & \sigma_{a}^{2} \boldsymbol{I}_{n^{2}} \end{array} \right) \boldsymbol{C}^{ts\top} \right).$$
(S86)

Q.E.D.

## S6.1.2. Proof of Corollary 1

PROOF: From Theorem 1, we have

$$\begin{pmatrix} \delta_{\boldsymbol{u}}^{ts} \\ \delta_{\boldsymbol{v}}^{ts} \\ \delta_{\boldsymbol{A}}^{ts} \end{pmatrix} \approx \begin{pmatrix} \boldsymbol{C}_{12}^{ts} \\ \boldsymbol{C}_{22}^{ts} \\ \boldsymbol{C}_{32}^{ts} \end{pmatrix} \operatorname{vec}(\boldsymbol{E})$$
(S87)

where  $C_{12}^{ts} = (dn)^{-1} \boldsymbol{v}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}), C_{22}^{ts} = (dn)^{-1} (\boldsymbol{u}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}})) \boldsymbol{K}$  and  $C_{32}^{ts} = [(\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \otimes \boldsymbol{P}_{\boldsymbol{u}} + \boldsymbol{P}_{\boldsymbol{v}} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) + \boldsymbol{P}_{\boldsymbol{v}} \otimes \boldsymbol{P}_{\boldsymbol{u}}].$ 

(i) Rate of 
$$\hat{\boldsymbol{u}}^{ts}$$
,  $\hat{\boldsymbol{v}}^{ts}$ .

$$\frac{1}{n}\mathbb{E}\|\hat{\boldsymbol{u}}^{ts} - \boldsymbol{u}\|_{2}^{2} \approx \frac{1}{n}\mathrm{tr}(\sigma_{a}^{2}\boldsymbol{C}_{12}^{ts}\boldsymbol{C}_{12}^{ts\top})$$
(S88)

$$= \frac{1}{n} \frac{\sigma_a^2}{(dn)^2} \operatorname{tr}\left(\left(\boldsymbol{v}^\top \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}})\right) (\boldsymbol{v} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}))\right)$$
(S89)

$$=\frac{1}{n}\frac{\sigma_a^2}{(dn)^2}\mathrm{tr}\left(\boldsymbol{v}^{\top}\boldsymbol{v}\otimes(\boldsymbol{I}-\boldsymbol{P}\boldsymbol{u})\right)$$
(S90)

$$=\frac{1}{n}\frac{\sigma_a^2}{(dn)^2}\mathrm{tr}\left(n\boldsymbol{I}-\boldsymbol{u}\boldsymbol{u}^{\top}\right)\right)$$
(S91)

$$=\frac{\sigma_{a}^{2}}{d^{2}n^{2}}(n-1)$$
(S92)

$$=O(\kappa).$$
(S93)

Similarly for the rate of  $\hat{v}^{ts}$ ,

$$\frac{1}{n}\mathbb{E}\|\hat{\boldsymbol{v}}^{ts} - \boldsymbol{v}\|_{2}^{2} \approx \frac{1}{n}\mathrm{tr}(\sigma_{a}^{2}\boldsymbol{C}_{22}^{ts}\boldsymbol{C}_{22}^{ts\top})$$
(S94)

$$= \frac{1}{n} \frac{\sigma_a^2}{(dn)^2} \operatorname{tr}\left(\left(\boldsymbol{u}^\top \otimes (\boldsymbol{I} - \boldsymbol{P}\boldsymbol{v})\right) (\boldsymbol{u} \otimes (\boldsymbol{I} - \boldsymbol{P}\boldsymbol{v}))\right)$$
(S95)

$$=\frac{\sigma_a^2}{d^2 n^2}(n-1)$$
 (S96)

$$=O(\kappa) \tag{S97}$$

where the second equality is due to  $KK^{\top} = I$ .

(*ii*) Rate of  $\hat{A}^{ts}$ . Note that

$$[(\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \otimes \boldsymbol{P}_{\boldsymbol{u}}][(\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \otimes \boldsymbol{P}_{\boldsymbol{u}}]^{\top} = (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \otimes \boldsymbol{P}_{\boldsymbol{u}}$$
(S98)

$$[\mathbf{P}_{\boldsymbol{v}} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}})][\mathbf{P}_{\boldsymbol{v}} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}})]^{\top} = \mathbf{P}_{\boldsymbol{v}} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}})$$
(S99)

$$[\mathbf{P}_{\boldsymbol{v}} \otimes \mathbf{P}_{\boldsymbol{u}}][\mathbf{P}_{\boldsymbol{v}} \otimes \mathbf{P}_{\boldsymbol{u}}]^{\top} = \mathbf{P}_{\boldsymbol{v}} \otimes \mathbf{P}_{\boldsymbol{u}}$$
(S100)

$$[(\boldsymbol{I} - \boldsymbol{P}\boldsymbol{v}) \otimes \boldsymbol{P}\boldsymbol{u}][\boldsymbol{P}\boldsymbol{v} \otimes (\boldsymbol{I} - \boldsymbol{P}\boldsymbol{u})]^{\top} = \boldsymbol{0}$$
(S101)

$$[\boldsymbol{P}_{\boldsymbol{v}} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}})][(\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \otimes \boldsymbol{P}_{\boldsymbol{u}}]^{\top} = \boldsymbol{0}$$
(S102)

$$[(I - Pv) \otimes Pu][Pv \otimes Pu]^{\top} = 0$$
(S103)

$$[\mathbf{P}_{\boldsymbol{v}} \otimes (\mathbf{I} - \mathbf{P}_{\boldsymbol{u}})][\mathbf{P}_{\boldsymbol{v}} \otimes \mathbf{P}_{\boldsymbol{u}}]^{\top} = \mathbf{0}.$$
(S104)

Putting together, we have

$$\mathbb{E}\frac{\left\|\widehat{\boldsymbol{A}}^{ts} - \boldsymbol{A}_{0}\right\|_{F}^{2}}{\left\|\boldsymbol{A}_{0}\right\|_{F}^{2}} \approx \frac{1}{d^{2}n^{2}} \operatorname{tr}(\sigma_{a}^{2}\boldsymbol{C}_{32}^{ts}\boldsymbol{C}_{32}^{ts\top})$$
(S105)

$$=\frac{1}{d^2n^2}\sigma_a^2 \operatorname{tr}((\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{v}})\otimes\boldsymbol{P}_{\boldsymbol{u}}+\boldsymbol{P}_{\boldsymbol{v}}\otimes(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{u}})+\boldsymbol{P}_{\boldsymbol{v}}\otimes\boldsymbol{P}_{\boldsymbol{u}}) \quad (S106)$$

$$= \frac{\sigma_a^2}{d^2 n^2} \left[ 2(n-1) + 1 \right] = \frac{\sigma_a^2}{d^2 n^2} (2n-1) = O(\kappa).$$
(S107)

Q.E.D.

# S6.1.3. Proof of Corollary 2

PROOF: We now prove the rate of  $\hat{\beta}^{ts}$  of two-stage.

(1) Rate of  $\hat{\beta}_u^{ts}$  and  $\hat{\beta}_v^{ts}$ . Recall that

$$\begin{pmatrix} \delta_{\beta_{u}}^{ts} \\ \delta_{\beta_{v}}^{ts} \end{pmatrix} \approx \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} \begin{bmatrix} \boldsymbol{\epsilon} - (\delta_{\boldsymbol{u}}^{ts}\beta_{u} + \delta_{\boldsymbol{v}}^{ts}\beta_{v}) \end{bmatrix}$$
(S108)

$$= \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} \Big[ \boldsymbol{\epsilon} - \frac{1}{dn} \left( \beta_u \boldsymbol{v}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) + \beta_v \boldsymbol{u}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \boldsymbol{K} \right) \operatorname{vec}(\boldsymbol{E}) \Big] (S109)$$

$$\stackrel{def}{=} C_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} \left[ \epsilon + \boldsymbol{B} \operatorname{vec}(\boldsymbol{E}) \right]$$
(S110)

$$\stackrel{def}{=} \boldsymbol{D}_1 \boldsymbol{\epsilon} + \boldsymbol{D}_2 \operatorname{vec}(\boldsymbol{E}). \tag{S111}$$

Note that

$$\boldsymbol{D}_{1}\boldsymbol{D}_{1}^{\top} = \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} (\tilde{\boldsymbol{u}} \; \tilde{\boldsymbol{v}}) \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} = \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1}$$
(S112)

and

$$\boldsymbol{D}_{2}\boldsymbol{D}_{2}^{\top} = \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} \boldsymbol{B}\boldsymbol{B}^{\top} \left( \tilde{\boldsymbol{u}} \, \tilde{\boldsymbol{v}} \right) \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1}$$
(S113)

where

$$\boldsymbol{B}\boldsymbol{B}^{\top} = \frac{1}{(dn)^2} \Big[ \beta_u^2 (\boldsymbol{v}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}})) (\boldsymbol{v} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}})) + \beta_v^2 (\boldsymbol{u}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}})) (\boldsymbol{u} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}))$$
(S114)

$$+\beta_{u}\beta_{v}(\boldsymbol{v}^{\top}\otimes(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{u}}))\boldsymbol{K}^{\top}(\boldsymbol{u}\otimes(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{v}}))+\beta_{u}\beta_{v}(\boldsymbol{u}^{\top}\otimes(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{v}}))\boldsymbol{K}(\boldsymbol{v}\otimes(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{u}}))\Big|$$
S115)

$$= \frac{1}{d^2n} \left[ \beta_u^2 (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) + \beta_v^2 (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \right]$$
(S116)

since  $(\boldsymbol{v}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}))\boldsymbol{K}^{\top}(\boldsymbol{u} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}})) = (\boldsymbol{v}^{\top}(\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}})\boldsymbol{u})\boldsymbol{K} = 0$  and  $(\boldsymbol{u}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}))\boldsymbol{K}(\boldsymbol{v} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}})) = (\boldsymbol{u}^{\top}(\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}})\boldsymbol{v})\boldsymbol{K} = 0$ .

Plugging in (S112), (S113) and (S116),

-

$$Cov \begin{pmatrix} \delta_{\beta_u}^{t_s} \\ \delta_{\beta_v}^{t_s} \end{pmatrix} \approx \sigma_y^2 D_1 D_1^{\top} + \sigma_a^2 D_2 D_2^{\top}$$
(S117)

$$=\sigma_y^2 C_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \tag{S118}$$

$$+\sigma_{a}^{2}\frac{1}{d^{2}n}C_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1}\left(\frac{\tilde{\boldsymbol{u}}^{\top}}{\tilde{\boldsymbol{v}}^{\top}}\right)\left[\beta_{u}^{2}(\boldsymbol{I}-\boldsymbol{P}\boldsymbol{u})+\beta_{v}^{2}(\boldsymbol{I}-\boldsymbol{P}\boldsymbol{v})\right]\left(\tilde{\boldsymbol{u}}\;\tilde{\boldsymbol{v}}\right)C_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1}$$
(S119)

$$=\sigma_y^2 \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} + \sigma_a^2 \frac{1}{d^2 n} \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \beta_v^2 \tilde{\boldsymbol{u}}^\top (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \tilde{\boldsymbol{u}} & 0\\ 0 & \beta_u^2 \tilde{\boldsymbol{v}}^\top (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) \tilde{\boldsymbol{v}} \end{pmatrix} \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1}.$$
(S120)

Therefore, we obtain the rate of  $\hat{\beta}_{u}^{ts}$  and  $\hat{\beta}_{v}^{ts}$  from the diagonal entries of  $Cov\begin{pmatrix}\delta_{\beta_{u}}^{ts}\\\delta_{\beta_{v}}^{ts}\end{pmatrix}$  as

$$\mathbb{E}(\hat{\beta}_{u}^{ts} - \beta_{u})^{2} = \frac{\sigma_{y}^{2}}{c} \tilde{\boldsymbol{v}}^{\top} \tilde{\boldsymbol{v}}$$
(S121)

$$+ \frac{\sigma_a^2}{c^2} \frac{1}{d^2 n} \left[ \beta_v^2 \tilde{\boldsymbol{v}}^\top \tilde{\boldsymbol{v}} \tilde{\boldsymbol{u}}^\top (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \tilde{\boldsymbol{u}} \tilde{\boldsymbol{v}}^\top \tilde{\boldsymbol{v}} + \beta_u^2 \tilde{\boldsymbol{u}}^\top \tilde{\boldsymbol{v}} \tilde{\boldsymbol{v}}^\top (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) \tilde{\boldsymbol{v}} \tilde{\boldsymbol{u}}^\top \tilde{\boldsymbol{v}} \right] (1 + o(1)) \quad (S122)$$

$$\mathbb{E}(\hat{\beta}_v^{ts} - \beta_v)^2 = \frac{\sigma_y^2}{c} \tilde{\boldsymbol{u}}^\top \tilde{\boldsymbol{u}}$$
(S123)

$$+ \frac{\sigma_a^2}{c^2} \frac{1}{d^2 n} \left[ \beta_u^2 \tilde{\boldsymbol{u}}^\top \tilde{\boldsymbol{u}} \tilde{\boldsymbol{v}}^\top (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) \tilde{\boldsymbol{v}} \tilde{\boldsymbol{u}}^\top \tilde{\boldsymbol{u}} + \beta_v^2 \tilde{\boldsymbol{v}}^\top \tilde{\boldsymbol{u}} \tilde{\boldsymbol{u}}^\top (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \tilde{\boldsymbol{u}} \tilde{\boldsymbol{v}}^\top \tilde{\boldsymbol{u}} \right] (1 + o(1)) \quad (S124)$$

where  $c = \tilde{\boldsymbol{u}}^{\top} \tilde{\boldsymbol{u}} \tilde{\boldsymbol{v}}^{\top} \tilde{\boldsymbol{v}} - (\tilde{\boldsymbol{u}}^{\top} \tilde{\boldsymbol{v}})^2$ .

(2) Rate of  $\hat{\beta}_x^{ts}$ . Recall that

$$\delta_{\boldsymbol{\beta}_{x}}^{ts} \approx (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}(\boldsymbol{\epsilon} - \boldsymbol{u}\delta_{\boldsymbol{\beta}_{u}}^{ts} - \delta_{\boldsymbol{u}}^{ts}\boldsymbol{\beta}_{u} - \boldsymbol{v}\delta_{\boldsymbol{\beta}_{v}}^{ts} - \delta_{\boldsymbol{v}}^{ts}\boldsymbol{\beta}_{v}).$$
(S125)

From (S110), we have

$$\boldsymbol{\epsilon} - \left(\delta_{\boldsymbol{u}}^{ts}\beta_{u} + \delta_{\boldsymbol{v}}^{ts}\beta_{v}\right) \stackrel{def}{=} \boldsymbol{\epsilon} + \boldsymbol{B}\operatorname{vec}(\boldsymbol{E}) \tag{S126}$$

and

$$\begin{pmatrix} \delta_{\beta_u}^{ts} \\ \delta_{\beta_v}^{ts} \end{pmatrix} \approx \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^\top \\ \tilde{\boldsymbol{v}}^\top \end{pmatrix} [\boldsymbol{\epsilon} + \boldsymbol{B}\operatorname{vec}(\boldsymbol{E})].$$
 (S127)

Then we have

$$\boldsymbol{u}\delta^{ts}_{\beta_{u}} + \boldsymbol{v}\delta^{ts}_{\beta_{v}} = \left(\boldsymbol{u} \; \boldsymbol{v}\right) \begin{pmatrix} \delta^{ts}_{\beta_{u}} \\ \delta^{ts}_{\beta_{v}} \end{pmatrix} \approx \left(\boldsymbol{u} \; \boldsymbol{v}\right) \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} \left[\boldsymbol{\epsilon} + \boldsymbol{B}\operatorname{vec}(\boldsymbol{E})\right].$$
(S128)

Let

$$\tilde{\boldsymbol{G}} = \left(\boldsymbol{u} \; \boldsymbol{v}\right) \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} \text{ and } \boldsymbol{G} = \left(\boldsymbol{u} \; \boldsymbol{v}\right) \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \boldsymbol{u}^{\top} \\ \boldsymbol{v}^{\top} \end{pmatrix}.$$
(S129)

Plugging (S110) and (S128) into (S125) yields

$$\delta_{\boldsymbol{\beta}_{x}}^{ts} \approx (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top} \left[ \boldsymbol{I} - \left(\boldsymbol{u} \; \boldsymbol{v}\right) \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} \right] (\boldsymbol{\epsilon} + \boldsymbol{B}\operatorname{vec}(\boldsymbol{E}))$$
(S130)

$$= (X^{\top}X)^{-1}X^{\top}(I - \tilde{G})(\epsilon + B\operatorname{vec}(E))$$
(S131)

$$\stackrel{def}{=} \boldsymbol{F}_1 \boldsymbol{\epsilon} + \boldsymbol{F}_2 \operatorname{vec}(\boldsymbol{E}). \tag{S132}$$

Hence, the variance-covariance matrix of  $\delta_{\pmb{\beta}_x}^{ts}$  is

$$Cov\left(\delta_{\boldsymbol{\beta}_{x}}^{ts}\right) \approx \sigma_{y}^{2} \boldsymbol{F}_{1} \boldsymbol{F}_{1}^{\top} + \sigma_{a}^{2} \boldsymbol{F}_{2} \boldsymbol{F}_{2}^{\top}$$
(S133)

where we will derive the explicit form of  $F_1F_1^{\top}$  and  $F_2F_2^{\top}$  in the following. (a)  $F_1F_1^{\top}$ . Since

$$(I - \tilde{G})(I - \tilde{G})^{\top} = I - \tilde{G} - \tilde{G}^{\top} + G$$
(S134)

and

$$\begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} \boldsymbol{X} = \boldsymbol{0} \text{ and thus } \tilde{\boldsymbol{G}}\boldsymbol{X} = \boldsymbol{0}, \tag{S135}$$

consequently

$$F_1 F_1^{\top} = (X^{\top} X)^{-1} + (X^{\top} X)^{-1} X^{\top} G X (X^{\top} X)^{-1}.$$
 (S136)

(b)  $\boldsymbol{F}_{2}\boldsymbol{F}_{2}^{\top}$ . Recall (S116) where

$$\boldsymbol{B}\boldsymbol{B}^{\top} = \frac{1}{d^2n} \left[ \beta_u^2 (\boldsymbol{I} - \boldsymbol{P}\boldsymbol{u}) + \beta_v^2 (\boldsymbol{I} - \boldsymbol{P}\boldsymbol{v}) \right]$$

Plugging in we have,

$$(I - \tilde{G})BB^{\top}(I - \tilde{G})^{\top}$$
(S137)

$$= \frac{1}{d^2 n} \left[ \beta_u^2 (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) + \beta_v^2 (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) - \tilde{\boldsymbol{G}} \boldsymbol{B} \boldsymbol{B}^\top - \boldsymbol{B} \boldsymbol{B}^\top \tilde{\boldsymbol{G}}^\top \right]$$
(S138)

$$+ (\boldsymbol{u}\,\boldsymbol{v}) \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \beta_{v}^{2} \tilde{\boldsymbol{u}}^{\top} (\boldsymbol{I} - \boldsymbol{P}_{v}) \tilde{\boldsymbol{u}} & \boldsymbol{0} \\ \boldsymbol{0} & \beta_{u}^{2} \tilde{\boldsymbol{v}}^{\top} (\boldsymbol{I} - \boldsymbol{P}_{u}) \tilde{\boldsymbol{v}} \end{pmatrix} \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \boldsymbol{u}^{\top} \\ \boldsymbol{v}^{\top} \end{pmatrix} \end{bmatrix}.$$
(S139)

Together with (S136) and using (S135), we obtain the variance-covariance matrix of  $\delta_{\beta_x}^{ts}$  as follows.

$$Cov\left(\hat{\boldsymbol{\beta}}_{x}^{ts}-\boldsymbol{\beta}_{x}\right) \tag{S140}$$

$$\approx \sigma_y^2 \left[ (\boldsymbol{X}^\top \boldsymbol{X})^{-1} + (\boldsymbol{X}^\top \boldsymbol{X})^{-1} \boldsymbol{X}^\top \left( \boldsymbol{u} \, \boldsymbol{v} \right) \boldsymbol{C}_{\hat{\boldsymbol{u}}\hat{\boldsymbol{v}}}^{-1} \begin{pmatrix} \boldsymbol{u}^\top \\ \boldsymbol{v}^\top \end{pmatrix} \boldsymbol{X} (\boldsymbol{X}^\top \boldsymbol{X})^{-1} \right]$$
(S141)

$$+\sigma_a^2 \frac{1}{d^2 n} (\boldsymbol{X}^\top \boldsymbol{X})^{-1} \boldsymbol{X}^\top \left[ \beta_u^2 (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) + \beta_v^2 (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \right]$$
(S142)

$$+ \left(\boldsymbol{u} \, \boldsymbol{v}\right) \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \beta_{\boldsymbol{v}}^{2} \tilde{\boldsymbol{u}}^{\top} (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \tilde{\boldsymbol{u}} & \boldsymbol{0} \\ \boldsymbol{0} & \beta_{\boldsymbol{u}}^{2} \tilde{\boldsymbol{v}}^{\top} (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) \tilde{\boldsymbol{v}} \end{pmatrix} \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \boldsymbol{u}^{\top} \\ \boldsymbol{v}^{\top} \end{pmatrix} \end{bmatrix} \boldsymbol{X} (\boldsymbol{X}^{\top} \boldsymbol{X})^{-1}.$$
(S143)

#### S6.1.4. Proof of Corollary 3

PROOF: For simplicity, let us first consider  $\boldsymbol{y} = \boldsymbol{u}\beta_u + \boldsymbol{v}\beta_v + \boldsymbol{\epsilon}$  where  $\boldsymbol{\epsilon} \sim N(\boldsymbol{0}, \boldsymbol{I}_n)$  independently. Denote the correlation between  $\boldsymbol{u}$  and  $\boldsymbol{v}$  as  $\rho$ .

The two-stage procedure first estimates  $\hat{\boldsymbol{u}}$  and  $\boldsymbol{v}$  from the network  $\boldsymbol{A}$  as  $\hat{\boldsymbol{u}}^{ts}$  and  $\hat{\boldsymbol{v}}^{ts}$ , and then regresses  $\boldsymbol{y}$  on  $\hat{\boldsymbol{u}}^{ts}$  and  $\hat{\boldsymbol{v}}^{ts}$  to obtain  $\hat{\boldsymbol{\beta}}^{ts\top} = (\hat{\beta}_{u}^{ts}, \hat{\beta}_{v}^{ts})$  in the second stage. From Theorem 1, we have  $\hat{\boldsymbol{u}}^{ts} = \boldsymbol{u} + \delta_{\boldsymbol{u}}^{ts}$  and  $\hat{\boldsymbol{v}}^{ts} = \boldsymbol{v} + \delta_{\boldsymbol{v}}^{ts}$  where

$$\begin{pmatrix} \delta_{\boldsymbol{u}}^{ts} \\ \delta_{\boldsymbol{v}}^{ts} \end{pmatrix} \approx (dn)^{-1} \begin{pmatrix} \boldsymbol{v}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P}\boldsymbol{u}) \\ (\boldsymbol{u}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P}\boldsymbol{v})) \boldsymbol{K} \end{pmatrix} \operatorname{vec}(\boldsymbol{E}) \stackrel{def}{=} \begin{pmatrix} \boldsymbol{C}_{12}^{ts} \\ \boldsymbol{C}_{22}^{ts} \end{pmatrix} \operatorname{vec}(\boldsymbol{E}),$$
(S144)

 $u \perp \delta_{u}^{ts}$  and  $v \perp \delta_{v}^{ts}$ . Furthermore  $\delta_{u}^{ts} \perp \delta_{v}^{ts}$  because

$$\delta_{\boldsymbol{u}}^{ts^{\top}} \delta_{\boldsymbol{v}}^{ts} = \frac{1}{(dn)^2} \operatorname{vec}\left(\boldsymbol{E}\right)^{\top} \left(\boldsymbol{v} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}})\right) \left( \left(\boldsymbol{u}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}})\right) \right) \boldsymbol{K} \operatorname{vec}\left(\boldsymbol{E}\right)$$
(S145)

$$= \frac{\sigma_a^2}{(dn)^2} \operatorname{tr}\left(\left(\boldsymbol{v}\boldsymbol{u}^{\top}\right) \otimes \left((\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}})(\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}})\right)\boldsymbol{K}\right)$$
(S146)

$$= \frac{\sigma_a^2}{(dn)^2} \operatorname{tr}\left( (\boldsymbol{I} - \boldsymbol{P}\boldsymbol{u})(\boldsymbol{I} - \boldsymbol{P}\boldsymbol{v})\boldsymbol{v}\boldsymbol{u}^{\mathsf{T}} \right)$$
(S147)

$$= 0.$$
 (S148)

Denote  $\widetilde{W} = W + \delta$  where W = (u v) and  $\delta^{ts} = (\delta^{ts}_{u} \delta^{ts}_{v})$ , then

$$\widehat{\boldsymbol{\beta}}^{ts} = \left(\widetilde{\boldsymbol{W}}^{\top}\widetilde{\boldsymbol{W}}\right)^{-1}\widetilde{\boldsymbol{W}}^{\top}\boldsymbol{y} = \left(\boldsymbol{W}^{\top}\boldsymbol{W} + \boldsymbol{\delta}^{\top}\boldsymbol{\delta}\right)^{-1}\boldsymbol{W}^{\top}\boldsymbol{W}\begin{pmatrix}\beta_{u}\\\beta_{v}\end{pmatrix}.$$
(S149)

Since the correlation between  $\boldsymbol{u}$  and  $\boldsymbol{v}$  is  $\rho$ , the covariance matrix of  $\boldsymbol{W}$  is  $\Sigma_{\boldsymbol{u}\boldsymbol{v}} = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$ . The covariance matrix of  $\boldsymbol{\delta}^{ts}$  is  $\Sigma_{\boldsymbol{\delta}^{ts}} = \begin{pmatrix} 0 & \kappa \\ \kappa & 0 \end{pmatrix}$  due to Corollary 1 and  $\delta_{\boldsymbol{u}}^{ts} \perp \delta_{\boldsymbol{v}}^{ts}$ . Taking plim of both sides of (S149), we have

$$\operatorname{plim} \widehat{\boldsymbol{\beta}}^{ts} = \left( \Sigma_{\boldsymbol{u}\boldsymbol{v}} + \Sigma_{\boldsymbol{\delta}^{ts}} \right)^{-1} \Sigma_{\boldsymbol{u}\boldsymbol{v}} \begin{pmatrix} \beta_u \\ \beta_v \end{pmatrix}$$
(S150)

$$=\frac{1}{\left(1+\kappa\right)^{2}-\rho^{2}}\begin{pmatrix}1+\kappa-\rho^{2}&\kappa\rho\\\kappa\rho&1+\kappa-\rho^{2}\end{pmatrix}\begin{pmatrix}\beta_{u}\\\beta_{v}\end{pmatrix}.$$
(S151)

Specifically, we obtain

$$\operatorname{plim} \hat{\beta}_{u}^{ts} = \frac{(1+\kappa-\rho^{2})\beta_{u}+\kappa\rho\beta_{v}}{(1+\kappa)^{2}-\rho^{2}},$$
(S152)

plim 
$$\hat{\beta}_{v}^{ts} = \frac{(1+\kappa-\rho^{2})\beta_{v}+\kappa\rho\beta_{u}}{(1+\kappa)^{2}-\rho^{2}}.$$
 (S153)

When  $\boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta}_x + \boldsymbol{u}\boldsymbol{\beta}_u + \boldsymbol{v}\boldsymbol{\beta}_v + \boldsymbol{\epsilon}$  and  $\operatorname{cov}(\boldsymbol{X}, (\boldsymbol{u}\boldsymbol{v})) = \boldsymbol{0} \in \mathbb{R}^{p \times 2}$ , the above results remain the same due to the property of ordinary least square estimator.

Q.E.D.

Q.E.D.

# S6.2.1. Proof of Theorem 2

\_

PROOF: In Section S1, we derive the estimates of  $(\hat{d}, \hat{u}, \hat{v}, \hat{\beta})$  as (S6)-(S10). Together they lead to the first order expansion

$$(\boldsymbol{u} + \delta \boldsymbol{u})^{\top} (\boldsymbol{u} + \delta \boldsymbol{u}) \approx n \tag{S154}$$

$$(\boldsymbol{v} + \delta \boldsymbol{v})^{\top} (\boldsymbol{v} + \delta \boldsymbol{v}) \approx n$$
 (S155)

$$\beta_{u} + \delta_{\beta_{u}} \approx \frac{1}{n} (\boldsymbol{u} + \delta_{\boldsymbol{u}})^{\top} (\boldsymbol{X}\boldsymbol{\beta}_{x} + \boldsymbol{u}\beta_{u} + \boldsymbol{v}\beta_{v} + \boldsymbol{\epsilon} - \boldsymbol{X}\boldsymbol{\beta}_{x} - \boldsymbol{X}\delta_{\boldsymbol{\beta}_{x}} - \boldsymbol{v}\beta_{v} - \boldsymbol{v}\delta_{\beta_{v}} - \delta_{\boldsymbol{v}}\beta_{v} (S156)$$

$$\beta_{v} + \delta_{\beta_{v}} \approx \frac{1}{n} (\boldsymbol{v} + \delta_{\boldsymbol{v}})^{\top} (\boldsymbol{X} \boldsymbol{\beta}_{x} + \boldsymbol{u} \beta_{u} + \boldsymbol{v} \beta_{v} + \boldsymbol{\epsilon} - \boldsymbol{X} \boldsymbol{\beta}_{x} - \boldsymbol{X} \delta_{\boldsymbol{\beta}_{x}} - \boldsymbol{u} \beta_{u} - \boldsymbol{u} \beta_{u} - \delta_{\boldsymbol{u}} \beta_{u} (\boldsymbol{s})$$

$$\boldsymbol{\beta}_{x} + \delta_{\boldsymbol{\beta}_{x}} \approx (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}(\boldsymbol{X}\boldsymbol{\beta}_{x} + \boldsymbol{\epsilon} - \boldsymbol{u}\delta_{\boldsymbol{\beta}_{u}} - \delta_{\boldsymbol{u}}\boldsymbol{\beta}_{u} - \boldsymbol{v}\delta_{\boldsymbol{\beta}_{v}} - \delta_{\boldsymbol{v}}\boldsymbol{\beta}_{v})$$
(S158)

$$d + \delta_d \approx (\boldsymbol{u} + \delta_{\boldsymbol{u}})^\top (d\boldsymbol{u}\boldsymbol{v}^\top + \boldsymbol{E})(\boldsymbol{v} + \delta_{\boldsymbol{v}})/n^2$$
(S159)

$$(\beta_u + \delta_{\beta_u})(-\boldsymbol{\epsilon} + \boldsymbol{X}\delta_{\boldsymbol{\beta}_x} + \boldsymbol{u}\delta_{\beta_u} + \delta_{\boldsymbol{u}}\beta_u + \boldsymbol{v}\delta_{\beta_v} + \delta_{\boldsymbol{v}}\beta_v)$$
(S160)

$$+\lambda(d+\delta_d)^2(\boldsymbol{u}+\delta_{\boldsymbol{u}}) - \lambda(d+\delta_d)(d\boldsymbol{u}\boldsymbol{v}^\top + \boldsymbol{E})(\boldsymbol{v}+\delta_{\boldsymbol{v}})/n \approx 0$$
(S161)

$$(\beta_v + \delta_{\beta_v})(-\boldsymbol{\epsilon} + \boldsymbol{X}\delta_{\boldsymbol{\beta}_u} + \boldsymbol{u}\delta_{\beta_u} + \delta_{\boldsymbol{u}}\beta_u + \boldsymbol{v}\delta_{\beta_v} + \delta_{\boldsymbol{v}}\beta_v)$$
(S162)

$$+\lambda(d+\delta_d)^2(\boldsymbol{v}+\delta_{\boldsymbol{v}}) - \lambda(d+\delta_d)(d\boldsymbol{v}\boldsymbol{u}^{\top} + \boldsymbol{E}^{\top})(\boldsymbol{u}+\delta_{\boldsymbol{u}})/n \approx 0$$
(S163)

After simplification and dropping the second order terms,

$$\boldsymbol{u}^{\mathsf{T}}\delta\boldsymbol{u}/n\approx0,\tag{S164}$$

$$\boldsymbol{v}^{\mathsf{T}}\delta\boldsymbol{v}/n\approx0,\tag{S165}$$

$$\boldsymbol{u}^{\mathsf{T}}\delta\boldsymbol{v}\beta_{v} + \boldsymbol{u}^{\mathsf{T}}\boldsymbol{X}\delta_{\boldsymbol{\beta}_{x}} + n\delta_{\beta_{u}} + \boldsymbol{u}^{\mathsf{T}}\boldsymbol{v}\delta_{\beta_{v}} \approx \boldsymbol{u}^{\mathsf{T}}\boldsymbol{\epsilon},\tag{S166}$$

$$\boldsymbol{v}^{\mathsf{T}}\delta_{\boldsymbol{u}}\beta_{\boldsymbol{u}} + \boldsymbol{v}^{\mathsf{T}}\boldsymbol{X}\delta_{\boldsymbol{\beta}_{\boldsymbol{u}}} + \boldsymbol{u}^{\mathsf{T}}\boldsymbol{v}\delta_{\boldsymbol{\beta}_{\boldsymbol{u}}} + n\delta_{\boldsymbol{\beta}_{\boldsymbol{v}}} \approx \boldsymbol{v}^{\mathsf{T}}\boldsymbol{\epsilon},$$
(S167)

$$\delta \boldsymbol{\beta}_{x} \approx (\boldsymbol{X}^{\top} \boldsymbol{X})^{-1} \boldsymbol{X}^{\top} (\boldsymbol{\epsilon} - \boldsymbol{u} \delta_{\beta_{u}} - \delta_{\boldsymbol{u}} \beta_{u} - \boldsymbol{v} \delta_{\beta_{v}} - \delta_{\boldsymbol{v}} \beta_{v}), \qquad (S168)$$

$$\delta_d \approx d\boldsymbol{u}^{\mathsf{T}} \delta \boldsymbol{u} + d\boldsymbol{v}^{\mathsf{T}} \delta \boldsymbol{v} + \boldsymbol{u}^{\mathsf{T}} \boldsymbol{E} \boldsymbol{v} \approx \boldsymbol{u}^{\mathsf{T}} \boldsymbol{E} \boldsymbol{v} / n^2, \qquad (S169)$$

$$(\beta_u^2 + \lambda d^2)\delta \boldsymbol{u} + \beta_u \beta_v \delta \boldsymbol{v} + \beta_u \boldsymbol{X} \delta_{\boldsymbol{\beta}_x} + \beta_u \boldsymbol{u} \delta_{\beta_u} + \beta_u \boldsymbol{v} \delta_{\beta_v} \approx \lambda d(\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) \boldsymbol{E} \boldsymbol{v} / n + \beta_u \boldsymbol{\epsilon}, (S170)$$

$$\beta_{u}\beta_{v}\delta\boldsymbol{u} + (\beta_{v}^{2} + \lambda d^{2})\delta\boldsymbol{v} + \beta_{v}\boldsymbol{X}\delta\boldsymbol{\beta}_{x} + \beta_{v}\boldsymbol{u}\delta_{\beta_{u}} + \beta_{v}\boldsymbol{v}\delta_{\beta_{v}} \approx \lambda d(\boldsymbol{I} - \boldsymbol{P}\boldsymbol{v})\boldsymbol{E}^{\top}\boldsymbol{u}/n + \beta_{v}\boldsymbol{\epsilon}(\boldsymbol{S}171)$$

where (S166)-(S167) have used (S164)-(S165).

Plugging (S168) into (S166)-(S167) and using (S164)-(S165)

$$\boldsymbol{u}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})\delta\boldsymbol{u}\beta_{\boldsymbol{u}}+\boldsymbol{u}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})\delta\boldsymbol{v}\beta_{\boldsymbol{v}}+\boldsymbol{u}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})\boldsymbol{u}\delta_{\beta_{\boldsymbol{u}}}+\boldsymbol{u}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})\boldsymbol{v}\delta_{\beta_{\boldsymbol{v}}}$$
(S172)  
$$\approx \boldsymbol{u}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})\boldsymbol{\epsilon},$$
(S173)

$$\boldsymbol{v}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})\delta\boldsymbol{u}\beta_{\boldsymbol{u}}+\boldsymbol{v}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})\delta\boldsymbol{v}\beta_{\boldsymbol{v}}+\boldsymbol{v}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})\boldsymbol{u}\delta_{\beta_{\boldsymbol{u}}}+\boldsymbol{v}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})\boldsymbol{v}\delta_{\beta_{\boldsymbol{v}}} \quad (S174)$$
  
$$\approx \boldsymbol{v}^{\top}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})\boldsymbol{\epsilon}. \quad (S175)$$

Recall  $\tilde{\boldsymbol{u}} = (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{X}})\boldsymbol{u}, \, \tilde{\boldsymbol{v}} = (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{X}})\boldsymbol{v}$  and  $\boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}} = \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top}\tilde{\boldsymbol{u}} \, \tilde{\boldsymbol{u}}^{\top}\tilde{\boldsymbol{v}} \\ \tilde{\boldsymbol{u}}^{\top}\tilde{\boldsymbol{v}} \, \tilde{\boldsymbol{v}}^{\top}\tilde{\boldsymbol{v}} \end{pmatrix}$ . Then (S172)-(S175) can be written as

$$\begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} (-\beta_u \boldsymbol{I}_n - \beta_v \boldsymbol{I}_n) \begin{pmatrix} \delta \boldsymbol{u} \\ \delta \boldsymbol{v} \end{pmatrix} + \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \delta_{\beta_u} \\ \delta_{\beta_v} \end{pmatrix} \approx \begin{pmatrix} \delta \boldsymbol{u} \\ \delta \boldsymbol{v} \end{pmatrix} \boldsymbol{\epsilon}.$$
 (S176)

Solving for  $\delta_{\beta_u}$ ,  $\delta_{\beta_v}$  gives

$$\begin{pmatrix} \delta_{\beta_{u}} \\ \delta_{\beta_{v}} \end{pmatrix} \approx \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} (-\beta_{u}\boldsymbol{I}_{n} - \beta_{v}\boldsymbol{I}_{n} \boldsymbol{I}_{n}) \begin{pmatrix} \delta_{\boldsymbol{u}} \\ \delta_{\boldsymbol{v}} \\ \boldsymbol{\epsilon} \end{pmatrix}.$$
(S177)

On the other hand, plugging (S168) into (S170, S171)

$$\begin{pmatrix} \beta_{u}^{2}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})+\lambda d^{2}\boldsymbol{I} & \beta_{u}\beta_{v}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})\\ \beta_{u}\beta_{v}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}}) & \beta_{v}^{2}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{X}})+\lambda d^{2}\boldsymbol{I} \end{pmatrix} \begin{pmatrix} \delta\boldsymbol{u}\\ \delta\boldsymbol{v} \end{pmatrix} + \begin{pmatrix} \beta_{u}\\ \beta_{v} \end{pmatrix} (\tilde{\boldsymbol{u}}\,\tilde{\boldsymbol{v}}) \begin{pmatrix} \delta\beta_{u}\\ \delta\beta_{v} \end{pmatrix}$$
(S178)

$$\approx \begin{pmatrix} \beta_u (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{X}}) & \lambda d\boldsymbol{v}^\top \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}})/n \\ \beta_v (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{X}}) & \lambda d \left( \boldsymbol{u}^\top \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}})/n \right) \boldsymbol{K} \end{pmatrix} \begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix}.$$
(S179)

Combining (S177)-(S179) and after some algebra, we have

$$\begin{pmatrix} \beta_{u}^{2} \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) + \lambda d^{2}\boldsymbol{I} & \beta_{u}\beta_{v} \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) \\ \beta_{u}\beta_{v} \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) & \beta_{v}^{2} \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) + \lambda d^{2}\boldsymbol{I} \end{pmatrix} \begin{pmatrix} \delta \boldsymbol{u} \\ \delta \boldsymbol{v} \end{pmatrix}$$
(S180)

$$\approx \left[ \lambda d^{2} \boldsymbol{I}_{2n} + \begin{pmatrix} \beta_{u}^{2} \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) & \beta_{u} \beta_{v} \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) \\ \beta_{u} \beta_{v} \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) & \beta_{v}^{2} \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) \end{pmatrix} \right] \begin{pmatrix} \delta \boldsymbol{u} \\ \delta \boldsymbol{v} \end{pmatrix}$$
(S181)

$$\approx \left[\lambda d^{2}\boldsymbol{I}_{2n} + \begin{pmatrix} \beta_{u}^{2} & \beta_{u}\beta_{v} \\ \beta_{u}\beta_{v} & \beta_{v}^{2} \end{pmatrix} \otimes \left(\boldsymbol{I} - \boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)\right] \begin{pmatrix} \delta\boldsymbol{u} \\ \delta\boldsymbol{v} \end{pmatrix}$$
(S182)

$$\approx \begin{pmatrix} \beta_{u} \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) & \lambda d\boldsymbol{v}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P}\boldsymbol{u})/n \\ \beta_{v} \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) & \lambda d \left( \boldsymbol{u}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P}\boldsymbol{v})/n \right) \boldsymbol{K} \end{pmatrix} \begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix}$$
(S183)

because  $(\tilde{u} \ \tilde{v}) C_{\tilde{u}\tilde{v}}^{-1} \begin{pmatrix} \tilde{u}^{\top} \\ \tilde{v}^{\top} \end{pmatrix} = P_{(\tilde{u}\tilde{v})}$  and  $I - P_X - P_{(\tilde{u}\tilde{v})} = I - P_{(Xuv)}$ . Using the Woodbury identity and the property of projection matrix,

$$\left(\lambda d^{2}\boldsymbol{I}_{2n} + \begin{pmatrix} \beta_{u}^{2} & \beta_{u}\beta_{v} \\ \beta_{u}\beta_{v} & \beta_{v}^{2} \end{pmatrix} \otimes \left(\boldsymbol{I} - \boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)\right)^{-1}$$
(S184)

$$= (\lambda d^2)^{-1} \left( \boldsymbol{I}_{2n} - (\lambda d^2 + \beta_u^2 + \beta_v^2)^{-1} \begin{pmatrix} \beta_u^2 & \beta_u \beta_v \\ \beta_u \beta_v & \beta_v^2 \end{pmatrix} \otimes \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) \right).$$
(S185)

Hence,

$$\begin{pmatrix} \delta \boldsymbol{u} \\ \delta \boldsymbol{v} \end{pmatrix} \approx (\lambda d^2)^{-1} \left( \boldsymbol{I}_{2n} - (\lambda d^2 + \beta_u^2 + \beta_v^2)^{-1} \begin{pmatrix} \beta_u^2 & \beta_u \beta_v \\ \beta_u \beta_v & \beta_v^2 \end{pmatrix} \otimes \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) \right)$$
(S186)

$$\begin{pmatrix} \beta_{u} \left( \boldsymbol{I} - \boldsymbol{P}_{\left( \boldsymbol{X} \boldsymbol{u} \boldsymbol{v} \right)} \right) & \lambda d\boldsymbol{v}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P} \boldsymbol{u})/n \\ \beta_{v} \left( \boldsymbol{I} - \boldsymbol{P}_{\left( \boldsymbol{X} \boldsymbol{u} \boldsymbol{v} \right)} \right) & \lambda d \left( \boldsymbol{u}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P} \boldsymbol{v})/n \right) \boldsymbol{K} \end{pmatrix} \begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix}$$
(S187)

$$\stackrel{def}{=} \begin{pmatrix} \boldsymbol{C}_{11} \ \boldsymbol{C}_{12} \\ \boldsymbol{C}_{21} \ \boldsymbol{C}_{22} \end{pmatrix} \begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix}.$$
(S188)

Let  $\hat{A} = \hat{d}\hat{u}\hat{v}^{\top}$ . Plugging in (S169) and (S188), the first order expansion leads to

$$\operatorname{vec}\left(\delta_{\boldsymbol{A}}\right) = \operatorname{vec}(\hat{d}\hat{\boldsymbol{u}}\hat{\boldsymbol{v}}^{\top}) - \operatorname{vec}(d\boldsymbol{u}\boldsymbol{v}^{\top})$$
(S189)

$$= \operatorname{vec}(d + \delta_d)(\boldsymbol{u} + \delta_{\boldsymbol{u}})(\boldsymbol{v} + \delta_{\boldsymbol{v}})^{\top} - \operatorname{vec}(d\boldsymbol{u}\boldsymbol{v}^{\top})$$
(S190)

$$= d \operatorname{vec}(\boldsymbol{u} \delta \boldsymbol{v}^{\mathsf{T}}) + d \operatorname{vec}(\delta \boldsymbol{u} \boldsymbol{v}^{\mathsf{T}}) + \boldsymbol{u}^{\mathsf{T}} \boldsymbol{E} \boldsymbol{v} \operatorname{vec}(\boldsymbol{u} \boldsymbol{v}^{\mathsf{T}}) / n^{2}$$
(S191)

$$= d\boldsymbol{K}\operatorname{vec}(\delta_{\boldsymbol{v}}\boldsymbol{u}^{\top}) + d\operatorname{vec}(\delta_{\boldsymbol{u}}\boldsymbol{v}^{\top}) + \boldsymbol{P}_{\boldsymbol{v}}\otimes\boldsymbol{P}_{\boldsymbol{u}}$$
(S192)

$$= d\mathbf{K}(\mathbf{u} \otimes \mathbf{I}_n) \delta_{\mathbf{v}} + d(\mathbf{v} \otimes \mathbf{I}_n) \delta_{\mathbf{u}} + \mathbf{P}_{\mathbf{v}} \otimes \mathbf{P}_{\mathbf{u}}$$
(S193)

$$= \left( d(\boldsymbol{v} \otimes \boldsymbol{I}_n) \ d\boldsymbol{K}(\boldsymbol{u} \otimes \boldsymbol{I}_n) \ \boldsymbol{P}\boldsymbol{v} \otimes \boldsymbol{P}\boldsymbol{u} \right) \begin{pmatrix} \boldsymbol{C}_{11} & \boldsymbol{C}_{12} \\ \boldsymbol{C}_{21} & \boldsymbol{C}_{22} \\ \boldsymbol{0}_{n^2 \times n} & \boldsymbol{I}_{n^2} \end{pmatrix} \begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix}$$
(S194)

$$\stackrel{def}{=} \left( \boldsymbol{C}_{31} \ \boldsymbol{C}_{32} \right) \begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix}.$$
(S195)

For  $\delta_{\beta_u}$  and  $\delta_{\beta_v}$ , plugging (S188) into (S177)

$$\begin{pmatrix} \delta_{\beta_{u}} \\ \delta_{\beta_{v}} \end{pmatrix} \approx \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} (-\beta_{u}\boldsymbol{I}_{n} - \beta_{v}\boldsymbol{I}_{n} \boldsymbol{I}_{n}) \begin{pmatrix} \delta \boldsymbol{u} \\ \delta \boldsymbol{v} \\ \boldsymbol{\epsilon} \end{pmatrix}$$
(S196)

$$= \boldsymbol{C}_{\boldsymbol{\tilde{u}}\boldsymbol{\tilde{v}}}^{-1} \begin{pmatrix} \boldsymbol{\tilde{u}}^{\top} \\ \boldsymbol{\tilde{v}}^{\top} \end{pmatrix} (-\beta_u \boldsymbol{I}_n - \beta_v \boldsymbol{I}_n \boldsymbol{I}_n) \begin{pmatrix} \boldsymbol{C}_{11} & \boldsymbol{C}_{12} \\ \boldsymbol{C}_{21} & \boldsymbol{C}_{22} \\ \boldsymbol{I}_n & \boldsymbol{0}_{n \times n^2} \end{pmatrix} \begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix}$$
(S197)

$$\stackrel{def}{=} \begin{pmatrix} \boldsymbol{C}_{41} \ \boldsymbol{C}_{42} \\ \boldsymbol{C}_{51} \ \boldsymbol{C}_{52} \end{pmatrix} \begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix}.$$
(S198)

Lastly for  $\delta_{\ensuremath{\boldsymbol{\beta}}_x}$ , we plug (S188) and (S198) into (S168)

$$\delta \boldsymbol{\beta}_{x} \approx (\boldsymbol{X}^{\top} \boldsymbol{X})^{-1} \boldsymbol{X}^{\top} (\boldsymbol{\epsilon} - \boldsymbol{u} \delta_{\beta_{u}} - \delta_{\boldsymbol{u}} \beta_{u} - \boldsymbol{v} \delta_{\beta_{v}} - \delta_{\boldsymbol{v}} \beta_{v})$$
(S199)

$$= (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top} (-\beta_{u}\boldsymbol{I}_{n} - \beta_{v}\boldsymbol{I}_{n} - \boldsymbol{u} - \boldsymbol{v} \boldsymbol{I}_{n}) \begin{pmatrix} \delta \boldsymbol{u} \\ \delta \boldsymbol{v} \\ \delta_{\beta_{u}} \\ \delta_{\beta_{v}} \\ \boldsymbol{\epsilon} \end{pmatrix}$$
(S200)

$$\approx (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}(-\beta_{u}\boldsymbol{I}_{n} - \beta_{v}\boldsymbol{I}_{n} - \boldsymbol{u} - \boldsymbol{v} \boldsymbol{I}_{n}) \begin{pmatrix} \boldsymbol{C}_{11} & \boldsymbol{C}_{12} \\ \boldsymbol{C}_{21} & \boldsymbol{C}_{22} \\ \boldsymbol{C}_{31} & \boldsymbol{C}_{32} \\ \boldsymbol{C}_{41} & \boldsymbol{C}_{42} \\ \boldsymbol{I}_{n} & \boldsymbol{0}_{n \times n^{2}} \end{pmatrix} \begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix}$$
(S201)

$$\stackrel{def}{=} \left( \boldsymbol{C}_{51} \, \boldsymbol{C}_{52} \right) \begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix}. \tag{S202}$$

Finally, putting (S188), (S195), (S198) and (S202) together

$$\begin{pmatrix} \hat{\boldsymbol{u}} - \boldsymbol{u} \\ \hat{\boldsymbol{v}} - \boldsymbol{v} \\ \operatorname{vec} \left( \hat{\boldsymbol{A}} - \boldsymbol{A}_0 \right) \\ \hat{\beta}_u - \beta_u \\ \hat{\beta}_v - \beta_v \\ \hat{\beta}_x - \beta_x \end{pmatrix} = \begin{pmatrix} \delta \boldsymbol{u} \\ \delta \boldsymbol{v} \\ \operatorname{vec} \left( \delta_{\boldsymbol{A}} \right) \\ \delta_{\beta_u} \\ \delta_{\beta_v} \\ \delta_{\beta_v} \\ \delta_{\boldsymbol{\beta}_x} \end{pmatrix} \approx \begin{pmatrix} C_{11} C_{12} \\ C_{21} C_{22} \\ C_{31} C_{32} \\ C_{41} C_{42} \\ C_{51} C_{52} \\ C_{61} C_{62} \end{pmatrix} \begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix} = \boldsymbol{C} \begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix}. \quad (S203)$$

Recall that we assume

$$\begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix} \sim N \begin{pmatrix} \boldsymbol{0}_{(n+n^2)\times 1}, \begin{pmatrix} \sigma_y^2 \boldsymbol{I}_n & \boldsymbol{0}_{n\times n^2} \\ \boldsymbol{0}_{n^2 \times n} & \sigma_a^2 \boldsymbol{I}_{n^2} \end{pmatrix} \end{pmatrix}.$$
(S204)

Therefore, the SuperCENT estimators converge to the following normal distribution asymptotically,

$$\begin{pmatrix} \hat{\boldsymbol{u}} - \boldsymbol{u} \\ \hat{\boldsymbol{v}} - \boldsymbol{v} \\ \operatorname{vec} \left( \hat{\boldsymbol{A}} - \boldsymbol{A}_0 \right) \\ \hat{\beta}_u - \beta_u \\ \hat{\beta}_v - \beta_v \\ \hat{\beta}_x - \beta_x \end{pmatrix} \xrightarrow{\mathcal{D}} N \Big( \boldsymbol{0}_{(2n+n^2+2+p)\times 1}, \boldsymbol{C} \left( \begin{array}{c} \sigma_y^2 \boldsymbol{I}_n & \boldsymbol{0}_{n\times n^2} \\ \boldsymbol{0}_{n^2 \times n} & \sigma_a^2 \boldsymbol{I}_{n^2} \end{array} \right) \boldsymbol{C}^\top \Big).$$
(S205)  
$$Q.E.D.$$

# S6.2.2. Proof of Corollary 4

PROOF: We first show the rate of  $\hat{u}$ ,  $\hat{v}$  then  $\hat{A}$ .

(1) Rate of  $\hat{u}$ ,  $\hat{v}$ . From Theorem 2, we have

$$\begin{pmatrix} \delta \boldsymbol{u} \\ \delta \boldsymbol{v} \end{pmatrix} \approx \begin{pmatrix} \boldsymbol{C}_{11} \ \boldsymbol{C}_{12} \\ \boldsymbol{C}_{21} \ \boldsymbol{C}_{22} \end{pmatrix} \begin{pmatrix} \boldsymbol{\epsilon} \\ \operatorname{vec}(\boldsymbol{E}) \end{pmatrix}$$
(S206)

where

$$\begin{pmatrix} \boldsymbol{C}_{11} \\ \boldsymbol{C}_{21} \end{pmatrix} = (\lambda d^2 + \beta_u^2 + \beta_v^2)^{-1} \begin{pmatrix} \beta_u \\ \beta_v \end{pmatrix} \otimes \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right)$$
(S207)

and

(

$$\begin{pmatrix} C_{12} \\ C_{22} \end{pmatrix} = \frac{1}{dn} \left[ \begin{pmatrix} \boldsymbol{v}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) \\ (\boldsymbol{u}^{\top} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}})) \boldsymbol{K} \end{pmatrix} \right]$$

$$- \frac{1}{(\lambda d^{2} + \beta_{u}^{2} + \beta_{v}^{2})} \begin{pmatrix} \beta_{u}^{2} \boldsymbol{v}^{\top} \otimes \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) + \beta_{u} \beta_{v} \left( \boldsymbol{u}^{\top} \otimes \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) \right) \boldsymbol{K} \\ \beta_{u} \beta_{v} \boldsymbol{v}^{\top} \otimes \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) + \beta_{v}^{2} \left( \boldsymbol{u}^{\top} \otimes \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) \right) \boldsymbol{K} \\ \end{bmatrix} \right]$$

$$S208)$$

$$S209)$$

For the rate of  $\hat{u}$ ,

$$\mathbb{E}\|\hat{\boldsymbol{u}} - \boldsymbol{u}\|^2 \approx \operatorname{tr}(\sigma_y^2 \boldsymbol{C}_{11} \boldsymbol{C}_{11}^\top + \sigma_a^2 \boldsymbol{C}_{12} \boldsymbol{C}_{12}^\top)$$
(S210)

$$=\frac{\sigma_y^2}{(\lambda d^2+\beta_u^2+\beta_v^2)^2}\operatorname{tr}\left(\beta_u^2\Big(\boldsymbol{I}-\boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\Big)\right)$$
(S211)

$$+\frac{\sigma_a^2}{n^2 d^2} \operatorname{tr}\left(\boldsymbol{v}^{\top} \boldsymbol{v}\right) \operatorname{tr}\left(\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}\right)$$
(S212)

$$+\frac{\sigma_a^2}{n^2 d^2 (\lambda d^2 + \beta_u^2 + \beta_v^2)^2} \operatorname{tr} \left(\beta_u^4 \boldsymbol{v}^\top \boldsymbol{v} \left(\boldsymbol{I} - \boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right) + \right)$$
(S213)

$$2\beta_{u}^{3}\beta_{v}\boldsymbol{u}\boldsymbol{v}^{\top}\otimes\left(\boldsymbol{I}-\boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)\boldsymbol{K}^{\top}+\beta_{u}^{2}\beta_{v}^{2}\boldsymbol{u}\boldsymbol{u}^{\top}\otimes\left(\boldsymbol{I}-\boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)\right)$$
(S214)

$$-\frac{2\sigma_a^2}{nd^2(\lambda d^2 + \beta_u^2 + \beta_v^2)} \operatorname{tr}\left(\beta_u^2/n\left(\boldsymbol{v}^\top \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}})\right)\left(\boldsymbol{v} \otimes \left(\boldsymbol{I} - \boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)\right)$$
(\$215)

$$\beta_u \beta_v / n\left( \boldsymbol{v}^\top \otimes (\boldsymbol{I} - \boldsymbol{P} \boldsymbol{u}) \right) \left( \boldsymbol{K}^\top \boldsymbol{u} \otimes \left( \boldsymbol{I} - \boldsymbol{P}_{\left( \boldsymbol{X} \boldsymbol{u} \boldsymbol{v} \right)} \right) \right) \right)$$
(S216)

$$=\frac{\sigma_{y}^{2}\beta_{u}^{2}(n-p-2)}{(\lambda d^{2}+\beta_{u}^{2}+\beta_{v}^{2})^{2}}+\left(\frac{\sigma_{a}^{2}}{d^{2}}-\frac{\sigma_{a}^{2}}{d^{2}n}\right)$$
(S217)

$$+\frac{\sigma_a^2\beta_u^2(n-p-2)}{n^2d^2(\lambda d^2+\beta_u^2+\beta_v^2)^2}\left(\beta_u^2n+2\beta_u\beta_v\operatorname{tr}(\boldsymbol{v}^{\top}\boldsymbol{u})+\beta_v^2n\right)$$
(S218)

$$-\frac{2\sigma_a^2(n-p-2)}{nd^2(\lambda d^2 + \beta_u^2 + \beta_v^2)}\beta_u^2$$
(S219)

$$= \left(\frac{\sigma_a^2}{d^2} - \frac{\sigma_a^2}{d^2n}\right) \tag{S220}$$

$$-\frac{\beta_{u}^{2}(n-p-2)}{(\lambda d^{2}+\beta_{u}^{2}+\beta_{v}^{2})^{2}}\left[\frac{2\lambda d^{2}+\beta_{u}^{2}+\beta_{v}^{2}}{d^{2}n}\sigma_{a}^{2}-\sigma_{y}^{2}\right].$$
(S221)

Then,

$$\frac{1}{n}\mathbb{E}\|\hat{\boldsymbol{u}} - \boldsymbol{u}\|^2 \approx \frac{\sigma_a^2(n-1)}{d^2n^2} - \frac{\beta_u^2(n-p-2)}{n(\lambda d^2 + \beta_u^2 + \beta_v^2)^2} \left[\frac{2\lambda d^2 + \beta_u^2 + \beta_v^2}{d^2n}\sigma_a^2 - \sigma_y^2\right]$$
(S222)

$$= O\left(\frac{\sigma_a^2}{d^2n} - \beta_u^2 \delta_{ts,sc}\right) \tag{S223}$$

where  $\delta_{ts,sc} = \frac{1}{(\lambda d^2 + \beta_u^2 + \beta_v^2)^2} \left[ \frac{2\lambda d^2 + \beta_u^2 + \beta_v^2}{d^2 n} \sigma_a^2 - \sigma_y^2 \right].$ 

To get the optimal  $\lambda$  in Remark 11, we take the partial derivative of  $\ell_u \triangleq (S221)$  with respect to  $\lambda$  yields

$$\frac{\partial \ell_u}{\partial \lambda} = \frac{\beta_u^2 (n-p-2)}{(\lambda d^2 + \beta_u^2 + \beta_v^2)^3} \left[ 2d^2 \sigma_y^2 + \frac{\sigma_a^2}{n} \left( 2(\lambda d^2 + \beta_u^2 + \beta_v^2) - 4\lambda d^2 - 2\beta_u^2 - 2\beta_v^2 \right) \right]$$
(S224)

$$=\frac{\beta_{u}^{2}(n-p-2)}{(\lambda d^{2}+\beta_{u}^{2}+\beta_{v}^{2})^{3}}\left[2d^{2}\sigma_{y}^{2}-\frac{2d^{2}}{n}\sigma_{a}^{2}\lambda\right].$$
(S225)

Setting  $\frac{\partial \ell_u}{\partial \lambda} = 0$  yields

$$\lambda_0 = \frac{n\sigma_y^2}{\sigma_a^2}.$$
(S226)

When  $\lambda \in (0, \lambda_0]$ ,  $\ell_u$  increases as  $\lambda$  increases;  $\lambda \in (\lambda_0, \infty)$ ,  $\ell_u$  decreases and converges to 0 as  $\lambda$  increases. The maximum of  $\ell_u$  is then taken at  $\lambda_0$ .

Similarly, we derive the rate of  $\hat{v}$  as

$$\frac{1}{n}\mathbb{E}\|\hat{\boldsymbol{v}} - \boldsymbol{v}\|^2 \approx \frac{\sigma_a^2(n-1)}{d^2n^2} - \frac{\beta_v^2(n-p-2)}{n(\lambda d^2 + \beta_u^2 + \beta_v^2)^2} \left[\frac{2\lambda d^2 + \beta_u^2 + \beta_v^2}{d^2n}\sigma_a^2 - \sigma_y^2\right]$$
(S227)

$$= O\left(\frac{\sigma_a^2}{d^2n} - \beta_v^2 \delta_{ts,sc}\right).$$
(S228)

(2) Rate of  $\hat{A}$ . From Theorem 2, we have

$$\operatorname{vec}\left(\delta_{\boldsymbol{A}}\right) \approx d\boldsymbol{K}(\boldsymbol{u} \otimes \boldsymbol{I}_{n})\delta_{\boldsymbol{v}} + d(\boldsymbol{v} \otimes \boldsymbol{I}_{n})\delta_{\boldsymbol{u}} + \boldsymbol{P}\boldsymbol{v} \otimes \boldsymbol{P}\boldsymbol{u}.$$
(S229)

Using (S188), we have

$$\operatorname{vec}\left(\delta_{\boldsymbol{A}}\right) \tag{S230}$$

$$\approx \frac{d}{\lambda d^{2} + \beta_{u}^{2} + \beta_{v}^{2}} \left( \beta_{v} \mathbf{K} \boldsymbol{u} \otimes \boldsymbol{P}_{(\boldsymbol{X} \boldsymbol{u} \boldsymbol{v})} + \beta_{u} \boldsymbol{v} \right) \otimes \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X} \boldsymbol{u} \boldsymbol{v})} \right) \boldsymbol{\epsilon}$$
(S231)

$$+ \left\{ \left[ (I - P_{v}) \otimes P_{u} + P_{v} \otimes (I - P_{u}) \right] \right\}$$
(S232)

$$-\frac{1}{\lambda d^{2}+\beta_{u}^{2}+\beta_{v}^{2}}\left[\beta_{v}^{2}\left(\boldsymbol{I}-\boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)\otimes\boldsymbol{P}\boldsymbol{u}+\beta_{u}^{2}\boldsymbol{P}\boldsymbol{v}\otimes\left(\boldsymbol{I}-\boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)\right]$$
(S233)

$$-\frac{\beta_{u}\beta_{v}}{\lambda d^{2}+\beta_{u}^{2}+\beta_{v}^{2}}\left[\boldsymbol{K}\boldsymbol{u}\boldsymbol{v}^{\top}/n\otimes\left(\boldsymbol{I}-\boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)+\boldsymbol{v}\boldsymbol{u}^{\top}/n\otimes\left(\boldsymbol{I}-\boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)\boldsymbol{K}\right] \quad (S234)$$

$$+P_{\boldsymbol{v}}\otimes P_{\boldsymbol{u}}$$
  $\left. \operatorname{vec}(\boldsymbol{E}) \right.$  (S235)

$$\stackrel{def}{=} \boldsymbol{H}_1 \boldsymbol{\epsilon} + (\boldsymbol{H}_2 + \boldsymbol{H}_3 + \boldsymbol{H}_4 + \boldsymbol{H}_5) \operatorname{vec}(\boldsymbol{E})$$
(S236)

because

$$d\boldsymbol{K}(\boldsymbol{u}\otimes\boldsymbol{I})\delta\boldsymbol{v} \tag{S237}$$

$$= d\mathbf{K}(\mathbf{u} \otimes \mathbf{I})[\mathbf{C}_{21}\boldsymbol{\epsilon} + \mathbf{C}_{22}\operatorname{vec}(\mathbf{E})]$$
(S238)

$$= d\boldsymbol{K}(\boldsymbol{u} \otimes \boldsymbol{I}) \frac{\beta_{v}}{\lambda d^{2} + \beta_{u}^{2} + \beta_{v}^{2}} \left( \boldsymbol{I} - \boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)} \right) \boldsymbol{\epsilon}$$
(S239)

$$+d\boldsymbol{K}(\boldsymbol{u}\otimes\boldsymbol{I})\frac{1}{dn}\boldsymbol{u}^{\top}\otimes(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{v}})\boldsymbol{K}\operatorname{vec}(\boldsymbol{E})$$
(S240)

$$-d\boldsymbol{K}(\boldsymbol{u}\otimes\boldsymbol{I})\frac{1}{dn}\frac{1}{(\lambda d^{2}+\beta_{u}^{2}+\beta_{v}^{2})}\left[\beta_{u}\beta_{v}\boldsymbol{v}^{\top}\otimes\left(\boldsymbol{I}-\boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)+\right]$$
(S241)

$$\beta_v^2 \left( \boldsymbol{u}^\top \otimes \left( \boldsymbol{I} - \boldsymbol{P}_{\left( \boldsymbol{X} \boldsymbol{u} \boldsymbol{v} \right)} \right) \right) \boldsymbol{K} \right] \operatorname{vec}(\boldsymbol{E}) \quad (S242)$$

$$=\frac{d\beta_{v}}{\lambda d^{2}+\beta_{u}^{2}+\beta_{v}^{2}}\boldsymbol{K}\left(\boldsymbol{u}\otimes\left(\boldsymbol{I}-\boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)\right)\boldsymbol{\epsilon}$$
(S243)

+ 
$$\left[ (\boldsymbol{I} - \boldsymbol{P}\boldsymbol{v}) \otimes \boldsymbol{P}\boldsymbol{u} - \frac{\beta_v^2}{\lambda d^2 + \beta_u^2 + \beta_v^2} \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) \otimes \boldsymbol{P}\boldsymbol{u} \right] \operatorname{vec}(\boldsymbol{E})$$
(S244)

$$-\frac{\beta_{u}\beta_{v}}{\lambda d^{2}+\beta_{u}^{2}+\beta_{v}^{2}}\boldsymbol{K}\left(\boldsymbol{u}\boldsymbol{v}^{\top}/n\otimes\left(\boldsymbol{I}-\boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)\right)\operatorname{vec}(\boldsymbol{E})$$
(S245)

and similarly

$$d(\boldsymbol{v}\otimes\boldsymbol{I})\delta\boldsymbol{u} \tag{S246}$$

$$=\frac{d\beta_{u}}{\lambda d^{2}+\beta_{u}^{2}+\beta_{v}^{2}}\boldsymbol{v}\otimes\left(\boldsymbol{I}-\boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)\boldsymbol{\epsilon}$$
(S247)

+ 
$$\left[ \boldsymbol{P}\boldsymbol{v} \otimes (\boldsymbol{I} - \boldsymbol{P}\boldsymbol{u}) - \frac{\beta_u^2}{\lambda d^2 + \beta_u^2 + \beta_v^2} \boldsymbol{P}\boldsymbol{v} \otimes \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) \right] \operatorname{vec}(\boldsymbol{E})$$
 (S248)

$$-\frac{\beta_{u}\beta_{v}}{\lambda d^{2}+\beta_{u}^{2}+\beta_{v}^{2}}\left(\boldsymbol{v}\boldsymbol{u}^{\top}/n\otimes\left(\boldsymbol{I}-\boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)\right)\boldsymbol{K}\operatorname{vec}(\boldsymbol{E}).$$
(S249)

In order to obtain the rate of  $\hat{A}$ , we need  $H_1H_1^{\top}$  and  $(H_2 + H_3 + H_4 + H_5)(H_2 + H_3 + H_4 + H_5)^{\top}$ :

$$\boldsymbol{H}_{1}\boldsymbol{H}_{1}^{\top} = \frac{d^{2}}{(\lambda d^{2} + \beta_{u}^{2} + \beta_{v}^{2})^{2}} \left\{ \beta_{v}^{2}\boldsymbol{K} \left( \boldsymbol{u}\boldsymbol{u}^{\top} \otimes \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) \right) \boldsymbol{K}^{\top} + \beta_{u}^{2}\boldsymbol{v}\boldsymbol{v}^{\top} \otimes \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) \right)$$
(S250)

$$+\beta_{u}\beta_{v}\boldsymbol{K}\left(\boldsymbol{u}\boldsymbol{v}^{\top}\otimes\left(\boldsymbol{I}-\boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})}\right)\right)+\beta_{u}\beta_{v}\left(\boldsymbol{v}\boldsymbol{u}^{\top}\otimes\left(\boldsymbol{I}-\boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})}\right)\right)\boldsymbol{K}^{\top}\right\}$$
(S251)

$$\boldsymbol{H}_{2}\boldsymbol{H}_{2}^{\top} = (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \otimes \boldsymbol{P}_{\boldsymbol{u}} + \boldsymbol{P}_{\boldsymbol{v}} \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}})$$
(S252)

$$\boldsymbol{H}_{3}\boldsymbol{H}_{3}^{\top} = \frac{1}{\left(\lambda d^{2} + \beta_{u}^{2} + \beta_{v}^{2}\right)^{2}} \left[\beta_{u}{}^{4}\boldsymbol{P}_{\boldsymbol{v}} \otimes \left(\boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})}\right) + \beta_{v}{}^{4}\left(\boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})}\right) \otimes \boldsymbol{P}_{\boldsymbol{u}}\right]$$
(S253)

$$\boldsymbol{H}_{4}\boldsymbol{H}_{4}^{\top} = \frac{\beta_{u}^{2}\beta_{v}^{2}}{(\lambda d^{2} + \beta_{u}^{2} + \beta_{v}^{2})^{2}} \left[ \boldsymbol{K}\boldsymbol{P}_{u} \otimes \left(\boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})}\right)\boldsymbol{K}^{\top} + \boldsymbol{P}_{v} \otimes \left(\boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})}\right) \right]$$
(S254)

$$\boldsymbol{H}_{5}\boldsymbol{H}_{5}^{\top} = \boldsymbol{P}_{\boldsymbol{v}} \otimes \boldsymbol{P}_{\boldsymbol{u}}$$
(S255)

$$\boldsymbol{H}_{2}\boldsymbol{H}_{3}^{\top} = \boldsymbol{H}_{3}\boldsymbol{H}_{2}^{\top} \tag{S256}$$

$$= -\frac{1}{\lambda d^{2} + \beta_{u}^{2} + \beta_{v}^{2}} \left[ \beta_{u}^{2} \boldsymbol{P}_{\boldsymbol{v}} \otimes \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) + \beta_{v}^{2} \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) \otimes \boldsymbol{P}_{\boldsymbol{u}} \right]$$
(S257)

$$\boldsymbol{H}_{2}\boldsymbol{H}_{4}^{\top} = \boldsymbol{H}_{4}\boldsymbol{H}_{2}^{\top} \tag{S258}$$

$$= -\frac{\beta_{u}\beta_{v}}{\lambda d^{2} + \beta_{u}^{2} + \beta_{v}^{2}} \left[ K\left(\frac{uv^{\top}}{n} \otimes \left(I - P_{(Xuv)}\right)\right) + \left(\frac{vu^{\top}}{n} \otimes \left(I - P_{(Xuv)}\right)\right) K^{\top} \right]$$
(S259)

$$\boldsymbol{H}_{3}\boldsymbol{H}_{4}^{\mathsf{T}} = \frac{\beta_{u}\beta_{v}}{(\lambda d^{2} + \beta_{u}^{2} + \beta_{v}^{2})^{2}} \left[ \beta_{u}^{2} \left( \frac{\boldsymbol{v}\boldsymbol{u}^{\mathsf{T}}}{n} \otimes \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) \right) \boldsymbol{K}^{\mathsf{T}} + \beta_{v}^{2} \boldsymbol{K} \left( \frac{\boldsymbol{u}\boldsymbol{v}^{\mathsf{T}}}{n} \otimes \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) \right) \right]$$
(S260)

$$\boldsymbol{H}_{4}\boldsymbol{H}_{3}^{\top} = \frac{\beta_{u}\beta_{v}}{(\lambda d^{2} + \beta_{u}^{2} + \beta_{v}^{2})^{2}} \left[ \beta_{u}^{2}\boldsymbol{K} \left( \frac{\boldsymbol{u}\boldsymbol{v}^{\top}}{n} \otimes \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) \right) + \beta_{v}^{2}\boldsymbol{K} \left( \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) \otimes \frac{\boldsymbol{v}\boldsymbol{u}^{\top}}{n} \right) \right]$$
(S261)

$$H_2 H_5^{\top} = H_3 H_5^{\top} = H_4 H_5^{\top} = H_5 H_2^{\top} = H_5 H_3^{\top} = H_5 H_4^{\top} = 0.$$
(S262)

Plugging in (S251) to (S262) yields the rate of  $\hat{A}$  as follows.

$$\mathbb{E}\left\|\hat{\boldsymbol{A}} - \boldsymbol{A}_0\right\|_F^2 \tag{S263}$$

$$\approx \sigma_y^2 \operatorname{tr} \left( \boldsymbol{H}_1 \boldsymbol{H}_1^{\mathsf{T}} \right) + \sigma_a^2 \operatorname{tr} \left( (\boldsymbol{H}_2 + \boldsymbol{H}_3 + \boldsymbol{H}_4 + \boldsymbol{H}_5) (\boldsymbol{H}_2 + \boldsymbol{H}_3 + \boldsymbol{H}_4 + \boldsymbol{H}_5)^{\mathsf{T}} \right)$$
(S264)

$$= \sigma_y^2 \frac{d^2}{(\lambda d^2 + \beta_u^2 + \beta_v^2)^2} \left(\beta_u^2 + \beta_v^2\right) n(n-p-2)$$
(S265)

$$+\sigma_a^2 [2(n-1)+1]$$
(S266)

$$+\sigma_{a}^{2}\frac{n-p-2}{(\lambda d^{2}+\beta_{u}^{2}+\beta_{v}^{2})^{2}}\left[\beta_{u}^{4}+\beta_{v}^{4}+2\beta_{u}^{2}\beta_{v}^{2}-2(\beta_{u}^{2}+\beta_{v}^{2})(\lambda d^{2}+\beta_{u}^{2}+\beta_{v}^{2})\right]$$
(S267)

$$=\sigma_{a}^{2}(2n-1)$$
(S268)  
$$n-n-2 \qquad (S268)$$

$$-\frac{n-p-2}{(\lambda d^2+\beta_u^2+\beta_v^2)^2}(\beta_u^2+\beta_v^2)\left[(2\lambda d^2+\beta_u^2+\beta_v^2)\sigma_a^2-nd^2\sigma_y^2\right].$$
(S269)

Therefore,

$$\mathbb{E}\frac{\left\|\hat{\boldsymbol{A}} - \boldsymbol{A}_{0}\right\|_{F}^{2}}{\left\|\boldsymbol{A}_{0}\right\|_{F}^{2}} \approx \frac{\sigma_{a}^{2}(2n-1)}{d^{2}n} - \frac{n-p-2}{n(\lambda d^{2}+\beta_{u}^{2}+\beta_{v}^{2})^{2}}(\beta_{u}^{2}+\beta_{v}^{2})\left[\frac{(2\lambda d^{2}+\beta_{u}^{2}+\beta_{v}^{2})}{d^{2}n}\sigma_{a}^{2} - \sigma_{y}^{2}\right]$$
(S270)

$$= O\left(\frac{\sigma_a^2}{d^2n} - (\beta_u^2 + \beta_v^2)\delta_{ts,sc}\right).$$
(S271)

# S6.2.3. Proof of Corollary 5

We now prove the rate of  $\hat{\beta}$  of SuperCENT.

(1) Rate of  $\hat{\beta}_u$  and  $\hat{\beta}_v$ . From Theorem 2, we have

$$\boldsymbol{\epsilon} - (\delta \boldsymbol{u} \beta_u + \delta \boldsymbol{v} \beta_v) \tag{S272}$$

$$= (\boldsymbol{I} - \beta_u \boldsymbol{C}_{11} + \beta_v \boldsymbol{C}_{12})\boldsymbol{\epsilon} + (\beta_u \boldsymbol{C}_{12} + \beta_v \boldsymbol{C}_{22})\operatorname{vec}(\boldsymbol{E})$$
(S273)

$$= \left[ \boldsymbol{I} - \frac{\beta_u^2 + \beta_v^2}{\lambda d^2 + \beta_u^2 + \beta_v^2} \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) \right] \boldsymbol{\epsilon}$$
(S274)

$$-\frac{1}{dn} \left[ \beta_u \boldsymbol{v}^\top \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) + \beta_v \left( \boldsymbol{u}^\top \otimes (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \right) \boldsymbol{K} \right] \operatorname{vec}(\boldsymbol{E})$$
(S275)

$$+\frac{\beta_u^2+\beta_v^2}{dn(\lambda d^2+\beta_u^2+\beta_v^2)} \left[\beta_u \boldsymbol{v}^\top \otimes \left(\boldsymbol{I}-\boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)+\right]$$
(S276)

$$\beta_{v}\left(\boldsymbol{u}^{\top}\otimes\left(\boldsymbol{I}-\boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)\boldsymbol{K}\right)\right]\operatorname{vec}(\boldsymbol{E})$$
 (S277)

$$\stackrel{def}{=} \boldsymbol{A}_1 \boldsymbol{\epsilon} + (\boldsymbol{C}_1 + \boldsymbol{C}_2) \operatorname{vec}(\boldsymbol{E}).$$
(S278)

Plugging (S278) into (S177), we obtain

$$\begin{pmatrix} \delta_{\beta_{u}} \\ \delta_{\beta_{v}} \end{pmatrix} \approx \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} \Big[ \boldsymbol{\epsilon} - (\delta_{\boldsymbol{u}}\beta_{u} + \delta_{\boldsymbol{v}}\beta_{v}) \Big]$$
(S279)

$$= C_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} [\boldsymbol{A}_{1}\boldsymbol{\epsilon} + (\boldsymbol{C}_{1} + \boldsymbol{C}_{2})\operatorname{vec}(\boldsymbol{E})]$$
(S280)

$$\stackrel{def}{=} \boldsymbol{B}_1 \boldsymbol{\epsilon} + \boldsymbol{B}_2 \operatorname{vec}(\boldsymbol{E}). \tag{S281}$$

To get the rate of  $\hat{\beta}_u$  and  $\hat{\beta}_v$ , we next calculate  $\boldsymbol{B}_1 \boldsymbol{B}_1^{\top}$  and  $\boldsymbol{B}_2 \boldsymbol{B}_2^{\top}$ . (a)  $\boldsymbol{B}_1 \boldsymbol{B}_1^{\top}$ . Since

$$\boldsymbol{A}_{1}\boldsymbol{A}_{1}^{\top} = \boldsymbol{I} - \frac{(\beta_{u}^{2} + \beta_{v}^{2})(2\lambda d^{2} + \beta_{u}^{2} + \beta_{v}^{2})}{(\lambda d^{2} + \beta_{u}^{2} + \beta_{v}^{2})^{2}} \left(\boldsymbol{I} - \boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)$$
(S282)
34

and

$$\left(I - P_{\left(Xuv\right)}\right)\tilde{u} = 0,$$
 (S283)

consequently

$$\boldsymbol{B}_{1}\boldsymbol{B}_{1}^{\top} = \boldsymbol{C}_{\boldsymbol{\tilde{\boldsymbol{u}}}\boldsymbol{\tilde{\boldsymbol{v}}}}^{-1} \begin{pmatrix} \boldsymbol{\tilde{\boldsymbol{u}}}^{\top} \\ \boldsymbol{\tilde{\boldsymbol{v}}}^{\top} \end{pmatrix} \boldsymbol{A}_{1}\boldsymbol{A}_{1}^{\top} \begin{pmatrix} \boldsymbol{\tilde{\boldsymbol{u}}} \ \boldsymbol{\tilde{\boldsymbol{v}}} \end{pmatrix} \boldsymbol{C}_{\boldsymbol{\tilde{\boldsymbol{u}}}\boldsymbol{\tilde{\boldsymbol{v}}}}^{-1} = \boldsymbol{C}_{\boldsymbol{\tilde{\boldsymbol{u}}}\boldsymbol{\tilde{\boldsymbol{v}}}}^{-1}.$$
(S284)

(b)  $\boldsymbol{B}_2 \boldsymbol{B}_2^\top$ . Since

$$C_1 C_1^{\top} \tag{S285}$$

$$=\frac{1}{(dn)^2}\Big[\beta_u^2(\boldsymbol{v}^\top\otimes(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{u}}))(\boldsymbol{v}\otimes(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{u}}))+\beta_v^2(\boldsymbol{u}^\top\otimes(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{v}}))(\boldsymbol{u}\otimes(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{v}}))$$
(S286)

$$+\beta_{u}\beta_{v}(\boldsymbol{v}^{\top}\otimes(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{u}}))\boldsymbol{K}^{\top}(\boldsymbol{u}\otimes(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{v}}))+\beta_{u}\beta_{v}(\boldsymbol{u}^{\top}\otimes(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{v}}))\boldsymbol{K}(\boldsymbol{v}\otimes(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{u}}))\Big] (S287)$$

$$=\frac{1}{d^2n}\left[\beta_u^2(\boldsymbol{I}-\boldsymbol{P}\boldsymbol{u})+\beta_v^2(\boldsymbol{I}-\boldsymbol{P}\boldsymbol{v})\right],\tag{S288}$$

$$\boldsymbol{C}_{2}\boldsymbol{C}_{2}^{\mathsf{T}} = n\left(\frac{\beta_{u}^{2} + \beta_{v}^{2}}{dn(\lambda d^{2} + \beta_{u}^{2} + \beta_{v}^{2})}\right)^{2} \left[\beta_{u}^{2}\left(\boldsymbol{I} - \boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right) + \beta_{v}^{2}\left(\boldsymbol{I} - \boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)\right] \quad (S289)$$

and

$$\boldsymbol{C}_{1}\boldsymbol{C}_{2}^{\top} = \boldsymbol{C}_{2}\boldsymbol{C}_{1}^{\top} = -\frac{1}{d^{2}n}\frac{\beta_{u}^{2} + \beta_{v}^{2}}{\lambda d^{2} + \beta_{u}^{2} + \beta_{v}^{2}} \left[\beta_{u}^{2} \left(\boldsymbol{I} - \boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right) + \beta_{v}^{2} \left(\boldsymbol{I} - \boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)\right], (S290)$$

consequently

$$(\boldsymbol{C}_{1} + \boldsymbol{C}_{2})(\boldsymbol{C}_{1} + \boldsymbol{C}_{2})^{\top}$$
(S291)  
=  $\frac{1}{d^{2}n} \left[ \left( \beta_{u}^{2} (\boldsymbol{I} - \boldsymbol{P}\boldsymbol{u}) + \beta_{v}^{2} (\boldsymbol{I} - \boldsymbol{P}\boldsymbol{v}) \right) - \frac{(\beta_{u}^{2} + \beta_{v}^{2})^{2} (2\lambda d^{2} + \beta_{u}^{2} + \beta_{v}^{2})}{(\lambda d^{2} + \beta_{u}^{2} + \beta_{v}^{2})^{2}} \left( \boldsymbol{I} - \boldsymbol{P}_{(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v})} \right) \right]$ (S292)

Therefore,

$$\boldsymbol{B}_{2}\boldsymbol{B}_{2}^{\top} = \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} (\boldsymbol{C}_{1} + \boldsymbol{C}_{2}) (\boldsymbol{C}_{1} + \boldsymbol{C}_{2})^{\top} (\tilde{\boldsymbol{u}} \; \tilde{\boldsymbol{v}}) \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1}$$
(S293)

$$= \frac{1}{d^2 n} \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} \left( \beta_u^2 (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) + \beta_v^2 (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \right) \left( \tilde{\boldsymbol{u}} \; \tilde{\boldsymbol{v}} \right) \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1}$$
(S294)

$$=\frac{1}{d^2n}\boldsymbol{C}_{\boldsymbol{\tilde{u}\tilde{v}}}^{-1}\begin{pmatrix}\beta_v^2\boldsymbol{\tilde{u}}^\top(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{v}})\boldsymbol{\tilde{u}}&0\\0&\beta_u^2\boldsymbol{\tilde{v}}^\top(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{u}})\boldsymbol{\tilde{v}}\end{pmatrix}\boldsymbol{C}_{\boldsymbol{\tilde{u}\tilde{v}}}^{-1}$$
(S295)

where the first equality is due to (\$283). Combining (\$284) and (\$295),

$$Cov \begin{pmatrix} \delta_{\beta_u} \\ \delta_{\beta_v} \end{pmatrix} \approx \sigma_y^2 \boldsymbol{B}_1 \boldsymbol{B}_1^\top + \sigma_a^2 \boldsymbol{B}_2 \boldsymbol{B}_2^\top$$
(S296)

$$=\sigma_y^2 C_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \tag{S297}$$

$$+\sigma_{a}^{2}\frac{1}{d^{2}n}\boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1}\begin{pmatrix}\tilde{\boldsymbol{u}}^{\top}\\\tilde{\boldsymbol{v}}^{\top}\end{pmatrix}\Big[\beta_{u}^{2}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{u}})+\beta_{v}^{2}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{v}})\Big]\begin{pmatrix}\tilde{\boldsymbol{u}}\\\tilde{\boldsymbol{v}}\end{pmatrix}\boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1}$$
(S298)

$$=\sigma_y^2 \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} + \sigma_a^2 \frac{1}{d^2 n} \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \beta_v^2 \tilde{\boldsymbol{u}}^\top (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \tilde{\boldsymbol{u}} & 0\\ 0 & \beta_u^2 \tilde{\boldsymbol{v}}^\top (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) \tilde{\boldsymbol{v}} \end{pmatrix} \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1}$$
(S299)

where is the same as  $Cov \begin{pmatrix} \delta_{\beta_u}^{ts} \\ \delta_{\beta_v}^{ts} \end{pmatrix}$  in (S120). Therefore,

$$\mathbb{E}(\hat{\beta}_u - \beta_u)^2 = \mathbb{E}(\hat{\beta}_u^{ts} - \beta_u)^2 \quad \text{and} \quad \mathbb{E}(\hat{\beta}_v - \beta_v)^2 = \mathbb{E}(\hat{\beta}_v^{ts} - \beta_v)^2.$$
(S300)

(2) Rate of  $\hat{\boldsymbol{\beta}}_x$ . Recall that

$$\tilde{\boldsymbol{G}} = \left(\boldsymbol{u} \, \boldsymbol{v}\right) \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \tilde{\boldsymbol{u}}^{\top} \\ \tilde{\boldsymbol{v}}^{\top} \end{pmatrix} \quad \text{and} \quad \boldsymbol{G} = \left(\boldsymbol{u} \, \boldsymbol{v}\right) \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \boldsymbol{u}^{\top} \\ \boldsymbol{v}^{\top} \end{pmatrix}.$$
(S301)

By plugging (S188) and (S198) into (S168), we have

$$\delta_{\boldsymbol{\beta}_{x}} \approx (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}(\boldsymbol{\epsilon} - \boldsymbol{u}\delta_{\beta_{u}} - \delta_{\boldsymbol{u}}\beta_{u} - \boldsymbol{v}\delta_{\beta_{v}} - \delta_{\boldsymbol{v}}\beta_{v})$$
(S302)

$$\approx (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top} \left[ \boldsymbol{I} - (\boldsymbol{u} \, \boldsymbol{v}) \, \boldsymbol{C}_{\boldsymbol{u}\boldsymbol{v}}^{-1} \begin{pmatrix} \boldsymbol{\tilde{u}}^{\top} \\ \boldsymbol{\tilde{v}}^{\top} \end{pmatrix} \right] [\boldsymbol{A}_{1}\boldsymbol{\epsilon} + (\boldsymbol{C}_{1} + \boldsymbol{C}_{2}) \operatorname{vec}(\boldsymbol{E})]$$
(S303)

$$= (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}(\boldsymbol{I}-\tilde{\boldsymbol{G}})[\boldsymbol{A}_{1}\boldsymbol{\epsilon}+(\boldsymbol{C}_{1}+\boldsymbol{C}_{2})\operatorname{vec}(\boldsymbol{E})]$$
(S304)

$$\stackrel{def}{=} \boldsymbol{F}_1 \boldsymbol{\epsilon} + \boldsymbol{F}_2 \operatorname{vec}(\boldsymbol{E}). \tag{S305}$$

To get the rate of  $\hat{\boldsymbol{\beta}}_x$ , we next calculate  $\boldsymbol{F}_1 \boldsymbol{F}_1^{\top}$  and  $\boldsymbol{F}_2 \boldsymbol{F}_2^{\top}$ . (a)  $\boldsymbol{F}_1 \boldsymbol{F}_1^{\top}$ . Since

$$(\boldsymbol{I} - \tilde{\boldsymbol{G}})(\boldsymbol{I} - \tilde{\boldsymbol{G}})^{\top} = \boldsymbol{I} - \tilde{\boldsymbol{G}} - \tilde{\boldsymbol{G}}^{\top} + \boldsymbol{G}$$
(S306)

and

$$\left(I - P_{\left(Xuv\right)}\right)\tilde{u} = 0 \text{ and } \left(I - P_{\left(Xuv\right)}\right)u = 0,$$
 (S307)

consequently

$$(I - \tilde{G})A_{1}A_{1}^{\top}(I - \tilde{G})^{\top} = (I - \tilde{G})(I - \tilde{G})^{\top} - \frac{(\beta_{u}^{2} + \beta_{v}^{2})(2\lambda d^{2} + \beta_{u}^{2} + \beta_{v}^{2})}{(\lambda d^{2} + \beta_{u}^{2} + \beta_{v}^{2})^{2}} \left(I - P_{(Xuv)}\right).$$
(S308)

Further because 
$$\left( \boldsymbol{I} - \boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)} \right) \boldsymbol{X} = \boldsymbol{0} \text{ and } \tilde{\boldsymbol{u}}^{\top} \boldsymbol{X} = \boldsymbol{0},$$
  
$$\boldsymbol{F}_{1} \boldsymbol{F}_{1}^{\top} = \left( \boldsymbol{X}^{\top} \boldsymbol{X} \right)^{-1} + \left( \boldsymbol{X}^{\top} \boldsymbol{X} \right)^{-1} \boldsymbol{X}^{\top} \left( \boldsymbol{u} \, \boldsymbol{v} \right) \boldsymbol{C}_{\tilde{\boldsymbol{u}} \tilde{\boldsymbol{v}}}^{-1} \left( \boldsymbol{u}^{\top} \right) \boldsymbol{X} \left( \boldsymbol{X}^{\top} \boldsymbol{X} \right)^{-1}.$$
(S309)

(b)  $\boldsymbol{F}_2 \boldsymbol{F}_2^{\top}$ . Note that due to (S307)

$$(\boldsymbol{I} - \tilde{\boldsymbol{G}}) \left( \boldsymbol{I} - \boldsymbol{P}_{\left(\boldsymbol{X} \boldsymbol{u} \boldsymbol{v}\right)} \right) (\boldsymbol{I} - \tilde{\boldsymbol{G}})^{\top} = \left( \boldsymbol{I} - \boldsymbol{P}_{\left(\boldsymbol{X} \boldsymbol{u} \boldsymbol{v}\right)} \right).$$
(S310)

Combining with (S292), we have

$$(\boldsymbol{I} - \tilde{\boldsymbol{G}})(\boldsymbol{C}_1 + \boldsymbol{C}_2)(\boldsymbol{C}_1 + \boldsymbol{C}_2)^\top (\boldsymbol{I} - \tilde{\boldsymbol{G}})^\top$$
(S311)

$$=\frac{1}{d^2n}\left[\beta_u^2(\boldsymbol{I}-\boldsymbol{P}\boldsymbol{u})+\beta_v^2(\boldsymbol{I}-\boldsymbol{P}\boldsymbol{v})\right]$$
(S312)

$$-\left(\beta_{v}^{2}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{v}})\tilde{\boldsymbol{u}}\ \beta_{u}^{2}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{u}})\tilde{\boldsymbol{v}}\right)\boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1}\left(\boldsymbol{u}^{\top}\right)-\left(\boldsymbol{u}\ \boldsymbol{v}\right)\boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1}\left(\beta_{u}^{2}(\boldsymbol{I}-\boldsymbol{P}_{\boldsymbol{v}})\tilde{\boldsymbol{u}}\right)\qquad(S313)$$

$$+ (\boldsymbol{u}\,\boldsymbol{v}) \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \beta_{v}^{2} \tilde{\boldsymbol{u}}^{\top} (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \tilde{\boldsymbol{u}} & 0\\ 0 & \beta_{u}^{2} \tilde{\boldsymbol{v}}^{\top} (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) \tilde{\boldsymbol{v}} \end{pmatrix} \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \boldsymbol{u}^{\top}\\ \boldsymbol{v}^{\top} \end{pmatrix}$$
(S314)

$$-\frac{(\beta_u^2+\beta_v^2)(2\lambda d^2+\beta_u^2+\beta_v^2)}{(\lambda d^2+\beta_u^2+\beta_v^2)^2} \left(\boldsymbol{I}-\boldsymbol{P}_{\left(\boldsymbol{X}\boldsymbol{u}\boldsymbol{v}\right)}\right)\right].$$
(S315)

Because  $(I - P_{(Xuv)})X = 0, X^{\top}\tilde{u} = 0$  and  $X^{\top}\tilde{v} = 0$ ,

$$\boldsymbol{F}_{2}\boldsymbol{F}_{2}^{\top} \tag{S316}$$

$$= \frac{1}{d^2 n} (\boldsymbol{X}^{\top} \boldsymbol{X})^{-1} \boldsymbol{X}^{\top} \bigg[ \beta_u^2 (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) + \beta_v^2 (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}})$$
(S317)

$$+ \left(\boldsymbol{u}\,\boldsymbol{v}\right) \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \beta_{\boldsymbol{v}}^{2} \tilde{\boldsymbol{u}}^{\top} (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \tilde{\boldsymbol{u}} & 0\\ 0 & \beta_{\boldsymbol{u}}^{2} \tilde{\boldsymbol{v}}^{\top} (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) \tilde{\boldsymbol{v}} \end{pmatrix} \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \boldsymbol{u}^{\top}\\ \boldsymbol{v}^{\top} \end{pmatrix} \end{bmatrix} \boldsymbol{X} \left(\boldsymbol{X}^{\top} \boldsymbol{X}\right)^{-1}.$$
(S318)

Together with (S309) and (S318), we obtain the variance-covariance matrix of  $\delta_{\beta_x}$  as follows.

$$Cov\left(\delta_{\boldsymbol{\beta}_{x}}\right) = \sigma_{y}^{2} \left[ (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1} + (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top} \left(\boldsymbol{u} \ \boldsymbol{v}\right) \boldsymbol{C}_{\boldsymbol{\tilde{u}}\boldsymbol{\tilde{v}}}^{-1} \left( \boldsymbol{u}^{\top} \right) \boldsymbol{X} (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1} \right] +$$
(S319)

$$\sigma_a^2 \frac{1}{d^2 n} (\boldsymbol{X}^\top \boldsymbol{X})^{-1} \boldsymbol{X}^\top \left[ \beta_u^2 (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{u}}) + \beta_v^2 (\boldsymbol{I} - \boldsymbol{P}_{\boldsymbol{v}}) \right]$$
(S320)

+ 
$$(\boldsymbol{u} \boldsymbol{v}) \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \beta_{v}^{2} \tilde{\boldsymbol{u}}^{\top} (\boldsymbol{I} - \boldsymbol{P}_{v}) \tilde{\boldsymbol{u}} & 0\\ 0 & \beta_{u}^{2} \tilde{\boldsymbol{v}}^{\top} (\boldsymbol{I} - \boldsymbol{P}_{u}) \tilde{\boldsymbol{v}} \end{pmatrix} \boldsymbol{C}_{\tilde{\boldsymbol{u}}\tilde{\boldsymbol{v}}}^{-1} \begin{pmatrix} \boldsymbol{u}^{\top}\\ \boldsymbol{v}^{\top} \end{pmatrix} \end{bmatrix} \boldsymbol{X} (\boldsymbol{X}^{\top} \boldsymbol{X})^{-1}, (S321)$$

which is the same as  $Cov\left(\delta_{\beta_x}^{ts}\right)$  in (S143).

## REFERENCES

Magnus, J. R. and Neudecker, H. (1979). The commutation matrix: some properties and applications. *The Annals of Statistics*, pages 381–394.