

A NEYMAN-ORTHOGONALIZATION APPROACH TO THE INCIDENTAL PARAMETER PROBLEM*

Stéphane Bonhomme[†]

Department of Economics, University of Chicago

Koen Jochmans[‡]

Toulouse School of Economics, Université Toulouse Capitole

Martin Weidner[§]

Department of Economics and Nuffield College, University of Oxford

December 12, 2024

Abstract

A popular approach to perform inference on a target parameter in the presence of nuisance parameters is to construct estimating equations that are orthogonal to the nuisance parameters, in the sense that their expected first derivative is zero. Such first-order orthogonalization may, however, not suffice when the nuisance parameters are very imprecisely estimated. Leading examples where this is the case are models for panel and network data that feature fixed effects. In this paper, we show how, in the conditional-likelihood setting, estimating equations can be constructed that are orthogonal to any chosen order. Combining these equations with sample splitting yields higher-order bias-corrected estimators of target parameters. In an empirical application we apply our method to a fixed-effect model of team production and obtain estimates of complementarity in production and impacts of counterfactual re-allocations.

JEL Classification: C13, C23, C55.

Keywords: Neyman-orthogonality, incidental parameter, higher-order bias correction, networks.

*We are grateful to Dmitry Arkhangelsky, Bo Honoré, Roger Moon, Whitney Newey, Andres Santos, Vira Semenova, and Vasilis Syrgkanis for comments and discussion.

Funded by the European Union (ERC-NETWORK-101044319) and by the French Government and the French National Research Agency under the Investissements d’Avenir program (ANR-17-EURE-0010). Views and opinions expressed are however those of the authors only and do not necessarily reflect those of the European Union or the European Research Council. Neither the European Union nor the granting authority can be held responsible for them.

[†]sbonhomme@uchicago.edu

[‡]koen.jochmans@tse-fr.eu

[§]martin.weidner@economics.ox.ac.uk

1 Introduction

Inference in the presence of nuisance parameters has received substantial attention. One fruitful way to proceed is to work with estimating equations that are orthogonal with respect to the nuisance parameters in the sense of [Neyman \(1959\)](#). Such equations underlie much of the results in semiparametric estimation ([Newey, 1994](#)) and are at the heart of recent advances on doubly-robust estimation and high-dimensional inference as discussed in [Chernozhukov, Chetverikov, Demirer, Duflo, Hansen, Newey, and Robins \(2018\)](#) and [Chernozhukov, Escanciano, Ichimura, Newey and Robins \(2022\)](#), among others. A key finding is that Neyman-orthogonality permits the construction of asymptotically unbiased estimators that converge at the usual $n^{-1/2}$ -rate provided the nuisance parameter has a convergence rate that is faster than $n^{-1/4}$, where n is the sample size.

The faster-than- $n^{-1/4}$ requirement applies in a variety of semiparametric estimation problems; [Robinson \(1988\)](#) and [Ichimura \(1993\)](#) are examples. However, it often fails in problems where the dimension of the nuisance parameter is large relative to the sample size. Panel data models with fixed effects are an example. There, we observe N units over T periods of time and the model includes both common parameters and unit-specific nuisance parameters. The latter are estimated at the rate $T^{-1/2}$. For an estimator of the former based on Neyman-orthogonalization to be successful we would therefore need that $T^{-1/2} = o((NT)^{-1/4})$, which translates into the requirement that $N = o(T)$. This is usually not a realistic condition in microeconomic applications. In fact, under this requirement the standard fixed-effect estimator would permit asymptotically-valid inference. Consequently, (first-order) Neyman-orthogonalization does not solve the incidental parameter problem in panel data.¹

The issue can be even more severe in high-dimensional regressions on network data. In such settings, the convergence rate of the estimator of the nuisance parameter depends

¹The problem is reminiscent of the poor performance of double machine-learning techniques in some settings, as recently documented by [Wüthrich and Zhu \(2021\)](#) and [Angrist and Frandsen \(2022\)](#). A related problem where the conventional approach fails is in a nonlinear version of the judge-leniency design, as discussed in [Hahn and Hausman \(2021\)](#).

on the connectivity structure of the network (Jochmans and Weidner, 2019). Examples include the estimation of teacher value-added (Jackson, Rockoff and Staiger, 2014), of the contributions of worker and firm heterogeneity to the variance of log wages and other covariance components (Abowd, Kramarz and Margolis, 1999, Kline, Saggio and Sølvssten, 2020), as well as of complementarity patterns in team production (Ahmadpoor and Jones, 2019, Bonhomme, 2021). Fixed effects in network-formation models are also poorly estimated. This is especially true in the prevalent case where the network is sparse (see, e.g., Graham, 2017, 2020).

Motivated by these concerns, we are interested in a higher-order generalization of Neyman-orthogonality, in the sense of Mackey, Syrgkanis and Zadik (2018). Moreover, we show how ensuring orthogonality to higher order can successfully reduce bias in several panel and network models. We say that an estimating equation is Neyman-orthogonal to order q when all q leading derivatives with respect to the nuisance parameter have zero expectation. When $q = 1$, this means that the expected Jacobian is zero, and so we recover the conventional definition of Neyman-orthogonality (to order one). Working with estimating equations that are Neyman-orthogonal to order q , when combined with sample splitting, allows one to construct asymptotically-linear estimators when nuisance parameters are estimated at a rate faster than $n^{-1/2(q+1)}$. As an example, in the panel data problem, this reduces the bias from $O(T^{-1})$ down to $O(T^{-q})$, yielding valid inference under the requirement that $N = o(T^{2q-1})$. We remark that combining orthogonalization with sample splitting (or cross-fitting) is important to achieve such an improvement. Moreover, orthogonalized estimating equations, by themselves, do not, in general, deliver estimators with improved sampling properties.

We show how to construct estimating equations that are orthogonal to any chosen order in a general conditional-likelihood setting. These estimating equations can be understood to be generalizations of the projected score of Small and McLeish (1989) and Waterman and Lindsay (1996). They have an interpretation as higher-order influence functions, as introduced in Robins, Li, Tchetgen Tchetgen and van der Vaart (2008). Our approach applies to general low-dimensional target parameters that satisfy some moment restric-

tions. This includes functions of the nuisance parameters such as average elasticities or other average effects. The conditional-likelihood framework allows us to orthogonalize a given estimating equation without introducing additional nuisance parameters. As is well known, this is not essential to achieve orthogonality to order one. However, the absence of additional nuisance parameters turns out to be very helpful in enabling the construction of higher-order orthogonalized estimating equations.

We illustrate the usefulness of our approach in several examples and in an empirical application to the estimation of nonlinear regressions on network data; a problem for which, at present, no alternative solutions exist. In this setting, we estimate a constant elasticity of substitution (CES) production function from the scientific output of research collaborations. As in [Ahmadpoor and Jones \(2019\)](#), the production function depends on researcher-specific fixed effects. Estimates of the parameters can be used to quantify the degree of complementarity among researchers within teams, and to compute the impact of counterfactual re-allocations in the spirit of earlier work by [Graham, Imbens and Ridder \(2014\)](#).

This problem is difficult because in the data that we use (taken from [Ductor, Fafchamps, Goyal and Van der Leij, 2014](#) and concerning publications in economics on EconLit), the number of collaborations per researcher is quite low. A conventional estimator is thus likely to suffer from bias. Our procedure uncovers the presence of complementarity among authors in the production of research articles. In a counterfactual exercise we also find that randomly pairing researchers would lead to a decrease in the average quality of articles. Our findings are corroborated in a simulation experiment targeted to our empirical application.

2 Problem statement and motivation

2.1 Setup

Let $Z_i = (Y_i, X_i)$ be random vectors, for $i = 1, \dots, N$. We consider a setting where the conditional density function of Y_i given X_i , $\ell(y|x; \theta_0, \eta_{i0})$, is known up to the parameters

θ_0 and η_{0i} . Throughout, we will treat $\eta_{10}, \dots, \eta_{N0}$ as nuisance parameters, and leave the marginal density of the conditioning variable, $\ell_{X_i}(x)$, unrestricted. We are interested in estimating a parameter μ_0 that is defined through the moment condition

$$\sum_{i=1}^N \mathbb{E}(u_i(Z_i; \theta_0, \eta_{i0}, \mu_0)) = 0, \quad (2.1)$$

where the expectations are over Z_i under $\ell(y|x; \theta_0, \eta_{i0}) \ell_{X_i}(x)$. We assume that, for all $i = 1, \dots, N$, Z_i contains n_i individual observations, and denote the total number of observations as $n = \sum_{i=1}^N n_i$. For example, in a balanced panel data setting with N units and T time periods, Z_i is the time series of unit i 's observations, $n_i = T$ for all i , and $n = NT$.

Our setup accommodates different types of target parameters. As an example, we can set $\mu_0 = \theta_0$. In this case, using $u_i(z; \theta, \eta_i)$ as a shorthand for $u_i(z; \theta, \eta_i, \theta)$, one possibility is to use the score,

$$u_i(z; \theta, \eta_i) = \frac{\partial \log \ell(y|x; \theta, \eta_i)}{\partial \theta}.$$

More generally, the moment condition (2.1) defines the target parameter

$$\mu_0 = \mu(\theta_0, \eta_{10}, \dots, \eta_{N0}, \ell_{X_1}, \dots, \ell_{X_N}),$$

which can be a function of the parameters θ_0 and η_{i0} describing the conditional distribution of Y_i given X_i , and of the marginal distribution of X_i . For example, we may be interested in an average marginal effect of the form

$$\mu_0 = \sum_{i=1}^N \int m_i(x; \theta_0, \eta_{i0}) \ell_{X_i}(x) dx,$$

where m_1, \dots, m_N are known functions.

To illustrate the setup we will refer to two leading examples.

Example: Neyman-Scott model. Our first example is the well-known [Neyman and Scott \(1948\)](#) model. Here,

$$Y_{ij} = \eta_{i0} + \varepsilon_{ij}, \quad \varepsilon_{ij} \sim \text{iid } \mathcal{N}(0, \sigma_0^2), \quad i = 1, \dots, N, \quad j = 1, \dots, T, \quad (2.2)$$

where the goal is to estimate $\theta_0 = \sigma_0^2$ in the presence of the nuisance parameters $\eta_{10}, \dots, \eta_{N0}$. Define, for all $i = 1, \dots, N$,

$$u_i(Y_i; \sigma^2, \eta_i) = -\frac{T}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{j=1}^T (Y_{ij} - \eta_i)^2,$$

where $Y_i = (Y_{i1}, \dots, Y_{iT})^\top$ has dimension $n_i = T$, and the total number of observations is $n = NT$. It is well-known that the maximum-likelihood estimator of σ_0^2 is on average too small, suffering from bias $-\sigma_0^2/T$. While in this panel data problem first-order orthogonality does not reduce the order of this bias, we demonstrate below that second-order orthogonalization fully removes it.

Example: CES production function. Consider an environment where we observe workers producing output in n teams of size 2. Moreover, let $k(j, 1)$ and $k(j, 2)$ denote the workers in team j , and write $\mathcal{K} = \{(k(j, 1), k(j, 2)) : j = 1, \dots, n\}$ for the set of workers in all teams; note that a given worker may be part of multiple teams. Consider a model for team production where team output is a CES aggregate of worker inputs (as in [Ahmadpoor and Jones, 2019](#)),

$$Y_j = \left(\frac{\eta_{k(j,1)0}^{\gamma_0} + \eta_{k(j,2)0}^{\gamma_0}}{2} \right)^{\frac{1}{\gamma_0}} \varepsilon_j^{\sigma_0}, \quad \log \varepsilon_j | \mathcal{K} \sim \text{iid } \mathcal{N}(0, 1), \quad j = 1, \dots, n. \quad (2.3)$$

In this model, one may be interested in estimating the substitution parameter γ_0 or the log error variance σ_0^2 , average elasticities, or effects of counterfactual re-allocations of workers to teams, for example.

To analyze this example we consider $N \leq n$ subsets of teams j , of size n_i each. Let Y_i denote the vector of team outcomes in subset i , and let η_i be the collection of all fixed effects of workers belonging to those teams. Finally, let $\theta = (\gamma, \sigma^2)^\top$. The scores with respect to γ and σ^2 take the form $u_i(Y_i; \theta, \eta_i)$, where the dependence of u_i on i reflects that the set of workers who belong to the subset i of teams generally differs from the workers belonging to another subset. In contrast to our previous example, the theoretical literature on network models such as (2.3) is scarce, and to our knowledge no approach has as yet been developed for achieving bias reduction in such a setting.

2.2 The role of first-order orthogonality and its limitations

In the remainder of this section, we motivate our approach in a setting where one wishes to estimate $\mu_0 = \theta_0$ based on a random sample Z_1, \dots, Z_n , taking u to be a univariate function and η_0 to be a scalar. Hence $N = 1$, and $n_1 = n$ is the total number of observations.

If $\mathbb{E}(\sum_{j=1}^n u(Z_j; \theta_0, \eta_0)) = 0$, a conventional estimator of θ_0 , say $\hat{\theta}$, would be the solution to

$$\sum_{j=1}^n u(Z_j; \theta, \hat{\eta}) = 0,$$

where $\hat{\eta}$ is a consistent estimator of η_0 obtained in a preliminary step. However, it is well known that such a “plug-in” estimator is sensitive to the quality of the preliminary estimator $\hat{\eta}$ used.

Assuming sufficient regularity, a standard argument based on a linearization around θ_0 yields

$$\left(\mathbb{E} \left(\frac{\partial u(Z_j; \theta_0, \eta_0)}{\partial \theta^\top} \right) + o_P(1) \right) (\hat{\theta} - \theta_0) = \frac{1}{n} \sum_{j=1}^n u(Z_j; \theta_0, \hat{\eta}),$$

so that the sampling properties of $\hat{\theta} - \theta_0$ are dictated by the sampling properties of the estimating equation. We have

$$\begin{aligned} \frac{1}{n} \sum_{j=1}^n u(Z_j; \theta_0, \hat{\eta}) &= \underbrace{\frac{1}{n} \sum_{j=1}^n u(Z_j; \theta_0, \eta_0)}_{(A)} \\ &+ \underbrace{\left(\frac{1}{n} \sum_{j=1}^n \frac{\partial u(Z_j; \theta_0, \eta_0)}{\partial \eta} - \mathbb{E} \left(\frac{\partial u(Z_j; \theta_0, \eta_0)}{\partial \eta} \right) \right)}_{(B)} (\hat{\eta} - \eta_0) \\ &+ \underbrace{\mathbb{E} \left(\frac{\partial u(Z_j; \theta_0, \eta_0)}{\partial \eta} \right)}_{(C)} (\hat{\eta} - \eta_0) \\ &+ O_P(|\hat{\eta} - \eta_0|^2). \end{aligned} \tag{2.4}$$

The (A) term in (2.4) is a zero-mean sample average to which a standard central-limit theorem can be applied. Hence, it is generally $O_P(n^{-1/2})$. The next two terms in the expansion capture the first-order effect of estimation noise in $\hat{\eta}$. The (B) term can generally

be ensured to be $o_P(n^{-1/2})$. A generic approach to achieve this is to compute $\hat{\eta}$ from data that are independent of Z_1, \dots, Z_n , for example using sample splitting. In that case, (B) is the product of a sample average of zero-mean random variables—which is $O_P(n^{-1/2})$ —and an $o_P(1)$ term—as $\hat{\eta}$ is consistent for η_0 —and, therefore, (B) is $o_P(n^{-1/2})$. The (C) term, however, features a non-random Jacobian that, in general, is non-zero. Hence, (C) is $O_P(|\hat{\eta} - \eta_0|)$, and will only be asymptotically negligible when $\hat{\eta}$ is superconsistent for η_0 , which is not usually the case.

Suppose now that u is first-order orthogonal, in the sense that

$$\mathbb{E} \left(\frac{\partial u(Z_j; \theta_0, \eta_0)}{\partial \eta} \right) = 0. \quad (2.5)$$

Then the (C) term vanishes from (2.4) and we obtain

$$\frac{1}{n} \sum_{j=1}^n u(Z_j; \theta_0, \hat{\eta}) = \frac{1}{n} \sum_{j=1}^n u(Z_j; \theta_0, \eta_0) + O_P(|\hat{\eta} - \eta_0|^2) + o_P(n^{-1/2}). \quad (2.6)$$

The requirement that $\hat{\eta} - \eta_0 = o_P(n^{-1/4})$ then guarantees that the impact of the estimation error in $\hat{\eta}$ on $\hat{\theta}$ is asymptotically negligible. While a given function u does not, in general, satisfy (2.5), Neyman (1959) proposed a general method to transform it into one that does. The resulting function is said to be Neyman-orthogonal.

Condition (2.5) has a long history in semiparametric estimation problems (Bickel, 1982, Schick, 1986, Newey, 1994). More recently, it has proved to be a fundamental ingredient in the literature on high-dimensional inference (see Chernozhukov, Chetverikov, Demirer, Duflo, Hansen, Newey, and Robins, 2018 or Chernozhukov, Escanciano, Ichimura, Newey and Robins, 2022). There are, however, instances where it is ineffective. To illustrate this it suffices to consider the simple panel data setting from the Neyman and Scott (1948) problem.

Example: Neyman-Scott model (continued). In this problem it is easy to verify that

$$\mathbb{E} \left(\frac{\partial u_i(Y_i; \sigma_0^2, \eta_{i0})}{\partial \eta_i} \right) = -\frac{1}{\sigma_0^4} \sum_{j=1}^T \mathbb{E} (Y_{ij} - \eta_{i0}) = 0,$$

and so the score is already first-order Neyman-orthogonal with respect to the fixed effects. Nevertheless, given preliminary estimators $\hat{\eta}_1, \dots, \hat{\eta}_N$, and letting $\nu_i = \hat{\eta}_i - \eta_{i0}$, the estimator

$$\hat{\sigma}^2 = \frac{1}{NT} \sum_{i=1}^N \sum_{j=1}^T (Y_{ij} - \hat{\eta}_i)^2,$$

has expectation $\sigma_0^2 - 2/N \sum_{i=1}^N \mathbb{E}(\bar{\varepsilon}_i \nu_i) + 1/N \sum_{i=1}^N \mathbb{E}(\nu_i^2)$, for $\bar{\varepsilon}_i = 1/T \sum_{j=1}^T \varepsilon_{ij}$. Thus, when using sample splitting, the bias is $1/N \sum_{i=1}^N \mathbb{E}(\nu_i^2)$, the mean squared error of the preliminary estimator. With cross-fitting this is, at best, $O(T^{-1})$. Hence, $\sqrt{NT}(\hat{\sigma}^2 - \sigma_0^2)$ will not have a correctly-centered limit distribution unless $N/T \rightarrow 0$. However, under this condition, the joint maximum-likelihood estimator of σ_0^2 and the fixed effects, too, is asymptotically unbiased. Hence, having a score that is Neyman-orthogonal, even when combined with sample splitting, does not suffice to resolve the incidental parameter problem in panel data problems.

2.3 Higher-order orthogonality

To see how Neyman-orthogonality to a higher order can be helpful we now consider a further expansion of (2.4). Again assuming sufficient regularity, we have, for any integer $q \geq 1$,

$$\begin{aligned} \frac{1}{n} \sum_{j=1}^n u(Z_j; \theta_0, \hat{\eta}) &= \underbrace{\frac{1}{n} \sum_{j=1}^n u(Z_j; \theta_0, \eta_0)}_{(A)} \\ &+ \underbrace{\sum_{p=1}^q \frac{1}{p!} \left(\frac{1}{n} \sum_{j=1}^n \frac{\partial^p u(Z_j; \theta_0, \eta_0)}{\partial \eta^p} - \mathbb{E} \left(\frac{\partial^p u(Z_j; \theta_0, \eta_0)}{\partial \eta^p} \right) \right)}_{(B)} (\hat{\eta} - \eta_0)^p \\ &+ \underbrace{\sum_{p=1}^q \frac{1}{p!} \mathbb{E} \left(\frac{\partial^p u(Z_j; \theta_0, \eta_0)}{\partial \eta^p} \right)}_{(C)} (\hat{\eta} - \eta_0)^p \\ &+ O_P(|\hat{\eta} - \eta_0|^{q+1}). \end{aligned}$$

Here, the (A) term is the same as before. Also, with sample splitting we can again ensure that the (B) term will be asymptotically negligible. On the other hand, if the function u

satisfies the higher-order orthogonality condition

$$\mathbb{E} \left(\frac{\partial^p u(Z_j; \theta_0, \eta_0)}{\partial \eta^p} \right) = 0, \quad 1 \leq p \leq q, \quad (2.7)$$

the (C) term is equal to zero, and so

$$\frac{1}{n} \sum_{j=1}^n u(Z_j; \theta_0, \hat{\eta}) = \frac{1}{n} \sum_{j=1}^n u(Z_j; \theta_0, \eta_0) + O_P(|\hat{\eta} - \eta_0|^{q+1}) + o_P(n^{-1/2}). \quad (2.8)$$

Comparing (2.8) to (2.6) we see that the impact of estimation noise in $\hat{\eta}$ on our estimator of θ_0 has been reduced further. Moreover, for the impact of estimation error to be negligible, we now only require that $|\hat{\eta} - \eta_0|^{q+1} = o_P(n^{-1/2})$. It then follows from standard results that, as $n \rightarrow \infty$,

$$\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} \mathcal{N}(0, \Sigma_\theta)$$

for some Σ_θ , provided that

$$\hat{\eta} - \eta_0 = o_P(n^{-1/2(q+1)})$$

holds.

Example: Neyman-Scott model (continued) In the model of [Neyman and Scott \(1948\)](#),

$$\mathbb{E} \left(\frac{\partial^2 u_i(Y_i; \sigma_0^2, \eta_{i0})}{\partial \eta_i^2} \right) = \frac{T}{\sigma_0^4} \neq 0.$$

It thus follows that u_i is not orthogonal to second order (or to any order higher than two). Below we will show that a second-order Neyman-orthogonal score equation exists; its solution turns out to be

$$\hat{\sigma}^2 = \frac{1}{N(T-1)} \sum_{i=1}^N \sum_{j=1}^T (Y_{ij} - \bar{Y}_i)^2, \quad (2.9)$$

where $\bar{Y}_i = 1/T \sum_{j=1}^T Y_{ij}$. This estimator does not depend on the preliminary estimator of the fixed effects used. Furthermore, it performs the well-known degrees-of-freedom correction to the maximum-likelihood estimator, yielding an estimator that is exactly unbiased for any $T \geq 2$.

The notion of q th-order Neyman-orthogonality as in (2.7) was introduced by [Mackey, Syrgkanis and Zadik \(2018\)](#). In the context of our likelihood setup, we will give a general procedure to construct higher-order Neyman-orthogonal functions below.

3 Estimation based on orthogonalized functions

We now present our estimation approach in the general case where the target parameter μ_0 may be equal to θ_0 or may be a different parameter such as an average effect, and there are multiple, vector-valued nuisance parameters η_{i0} . We start by formally defining higher-order Neyman-orthogonality in this general setup and describe estimation based on higher-order Neyman-orthogonal moment functions. In the next section, we will then show how to construct such functions.

3.1 Definition of higher-order orthogonality

Let d_η be the dimension of η and write $\eta = (\eta_1, \dots, \eta_{d_\eta})$. For any non-negative integer p and a vector of integers $m = (m_1, \dots, m_p)$ satisfying $1 \leq m_s \leq d_\eta$ for all $1 \leq s \leq p$, define

$$D_\eta^m = \frac{\partial^p}{\partial \eta_{m_1} \cdots \partial \eta_{m_p}}. \quad (3.1)$$

For a given p , there are $d_p = \binom{d_\eta + p - 1}{p}$ unique such partial derivatives. Let $\nabla_\eta^{(p)}$ be the vector operator of dimension d_p that collects all these unique partial derivatives of order p . Finally, let ∇_η^q be the vector operator of dimension $\sum_{p=1}^q d_p$ obtained on stacking $\nabla_\eta^{(p)}$ for $p = 1, \dots, q$.

Neyman-orthogonality to order q can now be defined as follows ([Mackey, Syrgkanis and Zadik, 2018](#)).

Definition 1. *If the function u satisfies*

$$\mathbb{E} [\nabla_\eta^q u(Z; \theta_0, \eta_0, \mu_0)] = 0, \quad (3.2)$$

for some integer q , then we say that u is Neyman-orthogonal to order q .

In this definition, *all* possible partial derivatives of $u(Z; \theta, \eta, \mu)$ with respect to η up to order q have mean zero. Furthermore, Definition 1 is written for a generic function u . In later applications, we apply it to functions u_i that depend on some subsets of observations.

3.2 Estimation

Let μ_0 satisfy (2.1) for (possibly vector-valued) functions u_1, \dots, u_N . We assume that the u_i , for all $i = 1, \dots, N$, are Neyman-orthogonal to order q with respect to η_i , in the sense of Definition 1. Suppose that we have access to preliminary estimators $\hat{\eta}_1, \dots, \hat{\eta}_N$ of the nuisance parameters that are independent of the data Z_1, \dots, Z_N . If η_{i0} is defined as the solution to a moment condition involving the same data, estimation based on sample-splitting, combined with cross-fitting (see, e.g., Newey and Robins, 2017), can be applied. When the observations are independent this is conventional. For situations where the data are dependent, modified sample-splitting strategies are available (see, e.g., Semenova, Goldman, Chernozhukov and Taddy, 2023).

We estimate μ_0 by the GMM estimator

$$\hat{\mu} = \operatorname{argmin}_{\mu} \left\| \sum_{i=1}^N u_i(Z_i; \hat{\theta}, \hat{\eta}_i, \mu) \right\|_W, \quad (3.3)$$

where W is a chosen symmetric positive-definite matrix, $\|u\|_W = \sqrt{u^\top W u}$, and $\hat{\theta}$ is an estimator of θ_0 .

The estimator $\hat{\theta}$ will depend on the problem at hand. If θ_0 is defined through a moment condition of the form $\sum_{i=1}^N \mathbb{E}(\tilde{u}_i(Z_i; \theta_0, \eta_{i0})) = 0$, for functions \tilde{u}_i that are Neyman-orthogonal to order q , then our framework can be applied and we can use

$$\hat{\theta} = \operatorname{argmin}_{\theta} \left\| \sum_{i=1}^N \tilde{u}_i(Z_i; \theta, \hat{\eta}_i) \right\|_{\tilde{W}}, \quad (3.4)$$

where \tilde{W} is again a chosen weight matrix. In this case, we may equally combine (3.3) and (3.4) into a single GMM estimation procedure.

In Section 7, we provide conditions under which this approach yields estimators that are $n^{-1/2}$ -consistent and asymptotically normal, where n is the total number of observa-

tions. We will impose two key conditions. The first one is that, although their number may increase with the sample size, the dimension of each η_i remains bounded as n tends to infinity. This imposes a suitable sense of sparsity in the relationship between the nuisance parameters and the outcomes. This condition is trivially satisfied in the panel data and network problems with fixed effects that we consider. The second key condition we impose is that the convergence rates of the preliminary estimates $\hat{\eta}_i$ be faster than $n^{-1/2(q+1)}$. This ensures that, after having orthogonalized to order q , any remainder terms are asymptotically negligible.

4 Achieving higher-order Neyman-orthogonality

4.1 Main result

Let u be a moment function, such as one of the u_i in (2.1). We now show how to construct an orthogonalized counterpart of u , which we call u_q^* , that is Neyman-orthogonal to order q , where $q \geq 1$ is any arbitrary order.

It is convenient to introduce the [Bhattacharyya \(1946\)](#) basis v_1, v_2, \dots , where

$$v_p(z; \theta, \eta) = \frac{\nabla_{\eta}^{(p)} \ell(y | x; \theta, \eta)}{\ell(y | x; \theta, \eta)}.$$

[McLeish and Small \(1994\)](#) discuss several properties of this basis. One important property for our purposes is that

$$\mathbb{E}_{\theta, \eta}(v_p(Z; \theta; \eta) | X = x) = \int v_p(z; \theta, \eta) \ell(y | x; \theta, \eta) dy = 0 \quad (4.1)$$

for any p , so all elements of the Battacharryya basis have (conditional) mean equal to zero. In (4.1), and throughout this section, $\mathbb{E}_{\theta, \eta}(\cdot | X = x)$ denotes the conditional expectation under $\ell(y | x; \theta, \eta)$.

The low-order basis functions are familiar from likelihood theory. For example,

$$\begin{aligned} v_1(z; \theta, \eta) &= \frac{\partial \log \ell(y | x; \theta, \eta)}{\partial \eta}, \\ v_2(z; \theta, \eta) &= \frac{\partial \log \ell(y | x; \theta, \eta)}{\partial \eta} \frac{\partial \log \ell(y | x; \theta, \eta)}{\partial \eta^\top} + \frac{\partial^2 \log \ell(y | x; \theta, \eta)}{\partial \eta \partial \eta^\top}. \end{aligned}$$

The fact that these functions have mean zero follows from the unbiasedness of the score and from the information equality, respectively.

Stacking the leading q basis functions, we obtain

$$w_q(z; \theta, \eta) = \frac{\nabla_\eta^q \ell(y | x; \theta, \eta)}{\ell(y | x; \theta, \eta)}.$$

The vectors w_q are mean-zero “generalized score functions”.

Next, let us define the matrices

$$\Sigma_{w_q w_q}(x; \theta, \eta) = \mathbb{E}_{\theta, \eta}(w_q(Z; \theta, \eta) w_q(Z; \theta, \eta)^\top | X = x),$$

and

$$\Sigma_{w_q u}(x; \theta, \eta, \mu) = \mathbb{E}_{\theta, \eta}(w_q(Z; \theta, \eta) u(Z; \theta, \eta, \mu)^\top | X = x),$$

which are, respectively, the (conditional) covariance matrix of the first q members of the Bhattacharyya basis, and the covariance matrix of the same q basis functions with the vector function u . Finally, let

$$b_q(x; \theta, \eta, \mu) = \nabla_\eta^q \mathbb{E}_{\theta, \eta}(u(Z; \theta, \eta, \mu)^\top | X = x).$$

Note that b_q is zero when u is the score for θ , i.e., $\frac{\partial \log \ell(y | x; \theta, \eta)}{\partial \theta}$. In general, however, b_q will be non-zero. Here we assume that $u(z; \theta, \eta, \mu)$ and $\ell(y | x; \theta, \eta)$ are sufficiently often differentiable in η , and that the expectations in the definitions of $\Sigma_{w_q w_q}$, $\Sigma_{w_q u}$, and b_q are well-defined.

The proof of the following result is in Appendix [A](#).

Theorem 1. *Suppose that $\Sigma_{w_q w_q}(x; \theta, \eta)$ is invertible and let*

$$A(x; \theta, \eta, \mu) = \Sigma_{w_q w_q}(x; \theta, \eta)^{-1} [\Sigma_{w_q u}(x; \theta, \eta, \mu) - b_q(x; \theta, \eta, \mu)].$$

Then the function

$$u_q^*(z; \theta, \eta, \mu) = u(z; \theta, \eta, \mu) - A(x; \theta, \eta, \mu)^\top w_q(z; \theta, \eta)$$

satisfies $\mathbb{E}_{\theta, \eta} [\nabla_\eta^q u_q^(Z; \theta, \eta, \mu) | X = x] = 0$. This implies that u_q^* is Neyman-orthogonal to order q , as defined above.*

Theorem 1 generalizes the projected-score construction of [Small and McLeish \(1989\)](#) and [Waterman and Lindsay \(1996\)](#). To see this, consider the case where $\mu_0 = \theta_0$, and u is the score function for θ . Then $b_q = 0$ and Theorem 1 yields

$$u_q^*(z; \theta, \eta, \mu) = u(z; \theta, \eta, \mu) - [\Sigma_{w_q w_q}(x; \theta, \eta)^{-1} \Sigma_{w_q u}(x; \theta, \eta)]^\top w_q(z; \theta, \eta),$$

which is the projected score of order q .² However, our result covers other estimating equations as well as more general parameters of interest, such as average elasticities or counterfactual quantities.

We remark that Theorem 1 requires the matrix $\Sigma_{w_q w_q}(x; \theta, \eta)$ to be invertible. In the standard case of first-order Neyman-orthogonality this corresponds to non-singularity of the information matrix of the nuisance parameters. For higher-order Neyman-orthogonality this requirement imposes further restrictions.

Example: CES production function (continued). Consider the team production model (2.3). Suppose we work with $N = n$ subsets that all contain a single team, and let $\eta_i = (\eta_{k(i,1)}, \eta_{k(i,2)})^\top$ denote the 2×1 vector of worker effects in team i . The 2×2 matrix

$$\Sigma_{w_1 w_1}(\theta, \eta_s) = \mathbb{E}_{\theta, \eta_i} \left(\frac{\partial \log \ell(Y_i; \theta, \eta_i)}{\partial \eta_i} \frac{\partial \log \ell(Y_i; \theta, \eta_i)}{\partial \eta_i^\top} \right)$$

is not invertible, as only the sum $\eta_{k(i,1)}^\gamma + \eta_{k(i,2)}^\gamma$ can be identified. In our application in Section 6, we will tackle this issue by combining data on teams of size 2 with single-author production, and working with subsets i of three teams each.

Example: Fixed-effect probit. Consider the standard binary-choice panel data model

$$\mathbb{P}_{\theta, \eta_i}(Y_{ij} = 1 | X_i) = \Phi(\eta_i + X_{ij}^\top \theta), \quad i = 1, \dots, N, \quad j = 1, \dots, T,$$

²The projected score was originally developed as a tool to achieve E-ancillarity ([Small and McLeish, 1988](#)) and to approximate the conditional score for θ , when the latter exists ([Waterman and Lindsay, 1996](#)). The fact that it is Neyman-orthogonal is noted in passing (although a link with Neyman's work is not made) but is not exploited. Moreover, unlike the conditional score, the projected score still depends on η , and it will generally not have improved properties over the score itself. As we highlight here, it is the combination of higher-order versions of Neyman-orthogonality with sample splitting that allows one to improve over working with the original score.

for (conditionally-independent) binary outcomes Y_{ij} and covariates $X_i = (X_{i1}, \dots, X_{iT})^\top$. It is not difficult to see that the rank of $\Sigma_{w_q w_q}(x; \theta, \eta)$ is bounded by 2^T . As a result, $\Sigma_{w_q w_q}(x; \theta, \eta)$ is singular for all $q > 2^T$.

4.2 Intuition and discussion

To gain intuition into the construction in Theorem 1 it is useful to again consider the case where u is a univariate function, the nuisance parameter is a scalar, and one wishes to estimate $\mu_0 = \theta_0$ (as in Subsections 2.2 and 2.3).

First-order orthogonality. To relate our approach to the literature consider first $q = 1$. Let

$$u_1^*(z; \theta, \eta) = u(z; \theta, \eta) - a_1(x; \theta, \eta) v_1(z; \theta, \eta),$$

for some function a_1 . Note that, by virtue of (4.1), the term involving v_1 does not introduce any bias. We have

$$\frac{\partial u_1^*(z; \theta, \eta)}{\partial \eta} = \frac{\partial u(z; \theta, \eta)}{\partial \eta} - \frac{\partial a_1(x; \theta, \eta)}{\partial \eta} v_1(z; \theta, \eta) - a_1(x; \theta, \eta) \frac{\partial v_1(z; \theta, \eta)}{\partial \eta}.$$

Take conditional expectations and exploit (4.1) to see that

$$\mathbb{E}_{\theta, \eta} \left(\frac{\partial u_1^*(Z; \theta, \eta)}{\partial \eta} \middle| X = x \right) = 0$$

if and only if

$$\mathbb{E}_{\theta, \eta} \left(\frac{\partial u(Z; \theta, \eta)}{\partial \eta} \middle| X = x \right) - a_1(x; \theta, \eta) \mathbb{E}_{\theta, \eta} \left(\frac{\partial v_1(Z; \theta, \eta)}{\partial \eta} \middle| X = x \right) = 0.$$

This is achieved by setting

$$a_1(x; \theta, \eta) = \left(\mathbb{E}_{\theta, \eta} \left(\frac{\partial v_1(Z; \theta, \eta)}{\partial \eta} \middle| X = x \right) \right)^{-1} \mathbb{E}_{\theta, \eta} \left(\frac{\partial u(Z; \theta, \eta)}{\partial \eta} \middle| X = x \right). \quad (4.2)$$

Iterating expectations shows that the resulting function u_1^* is Neyman-orthogonal to order $q = 1$. By the information matrix equality we have

$$\mathbb{E}_{\theta, \eta} \left(\frac{\partial v_1(Z; \theta, \eta)}{\partial \eta} \middle| X = x \right) = -\mathbb{E}_{\theta, \eta} (v_1(Z; \theta, \eta)^2 | X = x) = -\Sigma_{w_1 w_1}(x; \theta, \eta),$$

and

$$\mathbb{E}_{\theta,\eta} \left(\frac{\partial u(Z; \theta, \eta)}{\partial \eta} \middle| X = x \right) = -\mathbb{E}_{\theta,\eta} (v_1(Z; \theta, \eta)u(Z; \theta, \eta) | X = x) = -\Sigma_{w_1 u}(x; \theta, \eta),$$

leading to the representation of the function u_1^* as in the theorem.

The above derivation of (4.2) is well-known. Furthermore, it does not hinge on the likelihood structure. Indeed, recent work exploiting orthogonality, such as that surveyed in Chernozhukov, Chetverikov, Demirer, Duflo, Hansen, Newey, and Robins (2018), does so in the context of moment conditions. In our setup, as in Neyman's (1959) original work, the likelihood setup implies that a_1 is known up to the model parameters θ and η (conditional on the regressors). Outside of this framework, in contrast, a_1 needs to be treated as an additional nuisance parameter. This is possible because, as u_1^* is linear in a_1 , it is automatically first-order Neyman-orthogonal to it by virtue of (4.1). This logic, however, does not extend to higher order, as the implied system of equations becomes inconsistent, so that no solution exists.

Higher-order orthogonality. It suffices to look at the case $q = 2$. We again consider a linear transformation of u , now involving the leading two Bhattacharyya basis functions. This gives

$$u_2^*(z; \theta, \eta) = u(z; \theta, \eta) - \begin{pmatrix} a_{21}(x; \theta, \eta) \\ a_{22}(x; \theta, \eta) \end{pmatrix}^\top \begin{pmatrix} v_1(z; \theta, \eta) \\ v_2(z; \theta, \eta) \end{pmatrix}. \quad (4.3)$$

Taking first-derivatives with respect to the nuisance parameter, and proceeding as in the first-order case, gives

$$\mathbb{E}_{\theta,\eta} \left(\frac{\partial u(Z; \theta, \eta)}{\partial \eta} \middle| X = x \right) = \begin{pmatrix} a_{21}(x; \theta, \eta) \\ a_{22}(x; \theta, \eta) \end{pmatrix}^\top \begin{pmatrix} \mathbb{E}_{\theta,\eta} \left(\frac{\partial v_1(Z; \theta, \eta)}{\partial \eta} \middle| X = x \right) \\ \mathbb{E}_{\theta,\eta} \left(\frac{\partial v_2(Z; \theta, \eta)}{\partial \eta} \middle| X = x \right) \end{pmatrix}.$$

Solving this equation for a_{21} for given a_{22} yields

$$a_{21}(x; \theta, \eta) = a_1(x; \theta, \eta) - c_1(x; \theta, \eta) a_{22}(x; \theta, \eta), \quad (4.4)$$

where a_1 is given by (4.2) and

$$c_1(x; \theta, \eta) = \left(\mathbb{E}_{\theta,\eta} \left(\frac{\partial v_1(z; \theta, \eta)}{\partial \eta} \middle| X = x \right) \right)^{-1} \mathbb{E}_{\theta,\eta} \left(\frac{\partial v_2(z; \theta, \eta)}{\partial \eta} \middle| X = x \right).$$

The coefficient c_1 has the same form as a_1 , except that it features v_2 instead of u . Moreover, plugging (4.4) back into (4.3) yields

$$u_2^*(z; \theta, \eta) = u_1^*(z; \theta, \eta) - a_{22}(x; \theta, \eta) v_2^*(z; \theta, \eta),$$

where $v_2^*(z; \theta, \eta) = v_2(z; \theta, \eta) - c_1(x; \theta, \eta) v_1(z; \theta, \eta)$. Note that v_2^* is Neyman-orthogonal to order 1, that is,

$$\mathbb{E}_{\theta, \eta} \left(\frac{\partial v_2^*(Z; \theta, \eta)}{\partial \eta} \middle| X = x \right) = 0.$$

It follows that u_2^* is Neyman-orthogonal to order 1 for any a_{22} . We will now choose a_{22} such that u_2^* is Neyman-orthogonal to order 2.

Next, differentiating u_2^* with respect to η twice gives

$$\begin{aligned} \frac{\partial^2 u_2^*(z; \theta, \eta)}{\partial \eta^2} &= \frac{\partial^2 u_1^*(z; \theta, \eta)}{\partial \eta^2} + a_{22}(x; \theta, \eta) \frac{\partial^2 v_2^*(z; \theta, \eta)}{\partial \eta^2} \\ &\quad + \frac{\partial^2 a_{22}(x; \theta, \eta)}{\partial \eta^2} v_2^*(z; \theta, \eta) + 2 \frac{\partial a_{22}(x; \theta, \eta)}{\partial \eta} \frac{\partial v_2^*(z; \theta, \eta)}{\partial \eta}. \end{aligned}$$

Since v_2^* is orthogonal to order 1, the terms involving the first and second derivative of a_{22} drop out when taking expectations. It follows that u_2^* in (4.3) is Neyman-orthogonal to order 2 when one sets a_{21} to its expression in (4.4), and a_{22} to

$$a_{22}(x; \theta, \eta) = \left(\mathbb{E}_{\theta, \eta} \left(\frac{\partial^2 v_2^*(Z; \theta, \eta)}{\partial \eta^2} \middle| X = x \right) \right)^{-1} \mathbb{E}_{\theta, \eta} \left(\frac{\partial^2 u_1^*(Z; \theta, \eta)}{\partial \eta^2} \middle| X = x \right). \quad (4.5)$$

Note that this construction amounts to solving a system of linear equations. The fact that the solution in (4.4)–(4.5) coincides with the expression in Theorem 1 may then again be verified by using Bartlett identities.

5 Examples

5.1 Panel data models

Consider an $N \times T$ panel data model with individual effects. Here, the likelihood factors across the cross-sectional observations and the likelihood contribution of unit i takes the

form

$$\prod_{j=1}^T \log f(Y_{ij} | X_{ij}; \theta_0, \eta_{i0}).$$

The maximum-likelihood estimator is well-known to suffer from a bias that is $O(T^{-1})$; see [Hahn and Newey \(2004\)](#) and [Hahn and Kuersteiner \(2011\)](#) for derivations of this bias in static and dynamic models, respectively. Consider the estimation of θ_0 . The bias in the estimator comes from bias in the score stemming from estimation noise in the fixed effects. Taking η_i to be scalar for notational simplicity, and letting $\hat{\eta}_i$ be an estimator of η_{i0} , an expansion of the (normalized) score³

$$u_i(Z_i; \theta_0, \hat{\eta}_i) = \frac{1}{T} \sum_{j=1}^T \frac{\partial \log f(Y_{ij} | X_{ij}; \theta_0, \hat{\eta}_i)}{\partial \theta}$$

yields

$$\begin{aligned} u_i(Z_i; \theta_0, \hat{\eta}_i) &= u_i(Z_i; \theta_0, \eta_{i0}) + \frac{\partial u_i(Z_i; \theta_0, \eta_{i0})}{\partial \eta_i} (\hat{\eta}_i - \eta_{i0}) + \frac{1}{2} \frac{\partial^2 u_i(Z_i; \theta_0, \eta_{i0})}{\partial \eta_i^2} (\hat{\eta}_i - \eta_{i0})^2 \\ &\quad + o_P(|\hat{\eta}_i - \eta_{i0}|^2). \end{aligned}$$

Taking expectations and re-arranging shows that

$$\begin{aligned} \mathbb{E}(u_i(Z_i; \theta_0, \hat{\eta}_i)) &= \text{cov} \left(\frac{\partial u_i(Z_i; \theta_0, \eta_{i0})}{\partial \eta_i}, \hat{\eta}_i - \eta_{i0} \right) \\ &\quad + \mathbb{E} \left(\frac{\partial u_i(Z_i; \theta_0, \eta_{i0})}{\partial \eta_i} \right) \mathbb{E}(\hat{\eta}_i - \eta_{i0}) \\ &\quad + \frac{1}{2} \mathbb{E} \left(\frac{\partial^2 u_i(Z_i; \theta_0, \eta_{i0})}{\partial \eta_i^2} \right) \mathbb{E}((\hat{\eta}_i - \eta_{i0})^2) + o(\mathbb{E}(|\hat{\eta}_i - \eta_{i0}|^2)). \end{aligned}$$

If we set $\hat{\eta}_i = \hat{\eta}_i(\theta_0) = \arg \max_{\eta} \prod_{j=1}^T \log f(Y_{ij} | X_{ij}; \theta_0, \eta)$, the maximum-likelihood estimator given θ_0 , each one of these terms is $O(T^{-1})$. If we use an estimator $\hat{\eta}_i$ that is independent of the data the first term disappears. However, the remaining terms, which capture the nonlinearity bias and variance in the estimator of η_{i0} , remain. [Hahn and Newey \(2004\)](#), [Arellano and Hahn \(2007\)](#), and [Dhaene and Jochmans \(2015a,b\)](#) present estimators

³In this discussion we work with the score divided by T , to facilitate the comparison with the panel data literature.

of these terms based on the maximum-likelihood estimator that can be used to construct a bias-corrected estimator.

Lancaster (2002) and Woutersen (2002) integrate-out the fixed effects using a uniform prior after orthogonalizing to order 1 to obtain an estimator with bias $o(T^{-1})$; Arellano (2003) presents an alternative derivation of the same result. First-order orthogonality, by itself, does not suffice as it does not handle the third term in the expansion. Moreover, such an approach does not properly correct for the noise in the estimated fixed effects, yielding an estimator with a bias that is, at best, of the same order of magnitude as the (uncorrected) maximum-likelihood estimator itself. Li, Lindsay and Waterman (2003), building on Waterman and Lindsay (1996), show that their (second-order) projected score for θ , when evaluated at $\hat{\eta}_i(\theta)$, is a first-order unbiased estimating equation for θ . Thus, here, a sample-splitting procedure is not needed to achieve bias reduction. This is a consequence of the (second- or higher-order) projected score being orthogonal to the influence function of $\hat{\eta}_i(\theta)$, as a small calculation will allow to verify. While interesting, this property does not seem to extend to higher-order projections or to other parameters of interest, such as average marginal effects.

More generally, with $\hat{\eta}_i - \eta_{i0} = O_P(T^{-1/2})$, the score admits a higher-order expansion of the form,

$$\mathbb{E}(u_i(Z_i; \theta_0, \hat{\eta}_i)) = \frac{B_1}{T} + \frac{B_2}{T^2} + \dots + \frac{B_q}{T^q} + o(T^{-q})$$

for constants B_1, B_2, \dots, B_q . The maximum-likelihood estimator has $B_1 \neq 0$, in general, and so requires that $N/T \rightarrow 0$ to be asymptotically unbiased. The approaches to bias correction mentioned above remove B_1 but not the remaining terms. Approaches that estimate and subsequently remove all B_p , $1 \leq p \leq q$, are given by Dhaene and Jochmans (2015a,b). Likewise, an estimator based on Neyman-orthogonalization, combined with a sample-splitting estimator that uses preliminary estimators that satisfy $\hat{\eta}_i - \eta_{i0} = O_P(T^{-1/2})$ can be used to obtain the same result.

Example: Neyman-Scott model (continued). Recall that the (un-normalized) unit-specific score for σ^2 is

$$u_i(Y_i; \sigma^2, \eta_i) = -\frac{T}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{j=1}^T (Y_{ij} - \eta_i)^2.$$

The leading two elements of the Bhattacharyya basis for η_i are

$$v_{i,1}(Y_i; \sigma^2, \eta_i) = \sum_{j=1}^T \frac{Y_{ij} - \eta_i}{\sigma^2}, \quad v_{i,2}(Y_i; \sigma^2, \eta_i) = -\frac{T}{\sigma^2} + \left(\sum_{j=1}^T \frac{Y_{ij} - \eta_i}{\sigma^2} \right)^2.$$

We apply Theorem 1. A small calculation yields $A(\sigma^2, \eta_i) = (0, 1/2T)^\top$ and, after rearranging,

$$u_{i,2}^*(Y_i; \sigma^2, \eta_i) = \frac{1}{2\sigma^2} \left(\frac{\sum_{j=1}^T (Y_{ij} - \bar{Y}_i)}{\sigma^2} - (T - 1) \right),$$

which does not depend on η_i . Summing over the cross-sectional units gives the second-order orthogonalized score equation for σ^2 as

$$\sum_{i=1}^N u_{i,2}^*(Y_i; \sigma^2, \eta_i) = \frac{1}{2\sigma^2} \left(\frac{\sum_{i=1}^N \sum_{j=1}^T (Y_{ij} - \bar{Y}_i)}{\sigma^2} - N(T - 1) \right) = 0,$$

which yields the degrees-of-freedom corrected estimator $\hat{\sigma}^2$ in (2.9).

Another parameter of interest in this problem would be $\mu = 1/N \sum_{i=1}^N \eta_i^2$. This fits our framework with

$$u_i(Y_i; \sigma^2, \eta_i, \mu) = \eta_i^2 - \mu.$$

Here, the solution to the second-order orthogonalized moment equation for a given σ^2 is $1/N \sum_{i=1}^N \bar{Y}_i^2 - \sigma^2/T$. An unbiased estimator based on this equation then is $1/N \sum_{i=1}^N \bar{Y}_i^2 - \hat{\sigma}^2/T$.

To complement this example we consider in Appendix B the normal regression model

$$Y_i = X_i^\top \eta_{i0} + \varepsilon_i, \quad \varepsilon_i | X \sim \text{iid } \mathcal{N}(0, \sigma_0^2), \quad (5.1)$$

and show that second-order Neyman-orthogonalization similarly delivers exactly unbiased estimators of σ_0^2 , and of quadratic forms in $\eta_{10}, \dots, \eta_{N0}$.

Example: Linear autoregression. As another panel data example we provide results for the linear autoregressive model

$$Y_{ij} = \eta_{i0} + \rho_0 Y_{i,j-1} + \varepsilon_{ij}, \quad \varepsilon_{ij} \sim \text{iid } \mathcal{N}(0, \sigma_0^2).$$

Here $\theta = (\rho, \sigma^2)^\top$. We focus on ρ , since the analysis for σ^2 is similar to the previous example. The score for ρ for unit i , conditional on the first observation, is

$$u_i(Y_i; \theta, \eta_i) = \sum_{j=1}^T \frac{Y_{i,j-1}(Y_{ij} - \eta_i - \rho Y_{i,j-1})}{\sigma^2}$$

while

$$v_{i,1}(Y_i; \theta, \eta_i) = \sum_{j=1}^T \frac{(Y_{ij} - \eta_i - \rho Y_{i,j-1})}{\sigma^2}, \quad v_{i,2}(Y_i; \theta, \eta_i) = -\frac{T}{\sigma^2} + \left(\sum_{j=1}^T \frac{(Y_{ij} - \eta_i - \rho Y_{i,j-1})}{\sigma^2} \right)^2.$$

We find

$$A(\theta, \eta_i) = (\eta_i, \sigma^2/T)^\top c(\rho), \quad c(\rho) = \frac{1}{1-\rho} \left(1 - \frac{1}{T} \frac{1-\rho^T}{1-\rho} \right).$$

After some re-arrangement we obtain that the second-order Neyman-orthogonalized score equation takes the form

$$\frac{\sum_{i=1}^N \sum_{j=1}^T Y_{i,j-1}(Y_{ij} - \eta_i - \rho Y_{i,j-1})}{\sigma^2} + Nc(\rho) + NTc(\rho) \hat{\eta}_i(\rho)(\eta_i - \hat{\eta}_i(\rho)),$$

where $\hat{\eta}_i(\rho) = \bar{Y}_i - \rho \bar{Y}_{i-}$ with $\bar{Y}_i = 1/T \sum_{j=1}^T Y_{ij}$ and $\bar{Y}_{i-} = 1/T \sum_{j=1}^T Y_{i,j-1}$. This equation still depends on the η_i . However, at $\eta_i = \hat{\eta}_i(\rho)$ we obtain the adjusted score equation of [Lancaster \(2002\)](#) and [Dhaene and Jochmans \(2016\)](#), which is known to be exactly unbiased for any $T \geq 2$.

5.2 Nonlinear network regression

Our next example is the nonlinear regression model with $S \geq 1$ outcomes,

$$Y_i = m(X_i; \theta_0, \eta_{i0}) + \sigma(X_i; \theta_0) \varepsilon_i, \quad \varepsilon_i | X \sim \text{iid } \mathcal{N}(0, I_d), \quad (5.2)$$

where $m(x; \theta, \eta_i)$ is a $d \times 1$ vector, $\sigma(x; \theta)$ is an $d \times d$ diagonal matrix, and m and σ are known functions. We will show below that our CES production function example, in logarithms, fits into this framework.

For this model there are no analytical solutions for the orthogonalized estimators. We thus proceed numerically. To construct Neyman-orthogonal moment functions according to Theorem 1 we need to compute $\Sigma_{w_q w_q}(x; \theta, \eta)$, $\Sigma_{w_q u}(x; \theta, \eta, \mu)$, and $b_q(x; \theta, \eta, \mu)$, which involve higher-order derivatives of the conditional likelihood. To compute these derivatives, it is convenient to introduce the operator ∇_m^q that collects all derivatives with respect to m up to order q . By the chain rule,

$$\nabla_{\eta_i}^q \ell(y | x; \theta, \eta_i) = M(x, \theta, \eta_i) \nabla_m^q \ell(y | x; \theta, \eta_i),$$

where the matrix M has an analytical expression given by the multivariate Faà di Bruno formula (Constantine and Savits, 1996). Given the matrix M it is easy to compute $\Sigma_{w_q w_q}$, $\Sigma_{w_q u}$, and b_q . For example,

$$\begin{aligned} & \Sigma_{w_q w_q}(x; \theta, \eta_i) \\ &= M(x, \theta, \eta_i) \mathbb{E}_{\theta, \eta_i} \left(\frac{\nabla_m^q \ell(Y_i | X_i; \theta, \eta_i)}{\ell(Y_i | X_i; \theta, \eta_i)} \frac{\nabla_m^q \ell(Y_i | X_i; \theta, \eta_i)^\top}{\ell(Y_i | X_i; \theta, \eta_i)} \Big| X_i = x \right) M(x, \theta, \eta_i)^\top, \end{aligned}$$

where the expectation on the right-hand can be readily computed by relying on formulas for moments of Hermite polynomials. We relegate further details to Appendix C. In the next section we present simulations and an empirical application based on a version of (5.2) designed to study team production.

Example: CES production function (continued). Consider the team production model

$$Y_j = \beta_0(s_j) \left(\frac{1}{s_j} \sum_{r=1}^{s_j} \eta_{k(j,r)0}^{\gamma_0(s_j)} \right)^{\frac{1}{\gamma_0(s_j)}} \varepsilon_j^{\sigma_0(s_j)}, \quad \log \varepsilon_j | \mathcal{K} \sim \text{iid } \mathcal{N}(0, 1), \quad (5.3)$$

where s_j is the size of team $j = 1, \dots, n$, $(k(j, 1), \dots, k(j, s_j))$ are the s_j workers in team j , and the set $\mathcal{K} = \{k(j, r) : r = 1, \dots, s_j, j = 1, \dots, n\}$ collects the workers in all teams. Model (5.3) generalizes Model (2.3) by allowing for teams of varying sizes. Here we focus on teams of size 1 and 2, as in our application, and impose the normalization $\beta_0(1) = 1$. For simplicity we will denote $\beta_0 = \beta_0(2)$ and $\gamma_0 = \gamma_0(2)$, which are the team size and substitution parameters, respectively, in teams of size 2.

We now explain how (5.3) can be written as a special case of (5.2), for a suitable choice of subsets of observations. To any team j of size 2 involving workers k and k' , we associate a team $j_1(j)$ of size 1 only involving worker k , and a team $j_2(j)$ of size 1 only involving worker k' . This construction results in N subsets of three teams each. We then write the outcomes for these three teams, in logarithms, as

$$\log Y_j = \log \beta_0 + \frac{1}{\gamma_0} \log \left(\frac{\eta_{k(j,1)0}^{\gamma_0} + \eta_{k(j,2)0}^{\gamma_0}}{2} \right) + \sigma_0(2) \log \varepsilon_j, \quad (5.4)$$

$$\log Y_{j_1(j)} = \log \eta_{k(j,1)0} + \sigma_0(1) \log \varepsilon_{j_1(j)}, \quad (5.5)$$

$$\log Y_{j_2(j)} = \log \eta_{k(j,2)0} + \sigma_0(1) \log \varepsilon_{j_2(j)}, \quad (5.6)$$

which takes the same form as (5.2), for $d = 3$, $\theta = (\beta_0, \gamma_0, \sigma_0^2(1), \sigma_0^2(2))^\top$, Y_i the vector of the three outcomes in (5.4)–(5.6) for subset i , and η_{i0} the 2×1 vector of worker-specific effects in the corresponding teams.

6 Application to team production

6.1 Model, data, and implementation

We wish to estimate the parameters of the team production model in (5.4)–(5.6). We will be especially interested in estimating the substitution parameter γ , which drives the nature of complementarities in the team of size 2, and the team size parameter β , which reflects the premium (or penalty) associated with working together relative to working alone. In addition to estimating production-function parameters, we will also report estimates of a counterfactual random re-allocation of workers to teams. Under random assignment, average output in teams of size 2 can be written as

$$\mathbb{E}^{\text{rand}}(Y_j) = \frac{2}{n_2(n_2 - 1)} \sum_{k_1 < k_2} \beta_0 \left[\frac{1}{2} (\eta_{k_1 0}^{\gamma_0} + \eta_{k_2 0}^{\gamma_0}) \right]^{\frac{1}{\gamma_0}} \exp \left(\frac{1}{2} \sigma_0^2(2) \right), \quad (6.1)$$

where n_2 denotes the number of teams of size 2. As this quantity is an average over the worker fixed effects, it can be orthogonalized with respect to them using our approach.

Ahmadpoor and Jones (2019) consider model (5.3) without the error term ε_j . Here our goal is to address the statistical challenge caused by the presence of a large number of possibly imprecisely estimated fixed effects. An alternative would be to specify a distribution for author heterogeneity conditional on the team network (i.e., for all the η_{i0} 's conditional on \mathcal{K}), as in Bonhomme (2021). An advantage of such a procedure would be that, under correct specification, estimates are consistent even in poorly connected networks. This random-effect approach requires, however, to model how authors sort and collaborate in teams. Our approach avoids the need to do so. On the other hand, a fixed-effect approach requires that the author effects can be consistently estimated. In less well-connected networks, the convergence rate will be slower. Orthogonalization to a higher-order allows us to reduce the impact of estimation noise.

We look at the production of academic work on economics. We use data from Ductor, Fafchamps, Goyal and Van der Leij (2014), drawn from the EconLit database. These data contain a large collection of articles, indicated by their ID, together with author identifiers and a measure of journal quality proposed by Kodrzycki and Yu (2006). This measure is a ranking between 0 and 100, which we net of multiplicative time effects and will use as our outcome variable. We restrict the sample to articles published between 1990 and 1999, written either alone or with a single co-author. We only include authors who produced at least two sole-authored articles during the sampling period.

Our sample contains 91,626 articles, 10% of which are co-authored, and 16,408 authors. Average journal quality differs greatly across authors, with the 10th percentile of the quality measure being 0.4, the median being 0.9, and the 90th percentile being 8.5. The between-author variance in journal quality is 42% of the overall variance. The distribution of journal quality, in turn, is skewed to the right, with a median of 0.6, a 90th percentile of 12, and a 99th percentile of 52. The number of publications per author varies substantially, with a 10th percentile of 2, a median of 4, and a 90th percentile of 13.

To implement our approach, we construct subsets of three papers, one co-authored (j) and two sole-authored ($j_1(j), j_2(j)$), as described in (5.4)–(5.6). The score for θ based on subset i then involves the three teams $j, j_1(j)$, and $j_2(j)$. Proceeding in this way is helpful

as it limits the dimension of the parameter η_i to two. This is not only in line with the assumptions we make in deriving asymptotics, but also helpful in terms of computation. Moreover, it reduces the number of derivatives that need to be computed. The number of derivatives nevertheless remains substantial, as we need to compute 9 derivatives at order 2, 19 at order 3, and 55 at order 5, for example. Yet, using the computational remarks from Section 5.2, this can be implemented quite fast.

Finally, we exploit the network structure of the data to perform our sample splitting. For every worker, we construct a preliminary estimator of her fixed effect (in logs) as the average quality of her single-authored papers, except for one that we select at random and use later in estimation. This strategy is feasible due to our sample restriction. For each subset i of three teams, we then stack the two worker fixed effects together to form our preliminary estimate $\hat{\eta}_i$. We next estimate the parameters $\beta_0(2), \gamma_0(2), \sigma_0^2(1), \sigma_0^2(2)$ on the sample from which all these single-authored articles have been removed. In the present case, \tilde{u}_i in (3.4) has four components that correspond to the score with respect to all the parameters, and the weight matrix \widetilde{W} is irrelevant since the problem is just-identified. In order to limit the variability due to the choice of split, we average parameter estimates across 100 random splits, through cross-fitting. The bias in the parameter estimates takes a complex form due to the team network environment. In Appendix D we assess the ability of our orthogonalization approach to alleviate this bias in a Monte Carlo simulation.

6.2 Empirical estimates

Table 1 shows the estimates of $\beta_0, \gamma_0, \sigma_0^2(2)$, and $\sigma_0^2(1)$ for various estimators. These are the plug-in estimator based on the preliminary estimates $\hat{\eta}_i$ and six estimators based on Neyman-orthogonalized moments, for $1 \leq q \leq 6$. In addition to point estimates, we report estimated standard errors based on the parametric bootstrap.⁴

Starting with the substitution parameter γ , the uncorrected estimate is 0.12, which is

⁴Bootstrap replications are based on Neyman-orthogonalized estimates of $\beta_0, \gamma_0, \sigma_0^2(2)$, and $\sigma_0^2(1)$ to order $q = 6$, together with the sample-split estimates $\hat{\eta}_i$ of author effects. Within each bootstrap replication, we cross-fit the estimates 10 times. Results are based on 200 bootstrap replications.

Table 1: Estimation results

	Substitution γ	Team size β	Variance $\sigma^2(2)$	Variance $\sigma^2(1)$
Plug-in	0.1233 (0.0466)	1.2890 (0.0217)	1.6230 (0.0249)	1.6404 (0.0204)
$q = 1$	-1.9979 (0.2123)	1.3344 (0.0279)	1.6797 (0.0264)	1.7379 (0.0265)
$q = 2$	0.7268 (0.2367)	1.3046 (0.0359)	1.4341 (0.0260)	1.4679 (0.0260)
$q = 3$	0.4467 (0.2034)	1.2936 (0.0369)	1.4399 (0.0254)	1.4423 (0.0236)
$q = 4$	0.3976 (0.1763)	1.2931 (0.0362)	1.4346 (0.0254)	1.4238 (0.0232)
$q = 5$	0.3947 (0.1730)	1.2930 (0.0361)	1.4328 (0.0254)	1.4209 (0.0231)
$q = 6$	0.3944 (0.1770)	1.2930 (0.0365)	1.4316 (0.0254)	1.4194 (0.0230)

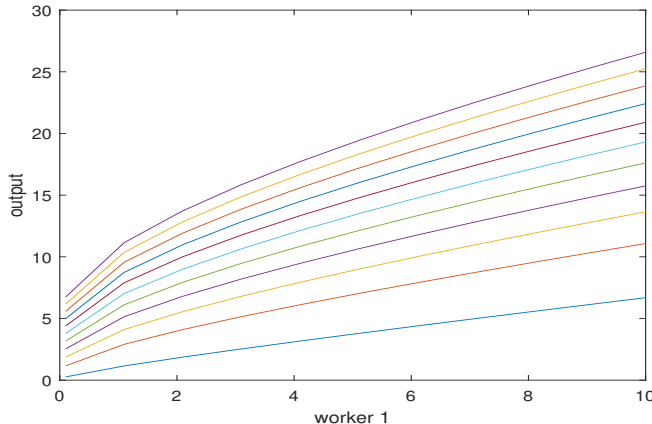
Notes: Point estimates based on q -ordered orthogonalized estimators, cross-fitted estimates (100 splits).

Parametric bootstrap standard errors in parentheses (200 replications).

close to the Cobb-Douglas case. The value of the first-order Neyman-orthogonalized estimate is quite different. However, since the preliminary estimates of the author fixed effects are based on very few observations, we do not expect this estimator to adequately correct for bias. This is confirmed by the fact that all other Neyman-orthogonalized estimates, for $q \in \{2, \dots, 6\}$, range between 0.39 and 0.73, which is higher than the plug-in estimate, and very different from the first-order orthogonalized estimate. Relative to the plug-in, the orthogonalized estimates with $q \geq 2$ all indicate somewhat less complementarity between authors in team production. Notice the stability of estimates for larger values of q . A substitution parameter $\gamma = 0.4$ corresponds to the case of imperfect complements; see Figure 1 for a graphical illustration.

Turning to the other parameters, the estimates of the team size parameter β are virtually unaffected by the orthogonalization. This suggests the bias is limited for this parameter. Its value is close to 1.3, implying that producing a paper with a co-author increases the paper's quality to some extent. Next, the log-error variance $\sigma^2(2)$ in teams of two coauthors

Figure 1: Production function estimate



Notes: Worker 1's type η_1 on the x-axis, average output Y_j on the y-axis. Each curve corresponds to a different worker 2's type η_2 . Figure based on the point estimates for $q = 6$ reported in Table 1.

is larger when using plug-in estimates (1.6) than when using orthogonalization with $q \geq 2$ (1.4), suggesting that the plug-in and first-order corrected estimates are biased upward. Lastly, the variance $\sigma^2(1)$ in teams of a single author is also larger under the plug-in estimator.

Model (5.3) implies some restrictions on the parameters $\gamma, \beta, \sigma^2(1), \sigma^2(2)$ that do not depend on the author-specific effects η_i . In Appendix E we exploit two types of restrictions as robustness checks. Our findings suggest that, while those restrictions seem broadly consistent with the higher-order orthogonalized estimates reported in Table 1, using them directly for estimation may lead to very imprecise estimates.

Lastly, we report estimates of average journal quality in a counterfactual scenario where authors are randomly assigned across teams of two co-authors, see (6.1). The first column in Table 2 shows estimates of the average output in the empirical allocation. This quantity can be estimated without bias as the sample mean of the journal quality variable, which is equal to 7.0. We see that the plug-in estimate is 8.4, larger than the empirical value. In comparison, Neyman-orthogonalized estimates for $q \geq 3$ range between 6.2 and 7.4, and estimates for $q = 5$ and $q = 6$ are closest to the empirical value. The second column in

Table 2: Empirical estimates: average output

	Average output	Counterfactual
Plug-in	8.4374 (0.2492)	7.1965 (0.1949)
$q = 1$	6.8469 (0.4837)	5.3147 (0.3936)
$q = 2$	9.0756 (0.5096)	8.5805 (0.7626)
$q = 3$	6.1566 (0.4452)	5.4774 (0.4354)
$q = 4$	7.4414 (0.4726)	6.6378 (0.6227)
$q = 5$	6.9854 (0.3889)	6.1694 (0.3896)
$q = 6$	7.1255 (0.3194)	6.3318 (0.4225)

Notes: Average output (the value in the data is 6.9995, standard error 0.2309), and counterfactual average output in a random allocation. Point estimates based on orthogonalized estimators to order q , cross-fitted estimates (100 splits). Parametric bootstrap standard errors in parentheses (200 replications).

Table 2 shows estimates of average article quality under random assignment of authors to teams, using the plug-in method and Neyman-orthogonalized estimates to order $q \geq 1$.⁵ The estimates vary with the order of orthogonalization. When taking $q \geq 4$, estimates range between 6.2 and 6.6. In addition, comparing the two columns of Table 2 shows that, irrespective of the order of orthogonalization, the estimates of average output are lower in the counterfactual scenario where workers are randomly allocated across teams.

The main takeaway from Table 2 is that randomly allocating authors among teams would tend to lower average paper quality. This is due to two economic forces. The first one is complementarity in production, as reflected by estimates of γ lower than 1. The second force is positive sorting. Indeed, the preliminary estimates of worker fixed effects are positively correlated within teams in the data. In the presence of complementarity,

⁵To speed up computation, we approximate (6.1) using a random subset of 1000 authors, for each random sample split (and each bootstrap replication).

decreasing assortative matching leads to lower output, which is what we find in Table 2.

7 Asymptotic properties

In this section, we show that, under higher-order orthogonality, the estimators $\hat{\theta}$ and $\hat{\mu}$ introduced in Section 3.2 are \sqrt{n} -consistent and asymptotically normal under appropriate assumptions, even if the convergence rate of $\hat{\eta}_i$ is slower than \sqrt{n} . We focus on deriving the asymptotic distribution of $\hat{\mu}$, assuming that we have already worked out the corresponding asymptotic result of $\hat{\theta}$. However, the corresponding theory for $\hat{\theta}$ is actually a special case of our results for $\hat{\mu}$, where θ is dropped from the arguments, μ is replaced by θ , and u_i are replaced by \tilde{u}_i . Thus, our focus on $\hat{\mu}$ is without loss of generality.

7.1 Notation

For the presentation of the asymptotic theory, it is useful to be explicit about which parameters depend on the sample size and which ones do not. Recall that n is the total number of observations in (Z_1, \dots, Z_N) , where each Z_i comprises n_i observations. In the asymptotic sequence, we let N and n_i depend on n , although we do not explicitly indicate this dependence. For example, in a panel data model, our assumptions allow both N and T to grow as the number NT of observations tends to infinity.

To indicate the dependence on the sample size, we will write η_n and μ_n instead of η and μ in this section. While the dimension of μ is not changing with n , the true parameter $\mu_{0,n}$ is implicitly defined as the solution of $\sum_{i=1}^N \mathbb{E}_{\theta_0, \eta_{0,n}} (u_i(Z_i; \theta_0, \eta_{0,n,i}, \mu)) = 0$, which may depend on n as well. By contrast, the parameter θ and its true value θ_0 are independent of n .

Remember also that $n = \sum_{i=1}^N n_i$, and note that if the observations within each unit i are independent, then we have $\ell(y_i | x_i; \theta, \eta_{n,i}) = \prod_{j=1}^{n_i} \ell(y_{ij} | x_{ij}; \theta, \eta_{n,i})$. Hence, if

$u_i(Z_i; \theta, \eta_{n,i}) = \frac{\partial \log \ell(y_i | x_i; \theta, \eta_{n,i})}{\partial \theta}$ and $u_{ij}(Z_{ij}; \theta, \eta_{n,i}) = \frac{\partial \log \ell(y_{ij} | x_{ij}; \theta, \eta_{n,i})}{\partial \theta}$, then

$$u_i(Z_i; \theta, \eta_{n,i}) = \sum_{j=1}^{n_i} u_{ij}(Z_{ij}; \theta, \eta_{n,i}).$$

More generally, whenever $n_i \rightarrow \infty$ we expect that u_i scales linearly with n_i , implying that $\frac{1}{n} \sum_{i=1}^N u_i$ is the correctly-scaled sample average of u_i , and also explaining the scaling of various other terms in Assumption 1 below.

7.2 A useful lemma

With this notation in hand, we now state our first assumption.

Assumption 1.

- (i) We have $\left[\frac{1}{n} \sum_{i=1}^N \frac{\partial u_i^\top(Z_i; \hat{\theta}, \hat{\eta}_{n,i}, \hat{\mu}_n)}{\partial \mu} \right] W \left[\frac{1}{\sqrt{n}} \sum_{i=1}^N u_i(Z_i; \hat{\theta}, \hat{\eta}_{n,i}, \hat{\mu}_n) \right] = o_P(1)$, for some non-random symmetric positive definite weight matrix W .
- (ii) As $n \rightarrow \infty$, $(\hat{\theta}, \hat{\eta}_n, \hat{\mu}_n)$ is contained in some convex neighborhood \mathcal{B}_n of $(\theta_0, \eta_{0,n}, \mu_{0,n})$. Let $\mathcal{B}_{n,i}$ be the convex neighborhood of $(\theta_0, \eta_{0,n,i}, \mu_{0,n})$ obtained by intersecting \mathcal{B}_n with the parameter parameter subspace for observation i .
- (iii) $\max_i \dim(\eta_{n,i}) = O(1)$.
- (iv) For every i , the function $u_i(Z_i, \theta, \eta_{n,i}, \mu)$ is $(q+1)$ times continuously differentiable in the parameters $(\theta, \eta_{n,i}, \mu)$, and we assume that for all its components all the partial derivatives of $u_i(Z_i; \theta, \eta_{n,i}, \mu)$ up to order $(q+1)$ are bounded in absolute value by $n_i C_{n,i}(Z_i) \geq 0$, uniformly in the neighborhood $\mathcal{B}_{n,i}$, such that $\frac{1}{n} \sum_{i=1}^N n_i \mathbb{E} [C_{n,i}(Z_i)^2] = O(1)$.
- (v) $\hat{\mu}_n - \mu_{0,n} = o_P(1)$ and $\frac{1}{n} \sum_{i=1}^N n_i \mathbb{E} \left(\|\hat{\eta}_{n,i} - \eta_{0,n,i}\|^{2(q+1)} \right) = o(n^{-1})$.
- (vi) $\hat{\theta} = \theta_0 + \frac{1}{n} \sum_{i=1}^N \psi_{n,i} + o_P(n^{-1/2})$, where $\mathbb{E}(\psi_{n,i}) = 0$ and $\frac{1}{n} \sum_{i=1}^N \mathbb{E} (\|\psi_{n,i}\|^2) = O(1)$.

(vii) *The probability limits*

$$G_\mu = \text{plim}_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^N \frac{\partial u_i(Z_i; \theta_0, \eta_{0,n,i}, \mu_{0,n})}{\partial \mu^\top}, \quad G_\theta = \text{plim}_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^N \frac{\partial u_i(Z_i; \theta_0, \eta_{0,n,i}, \mu_{0,n})}{\partial \theta^\top}$$

exist, and $\text{rank}(G_\mu) = \dim(\mu)$.

Part (i) in Assumption 1 is satisfied if $\widehat{\mu}_n$ is computed using GMM, see (3.3). In Part (ii), the neighborhood \mathcal{B}_n depends on the sample size n , partly because the number of nuisance parameters of $\eta_{n,i}$ generally depends on n . Part (iii) assumes that the maximal dimension of $\eta_{n,i}$ is bounded as $n \rightarrow \infty$. Part (iv) requires the derivatives of the moment functions (properly rescaled) to be suitably bounded. The first half of Part (v) is a high-level consistency assumption for $\widehat{\mu}_n$, which can be justified by guaranteeing that the objective function in (3.3) converges uniformly to a population counterpart that has a unique minimum at μ_0 . The second half of Part (v) is the rate requirement on the preliminary estimates $\widehat{\eta}_{n,i}$, imposing a rate faster than $n^{-1/2(q+1)}$. Part (vi) requires $\widehat{\theta}$ to be asymptotically linear, in particular requiring $\widehat{\theta} - \theta_0 = O_P(n^{-1/2})$. In the case where $\mu_{n,0} = \theta_0$ this condition is not needed. Lastly, Part (vii) assumes existence of Jacobian matrices and a rank condition.

In the statement of the following lemma, $D_{\eta_{n,i}}^m$ denote the derivative operator with respect to $\eta_{n,i}$.

Lemma 1. *Under Assumption 1 we have*

$$\begin{aligned} & \sqrt{n}(\widehat{\mu}_n - \mu_{0,n}) \\ &= - (G_\mu^\top W G_\mu)^{-1} G_\mu^\top W \left\{ \frac{1}{\sqrt{n}} \sum_{i=1}^N \left[u_i(Z_i; \theta_0, \eta_{0,n,i}, \mu_{0,n}) + G_\theta \psi_{n,i} \right] + R_n \right\} + o_P(1), \end{aligned}$$

where

$$R_n = \frac{1}{\sqrt{n}} \sum_{i=1}^N \sum_{m \in \mathbb{K}_{q,n,i}} \frac{1}{m!} \left[D_{\eta_{n,i}}^m u_i(Z_i; \theta_0, \eta_{0,n,i}, \mu_{0,n}) \right] (\widehat{\eta}_{n,i} - \eta_{0,n,i})^m,$$

and $\mathbb{K}_{q,n,i} = \left\{ m \in \mathbb{Z}^{\dim(\eta_{n,i})} : 1 \leq \sum_{r=1}^{\dim(\eta_{n,i})} m_r \leq q \right\}$.

7.3 Main result

We are now in position to establish the main result of this section, which concerns root- n consistency and asymptotic normality of estimators based on orthogonal equations. For this, we first state our second assumption.

Assumption 2.

- (i) The moment functions $u_i(Z_i; \theta, \eta_{n,i}, \mu)$ are Neyman-orthogonal to order q for all i , and $\sum_{i=1}^N \mathbb{E}(u_i(Z_i; \theta_0, \eta_{0,n,i}, \mu_{0,n})) = 0$.
- (ii) $\hat{\eta}_{n,i}$ are independent of (Z_1, \dots, Z_N) for all i .
- (iii) The Z_1, \dots, Z_N are independent across i .
- (iv) $\xi_{n,i} = u_i(Z_i; \theta_0, \eta_{0,n,i}, \mu_{0,n}) + G_\theta \psi_{n,i}$ satisfies Lindeberg's condition,⁶ and the following probability limit exists:

$$V_\xi = \text{plim}_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^N \text{Var}(\xi_{n,i}).$$

Part (i) in Assumption 2 requires u_i to be Neyman-orthogonal in the sense of Definition 1. Part (ii) requires the preliminary estimates to be independent from the estimation sample. With independent observations, this can be achieved by sample splitting. Part (iii) imposes independence between the Z_i 's. We impose this assumption to simplify the presentation. It is straightforward to modify the variance expression in Theorem 2 below to account for particular forms of dependence (e.g., clustered) by using an appropriate expression for the matrix V_ξ introduced in Part (iv).

The following theorem provides an asymptotic characterization of $\hat{\mu}_n$.

Theorem 2. *Let Assumptions 1 and 2 hold with the same value of $q \in \{1, 2, 3, \dots\}$. Then we have*

$$\sqrt{n} (\hat{\mu}_n - \mu_{0,n}) \xrightarrow{d} \mathcal{N}(0, (G_\mu^\top W G_\mu)^{-1} G_\mu^\top W V_\xi W G_\mu (G_\mu^\top W G_\mu)^{-1}).$$

⁶That is, for any $\epsilon > 0$, $\frac{1}{s_n^2} \sum_{i=1}^N \mathbb{E}[\xi_{n,i}^2 \cdot \mathbb{1}(|\xi_{n,i}| > \epsilon s_n)] \rightarrow 0$ as $n \rightarrow \infty$, where $s_n^2 = \sum_{i=1}^N \text{Var}(\xi_{n,i})$ and $\mathbb{1}$ is the indicator function.

8 Final remarks

In this paper we show how to construct higher-order Neyman-orthogonal moment functions in conditional-likelihood models. We use these functions, together with sample splitting, to reduce bias in estimation. Our application suggests that our higher-order corrections can be effective in network settings with fixed effects. There are several important avenues for future work. An area of application is to double/debiased machine learning with fixed effects, where the nuisance parameters contains some components, such as low-dimensional functions, for which first-order orthogonality may suffice. An open question is how to choose the degree q of orthogonality in practice. Finally, extending the approach to non-likelihood models is important, and we are working on a strategy that relies on independence assumptions and sample splitting.

References

- Abowd, J. M., F. Kramarz, and D. N. Margolis (1999). High wage workers and high wage firms. *Econometrica* 67, 251–333.
- Ahmadpoor, M. and B. F. Jones (2019). Decoding team and individual impact in science and invention. *Proceedings of the National Academy of Sciences* 116, 13885–13890.
- Andrews, M. J., L. Gill, T. Schank, and R. Upward (2008). High wage workers and low wage firms: negative assortative matching or limited mobility bias? *Journal of the Royal Statistical Society: Series A* 171, 673–697.
- Angrist, J. D. and B. Frandsen (2022). Machine labor. *Journal of Labor Economics* 40, S97–S140.
- Arellano, M. (2003). Discrete choices with panel data. *Investigaciones Economicas XXVII*, 423–458.
- Arellano, M. and S. Bonhomme (2012). Identifying distributional characteristics in random coefficients panel data models. *The Review of Economic Studies* 79, 987–1020.
- Arellano, M. and J. Hahn (2007). Understanding bias in nonlinear panel models: Some recent developments. In R. Blundell, W. K. Newey, and T. Persson (Eds.), *Advances In Economics and Econometrics*, Volume III. Econometric Society: Cambridge University Press.
- Bhattacharyya, A. (1946). On some analogues of the amount of information and their use in statistical estimation. *Sankhyā* 8, 1–14.
- Bickel, P. (1982). On adaptive estimation. *Annals of Statistics* 10, 647–671.
- Bonhomme, S. (2021). Teams: Heterogeneity, sorting, and complementarity. Mimeo.
- Chernozhukov, V., D. Chetverikov, M. Demirer, E. Duflo, C. Hansen, W. Newey, and J. Robins (2018). Double/debiased machine learning for treatment and structural parameters. *Econometrics Journal* 21, C1–C68.
- Chernozhukov, V., J. C. Escanciano, H. Ichimura, W. K. Newey, and J. M. Robins (2022). Locally robust semiparametric estimation. *Econometrica* 90, 1501–1535.
- Constantine, G. and T. Savits (1996). A multivariate Faa di Bruno formula with applica-

- tions. *Transactions of the American Mathematical Society* 348, 503–520.
- Dhaene, G. and K. Jochmans (2015a). Profile-score adjustments for incidental-parameter problems. Mimeo.
- Dhaene, G. and K. Jochmans (2015b). Split-panel jackknife estimation of fixed-effect models. *The Review of Economic Studies* 82, 991–1030.
- Dhaene, G. and K. Jochmans (2016). Likelihood inference in an autoregression with fixed effects. *Econometric Theory* 32, 1178–1215.
- Ductor, L., M. Fafchamps, S. Goyal, and M. J. Van der Leij (2014). Social networks and research output. *Review of Economics and Statistics* 96, 936–948.
- Ghazal, G. A. and H. Neudecker (2000). On second-order and fourth-order moments of jointly distributed random matrices: a survey. *Linear Algebra and its Applications* 321, 61–93.
- Graham, B. S. (2017). An econometric model of network formation with degree heterogeneity. *Econometrica* 85, 1033–1063.
- Graham, B. S. (2020). Sparse network asymptotics for logistic regression. *arXiv preprint arXiv:2010.04703*.
- Graham, B. S., G. W. Imbens, and G. Ridder (2014). Complementarity and aggregate implications of assortative matching: A nonparametric analysis. *Quantitative Economics* 5, 29–66.
- Hahn, J. and J. Hausman (2021). Problems with the control variable approach in achieving unbiased estimates in nonlinear models in the presence of many instruments. *Journal of Quantitative Economics* 19, 39–58.
- Hahn, J. and G. Kuersteiner (2011). Bias reduction for dynamic nonlinear panel models with fixed effects. *Econometric Theory* 27, 1152–1191.
- Hahn, J. and W. K. Newey (2004). Jackknife and analytical bias reduction for nonlinear panel models. *Econometrica* 72, 1295–1319.
- Ichimura, H. (1993). Semiparametric least squares (SLS) and weighted SLS estimation of single-index models. *Journal of Econometrics* 58, 71–120.
- Jackson, C. K., J. E. Rockoff, and D. O. Staiger (2014). Teacher effects and teacher related

- policies. *Annual Review of Economics* 6, 801–825.
- Jochmans, K. and M. Weidner (2019). Fixed-effect regressions on network data. *Econometrica* 87, 1543–1560.
- Kline, P., R. Saggio, and M. Sølvssten (2020). Leave-out estimation of variance components. *Econometrica* 88, 1859–1898.
- Kodrzycki, Y. K. and P. Yu (2006). New approaches to ranking economics journals. *The BE Journal of Economic Analysis & Policy* 5.
- Lancaster, T. (2002). Orthogonal parameters and panel data. *Review of Economic Studies* 69, 647–666.
- Li, H., B. Lindsay, and R. Waterman (2003). Efficiency of projected score methods in rectangular array asymptotics. *Journal of the Royal Statistical Society, Series B* 65, 191–208.
- Mackey, L., V. Syrgkanis, and I. Zadik (2018). Orthogonal machine learning: Power and limitations. In *International Conference on Machine Learning*, pp. 3375–3383. PMLR.
- Magnus, J. R. and H. Neudecker (1979). The commutation matrix: some properties and applications. *Annals of Statistics* 7, 381–394.
- Magnus, J. R. and H. Neudecker (1980). The elimination matrix: some lemmas and applications. *SIAM Journal on Algebraic Discrete Methods* 1(4), 422–449.
- McLeish, D. L. and C. G. Small (1994). *Hilbert Space Methods in Probability and Statistical Inference*. Wiley NY.
- Newey, W. K. (1994). The asymptotic variance of semiparametric estimators. *Econometrica* 62, 1349–1382.
- Newey, W. K. and J. M. Robins (2017). Cross-fitting and fast remainder rates for semiparametric estimation. Mimeo.
- Neyman, J. (1959). Optimal asymptotic tests of composite hypotheses. In U. Grenander (Ed.), *Probability and Statistics*, pp. 416–444. Wiley NY.
- Neyman, J. and E. L. Scott (1948). Consistent estimates based on partially consistent observations. *Econometrica* 16, 1–32.
- Robins, J., L. Li, E. Tchetgen Tchetgen, and A. van der Vaart (2008). Higher order influence

- functions and minimax estimation of nonlinear functionals. In *Probability and Statistics: Essays in Honor of David A. Freedman*, pp. 335–421. Institute of Mathematical Statistics.
- Robinson, P. M. (1988). Root- n -consistent semiparametric regression. *Econometrica* *56*, 931–954.
- Schick, A. (1986). On asymptotically efficient estimation in semiparametric models. *Annals of Statistics* *14*, 1139–1151.
- Semenova, V., M. Goldman, V. Chernozhukov, and M. Taddy (2023). Inference on heterogeneous treatment effects in high-dimensional dynamic panels under weak dependence. *Quantitative Economics* *14*, 471–510.
- Small, C. G. and D. L. McLeish (1988). Generalizations of ancillarity, completeness, and sufficiency in an inference function space. *Annals of Statistics* *16*, 534–551.
- Small, C. G. and D. L. McLeish (1989). Projection as a method for increasing sensitivity and eliminating nuisance parameters. *Biometrika* *76*, 693–703.
- Waterman, R. P. and B. G. Lindsay (1996). Projected score methods for approximating conditional scores. *Biometrika* *83*, 1–13.
- Woutersen, T. (2002). Robustness against incidental parameters. Technical report, Research Report.
- Wüthrich, K. and Y. Zhu (2021). Omitted variable bias of Lasso-based inference methods: A finite sample analysis. Forthcoming in *Review of Economics and Statistics*.

APPENDIX

A Proofs

A.1 Proof of Theorem 1

Before proving the theorem, it is useful to establish the following lemma.

Lemma 2. *Let $q \in \{1, 2, 3, \dots\}$, and let x be some realization of the covariates. Remember that ∇_η^q and $w_q(z; \theta, \eta)$ are vectors of dimension $k_q = \sum_{p=1}^q d_p$, and that $\Sigma_{w_q w_q}(x; \theta, \eta)$ is a $k_q \times k_q$ matrix. We assume that $\Sigma_{w_q w_q}(x; \theta, \eta)$ is invertible, and we define*

$$\tilde{w}_q(z; \theta, \eta) = \Sigma_{w_q w_q}(x; \theta, \eta)^{-1} w_q(z; \theta, \eta).$$

Then,

$$\mathbb{E}_{\theta, \eta} [(\nabla_\eta^q)^\top \tilde{w}_q(z; \theta, \eta) | X = x] = \begin{pmatrix} -\mathbb{I}_{d_1} & 0 & 0 & \cdots & 0 \\ 0 & +\mathbb{I}_{d_2} & 0 & \cdots & 0 \\ 0 & 0 & -\mathbb{I}_{d_3} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & (-1)^q \mathbb{I}_{d_q} \end{pmatrix}, \quad (\text{A.1})$$

where the diagonal $k_q \times k_q$ matrix on the right hand side is obtained by stacking $(-1)^p \mathbb{I}_{d_p}$ on the diagonal for $p = 1, \dots, q$, and \mathbb{I}_{d_p} is the identity matrix of dimensions d_p .

Proof of Lemma 2. By taking derivatives of $\int \ell(y|x; \theta, \eta) dy = 1$ with respect to η , we obtain

$$\int [\nabla_\eta^q \ell(y|x; \theta, \eta)] dy = 0,$$

which can also be written as $\mathbb{E}_{\theta, \eta}[w_q(Z; \theta, \eta) | X = x] = 0$. Since $\Sigma_{w_q w_q}(X; \theta, \eta)$ does not depend on Y we also have $\mathbb{E}_{\theta, \eta}[\tilde{w}_q(Z; \theta, \eta) | X = x] = 0$. Using this and the definition of \tilde{w}_q

we obtain

$$\begin{aligned}
\mathbb{I}_{k_q} &= \Sigma_{w_q w_q}(x; \theta, \eta)^{-1} \Sigma_{w_q w_q}(x; \theta, \eta) \\
&= \Sigma_{w_q w_q}(x; \theta, \eta)^{-1} \mathbb{E}_{\theta, \eta}[w_q(Z; \theta, \eta) w_q(Z; \theta, \eta)^\top | X = x] \\
&= \mathbb{E}_{\theta, \eta}[\tilde{w}_q(Z; \theta, \eta) w_q(Z; \theta, \eta)^\top | X = x] \\
&= \int \tilde{w}_q(z; \theta, \eta) [\nabla_\eta^q \ell(y|x; \theta, \eta)]^\top dy.
\end{aligned} \tag{A.2}$$

According its definition in Section 3.1, the elements of the k_q -vector operator ∇_η^q are given by

$$D_\eta^m = \frac{\partial^p}{\partial \eta_{m_1} \cdots \partial \eta_{m_p}},$$

and are uniquely labeled by vectors of integers $m = (m_1, \dots, m_p)$ of length $p \in \{1, \dots, q\}$ in the following set

$$\mathcal{C}_q = \bigcup_{p \in \{1, \dots, q\}} \{m = (m_1, \dots, m_p) : 1 \leq m_1 \leq \dots \leq m_p \leq d_\eta\}.$$

Analogously, we now introduce the notation $\tilde{w}_q^m(z; \theta, \eta)$, $m \in \mathcal{C}_q$, to uniquely denote the elements of the vector $\tilde{w}_q(z; \theta, \eta)$, which is also a vector of length $k_q = |\mathcal{C}_q|$. With that notation, the result in display (A.2) can equivalently be written as

$$\forall r, m \in \mathcal{C}_q : \int \tilde{w}_q^r(z; \theta, \eta) [D_\eta^m \ell(y|x; \theta, \eta)] dy = \mathbb{1}\{r = m\}. \tag{A.3}$$

Since $\mathbb{E}_{\theta, \eta}[\tilde{w}_q(Z; \theta, \eta) | X = x] = 0$, we also have

$$\int \tilde{w}_q^r(z; \theta, \eta) \ell(y|x; \theta, \eta) dy = 0.$$

For the empty vector $()$ of length zero we have $D_\eta^{\emptyset} \ell(y|x; \theta, \eta) = \ell(y|x; \theta, \eta)$. Using this notation we can combine the result in the last two displays to find that for all $r \in \mathcal{C}_q$ and all $v \in \mathcal{C}_q \cup \{()\}$ we have

$$\int \tilde{w}_q^r(z; \theta, \eta) [D_\eta^v \ell(y|x; \theta, \eta)] dy = \mathbb{1}\{r = v\}.$$

For $k \in \{1, 2, 3, \dots\}$, vector $t = (t_1, \dots, t_k) \in \mathcal{C}_q$, and a subset $S \subseteq \{1, \dots, k\}$, let t_S denote the vector formed by keeping only the indices in S , and let $t_{-S} = t_{\{1, \dots, k\} \setminus S}$ be the vector of

the remaining elements. Then, by applying D_η^t to the last display and using the product rule for differentiation we obtain

$$\sum_{S \subseteq \{1, \dots, k\}} \int [D_\eta^{t_S} \tilde{w}_q^r(z; \theta, \eta)] [D_\eta^{t-S} D_\eta^v \ell(y|x; \theta, \eta)] dy = 0. \quad (\text{A.4})$$

Of course, we have $D_\eta^t D_\eta^v = D_\eta^{(t,v)}$, and instead of distinguishing between t and v we can also just write m for (t, v) combined. The last display equation then implies that for any nonempty subset $T \subseteq \{1, 2, \dots, |m|\}$ we have (just set $t = m_T$ and $v = m_{-T}$ in the last display result—in doing so, it was important that above we allowed for v to be the empty vector):

$$\sum_{S \subseteq T} \int [D_\eta^{m_S} \tilde{w}_q^r(z; \theta, \eta)] [D_\eta^{m-S} \ell(y|x; \theta, \eta)] dy = 0.$$

Now, consider the following linear combination of the result in the last display,

$$\sum_{\substack{T \subseteq \{1, 2, \dots, |m|\} \\ T \neq \emptyset}} (-1)^{|T|} \sum_{S \subseteq T} \int [D_\eta^{m_S} \tilde{w}_q^r(z; \theta, \eta)] [D_\eta^{m-S} \ell(y|x; \theta, \eta)] dy = 0. \quad (\text{A.5})$$

For a fixed $S \subset \{1, 2, \dots, |m|\}$, the total coefficient of the term $\int [D_\eta^{m_S} \tilde{w}_q^r] [D_\eta^{m-S} \ell] dy$ in this linear combination is given by

$$\sum_{\substack{T: S \subseteq T \subseteq \{1, \dots, |m|\} \\ T \neq \emptyset}} (-1)^{|T|} = \begin{cases} -1 & \text{if } S = \emptyset, \\ (-1)^{|m|} & \text{if } S = \{1, \dots, |m|\}, \\ 0 & \text{otherwise.} \end{cases} \quad (\text{A.6})$$

To see that the result in the last display holds, notice first that for $S = \emptyset$ the sum is simply

$$\sum_{T: \emptyset \subseteq T \subseteq \{1, \dots, |m|\}} (-1)^{|T|} = -1 + \underbrace{\sum_{T \subseteq \{1, \dots, |m|\}} (-1)^{|T|}}_{=0} = -1,$$

where $\sum_{T \subseteq \{1, \dots, |m|\}} (-1)^{|T|} = \sum_{r=0}^{|m|} \binom{|m|}{r} (-1)^r = (1-1)^{|m|} = 0$ is a classic alternating sum result, which holds for all $|m| \geq 1$. Next, for $S = \{1, \dots, |m|\}$, the left hand side of (A.6) only sums over one element, $T = \{1, \dots, |m|\}$, and we thus get $(-1)^{|T|} = (-1)^{|m|}$ for the

sum. Finally, if $S \neq \emptyset$ and $S \neq \{1, \dots, |m|\}$, then the left hand side of (A.6) can be written as

$$\sum_{T: S \subseteq T \subseteq \{1, \dots, |m|\}} (-1)^{|T|} = \sum_{R \subseteq -S} (-1)^{|S|+|R|} = (-1)^{|S|} \underbrace{\sum_{R \subseteq -S} (-1)^{|R|}}_{=0} = 0$$

where we replaced the sum over T by a sum over R such that $T = S \cup R$, with $-S = \{1, \dots, |m|\} \setminus S$, and in the final step we used the alternating sum result again.

Using (A.6), our linear combination in (A.5) equals

$$- \int [D_\eta^m \tilde{w}_q^r(z; \theta, \eta)] \ell(y|x; \theta, \eta) dy + (-1)^{|m|} \int \tilde{w}_q^r(z; \theta, \eta) [D_\eta^m \ell(y|x; \theta, \eta)] dy = 0.$$

Together with (A.3) we thus find

$$\int [D_\eta^m \tilde{w}_q^r(z; \theta, \eta)] \ell(y|x; \theta, \eta) dy = (-1)^{|m|} \mathbb{1}\{r = m\}.$$

which in vector notation can be written as (A.1). \square

Proof of Theorem 1. Define $c_q(x; \theta, \eta, \mu) = [\Sigma_{w_q u}(x; \theta, \eta, \mu) - b_q(x; \theta, \eta, \mu)]^\top$. In the proof of Lemma 2 we introduced the notation $\tilde{w}_q^m(z; \theta, \eta)$, $m \in \mathcal{C}_q$, for the elements of the k_q -vector $\tilde{w}_q(z; \theta, \eta)$. Analogously, we now use $c_q^m(x; \theta, \eta, \mu)$ to denote the columns of the the $(\dim u) \times k_q$ -matrix $c(x; \theta, \eta, \mu)$, that is, $c_q^m(x; \theta, \eta, \mu)$ is a $(\dim u)$ -vector for every $m \in \mathcal{C}_q$. We have

$$\begin{aligned} c_q^m(x; \theta, \eta, \mu) &= \mathbb{E}_{\theta, \eta} \left[\frac{D_\eta^m \ell(y|x; \theta, \eta)}{\ell(y|x; \theta, \eta)} u(Z; \theta, \eta, \mu) \middle| X = x \right] - D_\eta^m \mathbb{E}_{\theta, \eta}(u(Z; \theta, \eta, \mu) | X = x) \\ &= \int [D_\eta^m \ell(y|x; \theta, \eta)] u(z; \theta, \eta, \mu) dy - D_\eta^m \int \ell(y|x; \theta, \eta) u(z; \theta, \eta, \mu) dy \\ &= - \sum_{S \subseteq \{1, \dots, |m|\}} \int [D_\eta^{m_S} \ell(y|x; \theta, \eta)] [D_\eta^{m-S} u(z; \theta, \eta, \mu)] dy, \end{aligned} \quad (\text{A.7})$$

where $z = (y, x)$ and in the last step we applied the product rule for differentiation as in (A.4) above, but the term for $S = \{1, \dots, |m|\}$ cancels with the term that stems from $\Sigma_{w_q u}(x; \theta, \eta, \mu)$, which explains why we only sum over subsets S that are different from $\{1, \dots, |m|\}$.

Next, by the definition of u_q^* and $A(x; \theta, \eta, \mu)$, we have

$$\begin{aligned} u_q^*(z; \theta, \eta, \mu) &= u(z; \theta, \eta, \mu) - A(x; \theta, \eta, \mu)^\top w_q(z; \theta, \eta) \\ &= u(z; \theta, \eta, \mu) - [\Sigma_{w_q u}(x; \theta, \eta, \mu) - b_q(x; \theta, \eta, \mu)]^\top \tilde{w}_q(z; \theta, \eta) \\ &= u(z; \theta, \eta, \mu) - \sum_{v \in \mathcal{C}_q} c_q^v(x; \theta, \eta, \mu) \tilde{w}_q^v(z; \theta, \eta). \end{aligned}$$

Let $m \in \mathcal{C}_q$. Applying the operator D_η^m to the last equation and again using the product rule for differentiation in the same way as before, we find

$$D_\eta^m u_q^*(z; \theta, \eta, \mu) = D_\eta^m u(z; \theta, \eta, \mu) - \sum_{S \subseteq \{1, \dots, |m|\}} \sum_{v \in \mathcal{C}_q} [D_\eta^{m-S} c_q^v(x; \theta, \eta, \mu)] [D_\eta^{m_S} \tilde{w}_q^v(z; \theta, \eta)].$$

Applying the conditional expectation operator to this and using Lemma 2 we obtain

$$\begin{aligned} &\mathbb{E}_{\theta, \eta} [D_\eta^m u_q^*(Z; \theta, \eta, \mu) | X = x] \\ &= \mathbb{E}_{\theta, \eta} [D_\eta^m u_q(Z; \theta, \eta, \mu) | X = x] \\ &\quad - \sum_{S \subseteq \{1, \dots, |m|\}} \sum_{v \in \mathcal{C}_q} [D_\eta^{m-S} c_q^v(x; \theta, \eta, \mu)] \underbrace{\mathbb{E}_{\theta, \eta} [D_\eta^{m_S} \tilde{w}_q^v(Z; \theta, \eta) | X = x]}_{=(-1)^{|S|} \mathbb{1}\{m_S=v\}} \\ &= \mathbb{E}_{\theta, \eta} [D_\eta^m u_q(Z; \theta, \eta, \mu) | X = x] - \sum_{\emptyset \neq S \subseteq \{1, \dots, |m|\}} (-1)^{|S|} [D_\eta^{m-S} c_q^{m_S}(x; \theta, \eta, \mu)]. \end{aligned}$$

where in the last step we used that for $S = \emptyset$ the indicator $\mathbb{1}\{m_S = v\}$ is always zero (because $v \in \mathcal{C}_q$ never has length zero), but for $S \neq \emptyset$ there is always exactly one $v \in \mathcal{C}_q$ that satisfies $m_S = v$, that is, in that second case we just remove the sum over v and replace v by m_S throughout. Next, plugging in the expression for $c_q^m(x; \theta, \eta, \mu)$ in equation (A.7) above we find

$$\begin{aligned} &\mathbb{E}_{\theta, \eta} [D_\eta^m u_q^*(Z; \theta, \eta, \mu) | X = x] \\ &= \mathbb{E}_{\theta, \eta} [D_\eta^m u(Z; \theta, \eta, \mu) | X = x] \\ &\quad + \sum_{\emptyset \neq S \subseteq \{1, \dots, |m|\}} (-1)^{|S|} D_\eta^{m-S} \left(\sum_{T \subsetneq S} \int [D_\eta^{m_T} \ell(y | x; \theta, \eta)] [D_\eta^{m_{S \setminus T}} u(z; \theta, \eta, \mu)] dy \right). \end{aligned}$$

By again using the product rule for differentiation to apply D_η^{m-S} to the product in the

last term, we obtain

$$\mathbb{E}_{\theta, \eta} [D_{\eta}^m u_q^*(Z; \theta, \eta, \mu) | X = x] \tag{A.8}$$

$$\begin{aligned} &= \int \ell(y | x; \theta, \eta) [D_{\eta}^m u(Z; \theta, \eta, \mu)] dy \\ &+ \sum_{\emptyset \neq S \subseteq \{1, \dots, |m|\}} (-1)^{|S|} \sum_{T \subsetneq S} \sum_{R \subseteq -S} \int [D_{\eta}^{m_{R \cup T}} \ell(y | x; \theta, \eta)] [D_{\eta}^{m_{(-S \setminus R) \cup (S \setminus T)}} u(z; \theta, \eta, \mu)] dy, \end{aligned} \tag{A.9}$$

where we write $-S$ for the set $\{1, \dots, |m|\} \setminus S$. All the terms on the right hand side of the last display equation are of the form $\int [D_{\eta}^{m_A} \ell(y | x; \theta, \eta)] [D_{\eta}^{m_{-A}} u(z; \theta, \eta, \mu)] dy$, for some $A \subseteq \{1, \dots, |m|\}$, and we can therefore write

$$\mathbb{E}_{\theta, \eta} [D_{\eta}^m u_q^*(Z; \theta, \eta, \mu) | X = x] = \sum_{A \subseteq \{1, \dots, |m|\}} \kappa_A \int [D_{\eta}^{m_A} \ell(y | x; \theta, \eta)] [D_{\eta}^{m_{-A}} u(z; \theta, \eta, \mu)] dy, \tag{A.10}$$

with

$$\kappa_A = \mathbb{1}\{A = \emptyset\} + \sum_{\emptyset \neq S \subseteq \{1, \dots, |m|\}} (-1)^{|S|} \sum_{T \subsetneq S} \sum_{R \subseteq -S} \mathbb{1}\{R \cup T = A\}.$$

Here, the indicator $\mathbb{1}\{A = \emptyset\}$ accounts for the term $\int \ell(y | x; \theta, \eta) [D_{\eta}^m u(Z; \theta, \eta, \mu)] dy$ in (A.9), while the second term in κ_A counts the contributions from the triple sum. Our goal is to show that $\kappa_A = 0$ for all $A \subseteq \{1, \dots, |m|\}$. We analyze two cases separately:

- **Case 1:** For $A = \emptyset$, we note that $R \cup T = \emptyset$ implies $R = T = \emptyset$. Thus, the indicator $\mathbb{1}\{R \cup T = \emptyset\}$ is non-zero only when $R = \emptyset$ and $T = \emptyset$, implying that

$$\kappa_{\emptyset} = 1 + \sum_{\emptyset \neq S \subseteq \{1, \dots, |m|\}} (-1)^{|S|} = \sum_{S \subseteq \{1, \dots, |m|\}} (-1)^{|S|} = 0,$$

where in the second step we used that $1 = (-1)^{|\emptyset|}$ to include that term into the sum over S , and the final step is the alternating sum result we already used in the proof of Lemma 2 above.

- **Case 2:** Next, consider $A \neq \emptyset$. For given A and S , we have

$$\sum_{T \subsetneq S} \sum_{R \subseteq -S} \mathbb{1}\{R \cup T = A\} = \begin{cases} 1 & \text{if } S \not\subseteq A, \\ 0 & \text{if } S \subseteq A. \end{cases} \quad (\text{A.11})$$

If $S \not\subseteq A$ (i.e. not $S \subseteq A$), then (A.11) holds because a solution to the conditions $R \cup T = A$, $T \subsetneq S$, $R \subseteq -S$ exists and is uniquely given by $T = A \cap S$ and $R = A \cap (-S)$. Uniqueness of the pair (T, R) implies that $\sum_{T \subsetneq S} \sum_{R \subseteq -S} \mathbb{1}\{R \cup T = A\} = 1$ in that case. However, if $S \subseteq A$ then no solution for the pair (T, R) exists (because $T = A \cap S$ implies $T = S$ in that case, which contradicts the condition $T \subsetneq S$), implying that the expression in (A.11) is indeed zero then. Using (A.11) we now find that

$$\begin{aligned} \kappa_A &= \sum_{\substack{\emptyset \neq S \subseteq \{1, \dots, |m|\} \\ S \not\subseteq A}} (-1)^{|S|} \\ &= \sum_{\emptyset \neq S \subseteq \{1, \dots, |m|\}} (-1)^{|S|} - \sum_{\emptyset \neq S \subseteq A} (-1)^{|S|} \\ &= \left[\sum_{S \subseteq \{1, \dots, |m|\}} (-1)^{|S|} - (-1)^{|\emptyset|} \right] - \left[\sum_{S \subseteq A} (-1)^{|S|} - (-1)^{|\emptyset|} \right] \\ &= [0 - 1] - [0 - 1] = 0. \end{aligned}$$

We have thus shown that $\kappa_A = 0$ for all $A \subseteq \{1, \dots, |m|\}$. By (A.10) we thus have $\mathbb{E}_{\theta, \eta} [D_\eta^m u_q^*(Z; \theta, \eta, \mu) | X = x] = 0$, which can also be written as

$$\mathbb{E}_{\theta, \eta} [\nabla_\eta^q u_q^*(Z; \theta, \eta, \mu) | X = x] = 0.$$

Plugging in the true parameter values θ_0 and η_0 we thus find

$$\mathbb{E} [\nabla_\eta^q u_q^*(Z; \theta_0, \eta_0, \mu) | X = x] = 0,$$

and by the law of iterated expectations also

$$\mathbb{E} [\nabla_\eta^q u_q^*(Z; \theta_0, \eta_0, \mu)] = 0.$$

Remarkably, this result holds for any value of μ .

□

A.2 Proof of Lemma 1

Let $\tau := (\theta, \mu)$. We write $u_{k,i}(\tau, \eta_i)$ for $u_{k,i}(Z_i; \theta, \eta_i, \mu)$, the k th component of the $\dim(u_i)$ -vector $u_i(Z_i; \theta, \eta_i, \mu)$. Furthermore, compared to the statement of the lemma we drop all subscripts n in the following derivations. In particular, for $\mathbb{K}_{q,n,i}$ we simply write $\mathbb{K}_{q,i}$. By a mean-value expansion of $\widehat{\eta}_i$ around η_{i0} we obtain

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^N u_{k,i}(Z_i; \widehat{\theta}, \widehat{\eta}_i, \widehat{\mu}) &= \frac{1}{n} \sum_{i=1}^N u_{k,i}(\widehat{\tau}, \widehat{\eta}_i) \\ &= \frac{1}{n} \sum_{i=1}^N u_{k,i}(\widehat{\tau}, \eta_{i0}) \\ &\quad + \frac{1}{n} \sum_{i=1}^N \sum_{m \in \mathbb{K}_{q,i}} \frac{1}{m!} [D_{\eta_i}^m u_{k,i}(\widehat{\tau}, \eta_{i0})] (\widehat{\eta}_i - \eta_{i0})^m \\ &\quad + \frac{1}{n} \sum_{i=1}^N \sum_{m \in \mathbb{K}_{q+1,i}} \frac{1}{m!} [D_{\eta_i}^m u_{k,i}(\widehat{\tau}, \widetilde{\eta}_i)] (\widehat{\eta}_i - \eta_{i0})^m, \end{aligned}$$

where $m! = \prod_r (m_r!)$, and $\widetilde{\eta}_i$ is some value between η_{i0} and $\widehat{\eta}_i$. Next, we perform a mean-value expansions in $\widehat{\tau}$ around τ_0 to obtain

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^N u_{k,i}(Z_i; \widehat{\theta}, \widehat{\eta}_i, \widehat{\mu}) &= \frac{1}{n} \sum_{i=1}^N u_{k,i}(\tau_0, \eta_{i0}) \\ &\quad + \underbrace{\frac{1}{n} \sum_{i=1}^N \left[\frac{\partial}{\partial \tau} u_{k,i}(\tau_0, \eta_{i0}) \right]^\top (\widehat{\tau} - \tau_0)}_{= [G_\mu(\widehat{\mu} - \mu_0) + G_\theta(\widehat{\theta} - \theta_0)]_k + o_P(\|\widehat{\tau} - \tau_0\|)} \\ &\quad + \underbrace{\frac{1}{2} (\widehat{\tau} - \tau_0)^\top \left\{ \frac{1}{n} \sum_{i=1}^N \left[\frac{\partial^2}{\partial \tau \partial \tau^\top} u_{k,i}(\widetilde{\tau}, \eta_{i0}) \right] \right\}}_{=: B_{1,k}} (\widehat{\tau} - \tau_0) \\ &\quad + \underbrace{\frac{1}{n} \sum_{i=1}^N \sum_{m \in \mathbb{K}_{q,i}} \frac{1}{m!} [D_{\eta_i}^m u_{k,i}(\tau_0, \eta_{i0})] (\widehat{\eta}_i - \eta_{i0})^m}_{= n^{-1/2} R_{n,k}, \text{ the } k\text{'th component of } R_n \text{ defined in the lemma.}} \\ &\quad + \underbrace{\frac{1}{n} \sum_{i=1}^N \sum_{m \in \mathbb{K}_{q,i}} \frac{1}{m!} \left[D_{\eta_i}^m \frac{\partial}{\partial \tau} u_{k,i}(\overline{\tau}, \eta_{i0}) \right]^\top (\widehat{\tau} - \tau_0) (\widehat{\eta}_i - \eta_{i0})^m}_{=: B_{2,k}} \end{aligned}$$

$$+ \underbrace{\frac{1}{n} \sum_{i=1}^N \sum_{m \in \mathbb{K}_{q+1,i}} \frac{1}{m!} [D_{\eta_i}^m u_{k,i}(\hat{\tau}, \tilde{\eta}_i)] (\hat{\eta}_i - \eta_{i0})^m}_{=: B_{3,k}},$$

where $\tilde{\tau}$ and $\bar{\tau}$ are values between $\hat{\tau}$ and τ_0 . Denote the dimensions of the parameters θ and μ by d_θ and d_μ , respectively. Our assumptions guarantee that

$$\begin{aligned} |B_{1,k}| &\leq \frac{(d_\theta + d_\mu)^2}{2} \|\hat{\tau} - \tau_0\|^2 \frac{1}{n} \sum_{i=1}^N n_i C(Z_i) \\ &\leq \frac{(d_\theta + d_\mu)^2}{2} \|\hat{\tau} - \tau_0\|^2 \underbrace{\left(\frac{1}{n} \sum_{i=1}^N n_i [C(Z_i)]^2 \right)^{1/2}}_{=O_P(1)} \underbrace{\left(\frac{1}{n} \sum_{i=1}^N n_i \right)^{1/2}}_{=O(1)} \\ &= O_P(\|\hat{\tau} - \tau_0\|^2), \\ |B_{2,k}| &\leq (d_\theta + d_\mu) \|\hat{\tau} - \tau_0\| \frac{1}{n} \sum_{i=1}^N n_i C(Z_i) \sum_{m \in \mathbb{K}_{q,i}} \frac{1}{m!} (\hat{\eta}_i - \eta_{i0})^m \\ &\leq (d_\theta + d_\mu) \|\hat{\tau} - \tau_0\| \underbrace{\left(\frac{1}{n} \sum_{i=1}^N n_i [C(Z_i)]^2 \right)^{1/2}}_{=O_P(1)} \underbrace{\left(\frac{1}{n} \sum_{i=1}^N n_i \left[\sum_{m \in \mathbb{K}_{q,i}} \frac{1}{m!} (\hat{\eta}_i - \eta_{i0})^m \right]^2 \right)^{1/2}}_{=o_P(1)} \\ &= o_P(\|\hat{\tau} - \tau_0\|), \\ |B_{3,k}| &\leq \frac{1}{n} \sum_{i=1}^N n_i C(Z_i) \|\hat{\eta}_i - \eta_{i0}\|^{q+1} \underbrace{\sum_{m \in \mathbb{K}_{q+1,i}} \frac{1}{m!}}_{=O(1)} \\ &= O(1) \underbrace{\left(\frac{1}{n} \sum_{i=1}^N n_i [C(Z_i)]^2 \right)^{1/2}}_{=O_P(1)} \underbrace{\left(\frac{1}{n} \sum_{i=1}^N n_i \|\hat{\eta}_i - \eta_{i0}\|^{2(q+1)} \right)^{1/2}}_{=o_P(n^{-1/2})}. \\ &= o_P(n^{-1/2}). \end{aligned}$$

Here, in addition to our assumption we also used the Cauchy-Schwarz inequality. We have thus shown that

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^N u_i(Z_i; \hat{\theta}, \hat{\eta}_i, \hat{\mu}) &= G_\mu(\hat{\mu} - \mu_0) + G_\theta(\hat{\theta} - \theta_0) + \frac{1}{n} \sum_{i=1}^N u_i(Z_i; \theta_0, \eta_{i0}, \mu_0) + n^{-1/2} R_n \\ &\quad + O_P(\|\hat{\tau} - \tau_0\|^2) + o_P(\|\hat{\tau} - \tau_0\|) + o_P(n^{-1/2}). \end{aligned}$$

Using our assumptions on the convergence of $\hat{\mu}$ and $\hat{\theta}$ we thus have

$$\begin{aligned} \frac{1}{\sqrt{n}} \sum_{i=1}^n u_i(Z_i; \hat{\theta}, \hat{\eta}_i, \hat{\mu}) &= G_\mu [\sqrt{n}(\hat{\mu} - \mu_0)] + o_P(\|\hat{\mu} - \mu_0\|) \\ &\quad + \frac{1}{\sqrt{n}} \sum_{i=1}^N [u_i(Z_i; \theta_0, \eta_{i0}, \mu_0) + G_\theta \psi_i] + R_n + o_P(1). \end{aligned} \quad (\text{A.12})$$

By Assumption 1(i) we have

$$\left[\frac{1}{n} \sum_{i=1}^N \frac{\partial u_i^\top(Z_i; \hat{\theta}, \hat{\eta}_i, \hat{\mu})}{\partial \mu} \right] W \left[\frac{1}{\sqrt{n}} \sum_{i=1}^N u_i(Z_i; \hat{\theta}, \hat{\eta}_i, \hat{\mu}) \right] = o_P(1).$$

By using Assumption 1(iv) and (vii) we thus have

$$G_\mu^\top W \left[\frac{1}{\sqrt{n}} \sum_{i=1}^N u_i(Z_i; \hat{\theta}, \hat{\eta}_i, \hat{\mu}) \right] = o_P(1).$$

Plugging the approximation in (A.12) into the last display gives

$$\begin{aligned} o_P(1) &= G_\mu^\top W \left[\frac{1}{\sqrt{n}} \sum_{i=1}^N u_i(Z_i; \hat{\theta}, \hat{\eta}_i, \hat{\mu}) \right] \\ &= (G_\mu^\top W G_\mu) [\sqrt{n}(\hat{\mu} - \mu_0)] + o_P(\|\hat{\mu} - \mu_0\|) \\ &\quad + G_\mu^\top W \left\{ \frac{1}{\sqrt{n}} \sum_{i=1}^N [u_i(Z_i; \theta_0, \eta_{i0}, \mu_0) + G_\theta \psi_i] + R_n \right\} + o_P(1). \end{aligned}$$

Since $G_\mu^\top W G_\mu$ is full rank, solving for $\sqrt{n}(\hat{\mu} - \mu_0)$ gives the statement of the lemma.

A.3 Proof of Theorem 2

We again drop all subscripts n in the derivations. Let $\xi_i = u_i(Z_i; \theta_0, \eta_{i0}, \mu_0) + G_\theta \psi_i$. From Lemma 1, we have

$$\sqrt{n}(\hat{\mu} - \mu_0) = (G_\mu^\top W G_\mu)^{-1} G_\mu^\top W \left\{ \frac{1}{\sqrt{n}} \sum_{i=1}^N \xi_i + R_n \right\} + o_P(1).$$

We will show that $R_n = o_P(1)$ under our assumptions. First, by Assumption 2(i), the moment function is Neyman-orthogonal to order q , which implies

$$\mathbb{E} [D_{\eta_i}^m u_i(Z_i; \theta_0, \eta_{i0}, \mu_0)] = 0$$

for all $m \in \mathbb{K}_{q,i}$. Therefore, R_n is a sum of mean-zero terms. Next, by Assumption 1(iv), the derivatives $D_{\eta_i}^m u_i(Z_i; \theta_0, \eta_{i0}, \mu_{0,n})$ are all bounded by $n_i C(Z_i)$ with $\frac{1}{n} \sum_{i=1}^N n_i \mathbb{E}[C(Z_i)^2] = O(1)$. Using this together with Assumption 1(v) and Assumption 2(ii), one obtains $\mathbb{E}[R_n^2] = o(1)$. By Chebyshev's inequality we thus have $R_n = o_P(1)$. Thus, we have

$$\sqrt{n}(\hat{\mu}_n - \mu_{0,n}) = (G_\mu^\top W G_\mu)^{-1} G_\mu^\top W \left\{ \frac{1}{\sqrt{n}} \sum_{i=1}^N \xi_i \right\} + o_P(1).$$

By Assumption 2(ii), (iii), the terms ξ_i are independent across i . Furthermore, by Assumption 2(iv), they satisfy Lindeberg's condition and have a well-defined variance limit V_ξ . Therefore, by the Lindeberg-Feller Central Limit Theorem:

$$\frac{1}{\sqrt{n}} \sum_{i=1}^N \xi_i \xrightarrow{d} \mathcal{N}(0, V_\xi).$$

The conclusion follows by the continuous mapping theorem, giving us

$$\sqrt{n}(\hat{\mu} - \mu_0) \xrightarrow{d} \mathcal{N}\left(0, (G_\mu^\top W G_\mu)^{-1} G_\mu^\top W V_\xi W G_\mu (G_\mu^\top W G_\mu)^{-1}\right).$$

ONLINE SUPPLEMENT

B Linear network regression

B.1 Model and results

Consider the linear regression model

$$Y = X\eta + \sigma\varepsilon, \quad \varepsilon | X \sim iid\mathcal{N}(0, I_n), \quad (2.1)$$

where n is the dimension of Y . We assume that $X^\top X$ is nonsingular with probability one. Model (2.1) nests the Neyman-Scott model (2.2), for $n = NT$ and $X = I_N \otimes \iota_T$, with I_N the $N \times N$ identity matrix and ι_T the $T \times 1$ vector of ones. Model (2.1) also nests settings where X is a network matrix, as in the log wage regression model of [Abowd, Kramarz and Margolis \(1999\)](#) based on linked worker-firm panel data, in which case η is a vector stacking worker and firm fixed-effects. Our goal is to estimate $\mu = \eta^\top Q\eta$ for some symmetric $r \times r$ matrix Q , where r denotes the dimension of η . Such quadratic forms are of interest in panel and network variance decompositions (e.g., [Arellano and Bonhomme, 2012](#), [Andrews, Gill, Schank and Upward, 2008](#), [Kline, Saggio and Sølvesten, 2020](#)).

Suppose to start with that σ^2 is known. Theorem 1 implies the following characterization of the first- and second-order estimating equations for μ , based on $u(y, x; \sigma^2, \eta, \mu) = \mu - \eta^\top Q\eta$.

Proposition 1.

$$\begin{aligned} u_1^*(y, x; \sigma^2, \eta, \mu) &= \mu - \eta^\top Q\eta - 2\eta^\top Q^\top (x^\top x)^{-1} x^\top (y - x\eta), \\ u_2^*(y, x; \sigma^2, \eta, \mu) &= \mu - y^\top x (x^\top x)^{-1} Q (x^\top x)^{-1} x^\top y + \sigma^2 \text{Trace}(Q(x^\top x)^{-1}). \end{aligned}$$

Hence, given a preliminary estimator $\hat{\eta}$, the associated first-order orthogonalized estimator of μ_0 is

$$\hat{\mu}_1 = \hat{\eta}^\top Q\hat{\eta} + 2\eta^\top Q^\top (x^\top x)^{-1} x^\top (y - x\hat{\eta}).$$

It is easy to see $\mathbb{E}_{\theta,\eta}[\widehat{\mu}_1] \neq \mu$. In turn, the second-order orthogonalized estimator is

$$\widehat{\mu}_2 = y^\top x(x^\top x)^{-1}Q(x^\top x)^{-1}x^\top y - \sigma^2 \text{Trace}(Q(x^\top x)^{-1}).$$

Note that $\widehat{\mu}_2$ does not depend on the preliminary estimate $\widehat{\eta}$, and that $\mathbb{E}_{\theta,\eta}[\widehat{\mu}_2] = \mu$. Hence, second-order Neyman-orthogonality leads to exact unbiased in this case. The expression coincides with the trace correction of [Andrews, Gill, Schank and Upward \(2008\)](#).

Turning to the estimation of σ^2 , we rely on the score

$$u(y, x; \sigma^2, \eta) = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4}(y - x\eta)^\top(y - x\eta).$$

Using [Theorem 1](#), we obtain the following characterization of the first- and second-order orthogonalized scores.

Proposition 2.

$$\begin{aligned} u_1^*(y, x; \sigma^2, \eta) &= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4}(y - x\eta)^\top(y - x\eta), \\ u_2^*(y, x; \sigma^2, \eta) &= -\frac{n - \text{Trace}(x(x^\top x)^{-1}x^\top)}{2\sigma^2} + \frac{1}{2\sigma^4}y^\top(I_n - x(x^\top x)^{-1}x^\top)y. \end{aligned}$$

As in the special case of the Neyman-Scott model, first-order orthogonalization is immaterial, and the first-order orthogonalized estimator of σ^2 is

$$\widehat{\sigma}_1^2 = \frac{(Y - X\widehat{\eta})^\top(Y - X\widehat{\eta})}{n},$$

and $\mathbb{E}_{\theta,\eta}[\widehat{\sigma}_1^2] \neq \sigma^2$. In turn, the second-order orthogonalized estimator is

$$\widehat{\sigma}^2 = \frac{Y^\top(I_n - X(X^\top X)^{-1}X^\top)Y}{n - \text{Trace}(X(X^\top X)^{-1}X^\top)}, \quad (2.2)$$

which is the familiar degree of freedom correction, exactly unbiased in this case, and independent of the preliminary estimator $\widehat{\eta}$.⁷ In the special case of the Neyman-Scott model, [\(2.2\)](#) simplifies to

$$\widehat{\sigma}^2 = \frac{1}{N(T-1)} \sum_{i=1}^N \sum_{t=1}^T (Y_{it} - \bar{Y}_i)^2,$$

where $\bar{Y}_i = \frac{1}{T} \sum_{t=1}^T Y_{it}$, which is exactly unbiased for fixed T and N .

⁷Note it is not necessary for $X^\top X$ to be non-singular for $\widehat{\sigma}^2$ to be well-defined, provided one replaces $(X^\top X)^{-1}$ by a generalized inverse.

B.2 Main proofs

Proof of Proposition 1. Let

$$u(Y, X; \theta, \eta, \mu) = \mu - \eta^\top Q \eta.$$

We have

$$\log \ell(Y | X; \theta, \eta) = -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (Y - X\eta)^\top (Y - X\eta).$$

Hence,

$$v_1(Y, X; \theta, \eta) = \frac{1}{\sigma^2} X^\top (Y - X\eta),$$

and

$$v_2(Y, X; \theta, \eta) = \text{vech} \left(-\frac{1}{\sigma^2} X^\top X + \frac{1}{\sigma^4} X^\top (Y - X\eta)(Y - X\eta)^\top X \right),$$

where $\text{vech}(C)$ denotes the half-vectorization of a symmetric matrix C .

By Theorem 1 we have

$$u_2^*(Y, X; \theta, \eta, \mu) = u(Y, X; \theta, \eta, \mu) - A^\top \begin{pmatrix} v_1(Y, X; \theta, \eta) \\ v_2(Y, X; \theta, \eta) \end{pmatrix},$$

where

$$A = - \left\{ \mathbb{E} \begin{bmatrix} v_1(Y, X; \theta, \eta)v_1(Y, X; \theta, \eta)^\top & v_1(Y, X; \theta, \eta)v_2(Y, X; \theta, \eta)^\top \\ v_2(Y, X; \theta, \eta)v_1(Y, X; \theta, \eta)^\top & v_2(Y, X; \theta, \eta)v_2(Y, X; \theta, \eta)^\top \end{bmatrix} \right\}^{-1} \begin{pmatrix} -2Q\eta \\ -2\text{vech}(Q) \end{pmatrix},$$

where for conciseness we omit the dependence of A on X , θ , and η from the notation, and we implicitly condition on X in all expectations.

Note

$$v_1(Y, X; \theta, \eta) = \frac{1}{\sigma^2} X^\top \varepsilon,$$

and

$$v_2(Y, X; \theta, \eta) = \text{vech} \left(-\frac{1}{\sigma^2} X^\top X + \frac{1}{\sigma^4} X^\top \varepsilon \varepsilon^\top X \right).$$

Hence

$$\begin{aligned}
\mathbb{E}[v_1(Y, X; \theta, \eta)v_1(Y, X; \theta, \eta)^\top] &= \frac{1}{\sigma^2}X^\top X, \\
\mathbb{E}[v_1(Y, X; \theta, \eta)v_2(Y, X; \theta, \eta)^\top] &= 0, \\
\mathbb{E}[v_2(Y, X; \theta, \eta)v_2(Y, X; \theta, \eta)^\top] \\
&= \mathbb{E} \left[\text{vech} \left(-\frac{1}{\sigma^2}X^\top X + \frac{1}{\sigma^4}X^\top \varepsilon \varepsilon^\top X \right) \text{vech} \left(-\frac{1}{\sigma^2}X^\top X + \frac{1}{\sigma^4}X^\top \varepsilon \varepsilon^\top X \right)^\top \right].
\end{aligned}$$

Let L_m denote the elimination matrix such that $\text{vech}(Q) = L\text{vec}(Q)$ (Magnus and Neudecker, 1980). Let K_n denote the commutation matrix such that $K_n \text{vec}(A) = \text{vec}(A^\top)$ (Magnus and Neudecker, 1979). Note that $K_n = K_n^\top$. We have the following result.

Lemma 3.

$$\mathbb{E}[v_2(Y, X; \theta, \eta)v_2(Y, X; \theta, \eta)^\top] = \frac{1}{\sigma^4}L_m(X^\top \otimes X^\top)[I_{n^2} + K_n](X \otimes X)L_m^\top.$$

It follows from the above that

$$\begin{aligned}
u_2^*(Y, X; \theta, \eta, \mu) &= \mu - \eta^\top Q \eta - 2\eta^\top Q^\top (X^\top X)^{-1} X^\top (Y - X\eta) \\
&\quad - 2\text{vech}(Q)^\top [L_m(X^\top \otimes X^\top)[I_{n^2} + K_n](X \otimes X)L_m^\top]^{-1} \\
&\quad \times \text{vech}(-\sigma^2 X^\top X + X^\top (Y - X\eta)(Y - X\eta)^\top X).
\end{aligned}$$

The following lemma is instrumental.

Lemma 4. *Let A and B be symmetric matrices. Then*

$$\text{vech}(A)^\top [L_m(X^\top \otimes X^\top)[I_{n^2} + K_n](X \otimes X)L_m^\top]^{-1} \text{vech}(B) = \frac{1}{2}\text{Trace}(A(X^\top X)^{-1}B(X^\top X)^{-1}).$$

By Lemma 4 applied to $A = Q$ and $B = X^\top (Y - X\eta)(Y - X\eta)^\top X - \sigma^2 X^\top X$, we then have

$$\begin{aligned}
u_2^*(Y, X; \theta, \eta, \mu) &= \mu - \eta^\top Q \eta - 2\eta^\top Q^\top (X^\top X)^{-1} X^\top (Y - X\eta) \\
&\quad - \text{Trace}(Q(X^\top X)^{-1}[X^\top (Y - X\eta)(Y - X\eta)^\top X - \sigma^2 X^\top X](X^\top X)^{-1}) \\
&= \mu - Y^\top X(X^\top X)^{-1}Q(X^\top X)^{-1}X^\top Y + \sigma^2 \text{Trace}(Q(X^\top X)^{-1}).
\end{aligned}$$

The associated second-order Neyman-orthogonal estimator is then

$$\hat{\mu} = Y^\top X(X^\top X)^{-1}Q(X^\top X)^{-1}X^\top Y - \sigma^2 \text{Trace}(Q(X^\top X)^{-1}),$$

which corresponds to the trace correction of [Andrews, Gill, Schank and Upward \(2008\)](#), for fixed σ^2 .

Proof of Proposition 2 Let

$$u(Y, X; \sigma^2, \eta) = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4}(Y - X\eta)^\top(Y - X\eta).$$

By [Theorem 1](#) we have

$$u_2^*(Y, X; \theta, \eta) = u(Y, X; \theta, \eta) - A^\top \begin{pmatrix} v_1(Y, X; \theta, \eta) \\ v_2(Y, X; \theta, \eta) \end{pmatrix},$$

where

$$A = \mathbb{E} \begin{bmatrix} v_1(Y, X; \theta, \eta)v_1(Y, X; \theta, \eta)^\top & v_1(Y, X; \theta, \eta)v_2(Y, X; \theta, \eta)^\top \\ v_2(Y, X; \theta, \eta)v_1(Y, X; \theta, \eta)^\top & v_2(Y, X; \theta, \eta)v_2(Y, X; \theta, \eta)^\top \end{bmatrix}^{-1} \\ \times \mathbb{E} \begin{bmatrix} v_1(Y, X; \theta, \eta)u(Y, X; \theta, \eta) \\ v_2(Y, X; \theta, \eta)u(Y, X; \theta, \eta) \end{bmatrix}.$$

We have the following result.

Lemma 5.

$$\mathbb{E}[v_1(Y, X; \theta, \eta)u(Y, X; \theta, \eta)] = 0,$$

and

$$\mathbb{E}[v_2(Y, X; \theta, \eta)u(Y, X; \theta, \eta)] = \frac{1}{2\sigma^4}L_m(X^\top \otimes X^\top)(I_{n^2} + K_n)\text{vec}(I_n).$$

Using [Lemma 5](#), we have

$$u_2^*(Y, X; \theta, \eta) = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4}(Y - X\eta)^\top(Y - X\eta) \\ - \frac{1}{2}\text{vec}(I_n)^\top(I_{n^2} + K_n)(X \otimes X)L_m^\top [L_m(X^\top \otimes X^\top)[I_{n^2} + K_n](X \otimes X)L_m^\top]^{-1} \\ \times L_m \text{vec} \left(-\frac{1}{\sigma^2}X^\top X + \frac{1}{\sigma^4}X^\top(Y - X\eta)(Y - X\eta)^\top X \right).$$

Lemma 6. *We equivalently have*

$$\begin{aligned} u_2^*(Y, X; \theta, \eta) &= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4}(Y - X\eta)^\top(Y - X\eta) \\ &\quad + \frac{1}{2\sigma^2}\text{Trace}(X(X^\top X)^{-1}X^\top) - \frac{1}{2\sigma^4}(Y - X\eta)^\top X(X^\top X)^{-1}X^\top(Y - X\eta). \end{aligned}$$

By Lemma 6, the second-order orthogonalized score is independent of η and is given by

$$u_2^*(Y, X; \theta, \eta) = -\frac{n - \text{Trace}(X(X^\top X)^{-1}X^\top)}{2\sigma^2} + \frac{1}{2\sigma^4}Y^\top(I_n - X(X^\top X)^{-1}X^\top)Y.$$

B.3 Proofs of intermediate lemmas

Proof of Lemma 3. We have

$$\begin{aligned} &\mathbb{E}[v_2(Y, X; \theta, \eta)v_2(Y, X; \theta, \eta)^\top] \\ &= \mathbb{E}\left[\text{vech}\left(-\frac{1}{\sigma^2}X^\top X + \frac{1}{\sigma^4}X^\top \varepsilon \varepsilon^\top X\right) \text{vech}\left(-\frac{1}{\sigma^2}X^\top X + \frac{1}{\sigma^4}X^\top \varepsilon \varepsilon^\top X\right)^\top\right] \\ &= L_m(X^\top \otimes X^\top) \mathbb{E}\left[\left(\frac{1}{\sigma^4}\varepsilon \otimes \varepsilon - \frac{1}{\sigma^2}\text{vec}(I_n)\right) \left(\frac{1}{\sigma^4}\varepsilon \otimes \varepsilon - \frac{1}{\sigma^2}\text{vec}(I_n)\right)^\top\right] (X \otimes X)L_m^\top \\ &= L_m(X^\top \otimes X^\top) \mathbb{E}\left[\frac{1}{\sigma^8}(\varepsilon \varepsilon^\top) \otimes (\varepsilon \varepsilon^\top) - \frac{1}{\sigma^4}\text{vec}(I_n)\text{vec}(I_n)^\top\right] (X \otimes X)L_m^\top. \end{aligned}$$

Now, by (4.3) in Ghazal and Neudecker (2000), we have, since $\varepsilon \varepsilon^\top \sim \mathcal{W}_n(\sigma^2 I_n, 1)$,

$$\mathbb{E}[(\varepsilon \varepsilon^\top) \otimes (\varepsilon \varepsilon^\top)] = \sigma^4 \text{vec}(I_n)\text{vec}(I_n)^\top + \sigma^4(I_{n^2} + K_n)(I_n \otimes I_n).$$

It follows that

$$\mathbb{E}[v_2(Y, X; \theta, \eta)v_2(Y, X; \theta, \eta)^\top] = \frac{1}{\sigma^4}L_m(X^\top \otimes X^\top)(I_{n^2} + K_n)(X \otimes X)L_m^\top.$$

This shows Lemma 3.

Proof of Lemma 4. Let D_m denote the duplication matrix, such, that for any symmetric matrix C , $D_m \text{vech}(C) = \text{vec}(C)$. We will make use of the following properties (Magnus

and Neudecker, 1980):

$$D_m = (I_{m^2} + K_m)L_m^\top (L_m(I_{m^2} + K_m)L_m^\top)^{-1}, \quad (2.3)$$

$$(I_{n^2} + K_n)(X \otimes X) = (X \otimes X)(I_{m^2} + K_m), \quad (2.4)$$

$$D_m L_m (I_{m^2} + K_m) = (I_{m^2} + K_m), \quad (2.5)$$

$$K_m D_m = D_m. \quad (2.6)$$

Note that $\text{vech}(A) = L_m D_m \text{vech}(A)$ and $\text{vech}(B) = L_m D_m \text{vech}(B)$. Hence

$$\begin{aligned} & \text{vech}(A)^\top [L_m(X^\top \otimes X^\top)[I_{n^2} + K_n](X \otimes X)L_m^\top]^{-1} \text{vech}(B) \\ &= \text{vech}(A)^\top D_m^\top L_m^\top [L_m(X^\top \otimes X^\top)[I_{n^2} + K_n](X \otimes X)L_m^\top]^{-1} L_m D_m \text{vech}(B) \\ &= \text{vech}(A)^\top (L_m(I_{m^2} + K_m)L_m^\top)^{-1} L_m(I_{m^2} + K_m)L_m^\top [L_m(X^\top \otimes X^\top)[I_{n^2} + K_n](X \otimes X)L_m^\top]^{-1} \\ & \quad \times L_m(I_{m^2} + K_m)L_m^\top (L_m(I_{m^2} + K_m)L_m^\top)^{-1} \text{vech}(B) \quad \text{by (2.3)} \\ &= \text{vech}(A)^\top (L_m(I_{m^2} + K_m)L_m^\top)^{-1} L_m(I_{m^2} + K_m)L_m^\top [L_m(X^\top \otimes X^\top)(X \otimes X)(I_{m^2} + K_m)L_m^\top]^{-1} \\ & \quad \times L_m(I_{m^2} + K_m)L_m^\top (L_m(I_{m^2} + K_m)L_m^\top)^{-1} \text{vech}(B) \quad \text{by (2.4)} \\ &= \text{vech}(A)^\top [L_m(X^\top \otimes X^\top)(X \otimes X)(I_{m^2} + K_m)L_m^\top]^{-1} \text{vech}(B). \end{aligned}$$

Now, we have

$$\begin{aligned} & L_m(X^\top \otimes X^\top)(X \otimes X)(I_{m^2} + K_m)L_m^\top D_m^\top ((X^\top X)^{-1} \otimes (X^\top X)^{-1})D_m \\ &= L_m(X^\top \otimes X^\top)(X \otimes X)(I_{m^2} + K_m)((X^\top X)^{-1} \otimes (X^\top X)^{-1})D_m \quad \text{by (2.5)} \\ &= L_m(I_{m^2} + K_m)(X^\top \otimes X^\top)(X \otimes X)((X^\top X)^{-1} \otimes (X^\top X)^{-1})D_m \quad \text{by (2.4)} \\ &= L_m(I_{m^2} + K_m)((X^\top X) \otimes (X^\top X))((X^\top X)^{-1} \otimes (X^\top X)^{-1})D_m \\ &= L_m(I_{m^2} + K_m)D_m \\ &= 2L_m D_m \quad \text{by (2.6)} \\ &= 2I_{m^2}. \end{aligned}$$

As a result,

$$[L_m(X^\top \otimes X^\top)(X \otimes X)(I_{m^2} + K_m)L_m^\top]^{-1} = \frac{1}{2}D_m^\top ((X^\top X)^{-1} \otimes (X^\top X)^{-1})D_m.$$

Hence

$$\begin{aligned}
& \text{vech}(A)^\top [L_m(X^\top \otimes X^\top)(X \otimes X)(I_{m^2} + K_m)L_m^\top]^{-1} \text{vech}(B) \\
&= \frac{1}{2} \text{vech}(A)^\top D_m^\top ((X^\top X)^{-1} \otimes (X^\top X)^{-1}) D_m \text{vech}(B) \\
&= \frac{1}{2} \text{vec}(A)^\top ((X^\top X)^{-1} \otimes (X^\top X)^{-1}) \text{vec}(B) \\
&= \frac{1}{2} \text{vec}(A)^\top \text{vec}((X^\top X)^{-1} B (X^\top X)^{-1}) \\
&= \frac{1}{2} \text{Trace} (A^\top (X^\top X)^{-1} B (X^\top X)^{-1}) \\
&= \frac{1}{2} \text{Trace} (A (X^\top X)^{-1} B (X^\top X)^{-1})
\end{aligned}$$

since A is symmetric. This shows Lemma 4.

Proof of Lemma 5. Since

$$u(Y, X; \theta, \eta) = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \varepsilon^\top \varepsilon,$$

we have

$$\mathbb{E}[v_1(Y, X; \theta, \eta)u(Y, X; \theta, \eta)] = 0,$$

and

$$\begin{aligned}
\mathbb{E}[v_2(Y, X; \theta, \eta)u(Y, X; \theta, \eta)] &= \mathbb{E} \left[\text{vech} \left(-\frac{1}{\sigma^2} X^\top X + \frac{1}{\sigma^4} X^\top \varepsilon \varepsilon^\top X \right) \left(-\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \varepsilon^\top \varepsilon \right) \right] \\
&= L_m \mathbb{E} \left[\text{vec} \left(\frac{1}{\sigma^4} X^\top \varepsilon \varepsilon^\top X \left(-\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \varepsilon^\top \varepsilon \right) \right) \right] \\
&= L_m \left(-\frac{n}{2\sigma^4} (X^\top \otimes X^\top) \text{vec}(I_n) + \frac{1}{2\sigma^8} (X^\top \otimes X^\top) \mathbb{E} [(\varepsilon \varepsilon^\top) \otimes (\varepsilon \varepsilon^\top)] \text{vec}(I_n) \right) \\
&= L_m \left(-\frac{n}{2\sigma^4} (X^\top \otimes X^\top) \text{vec}(I_n) \right. \\
&\quad \left. + \frac{1}{2\sigma^8} (X^\top \otimes X^\top) (\sigma^4 \text{vec}(I_n) \text{vec}(I_n)^\top + \sigma^4 (I_{n^2} + K_n) (I_n \otimes I_n)) \text{vec}(I_n) \right),
\end{aligned}$$

where we have used the expression for $\mathbb{E} [(\varepsilon \varepsilon^\top) \otimes (\varepsilon \varepsilon^\top)]$ from [Ghazal and Neudecker \(2000\)](#) as in the proof of Lemma 4. It follows that

$$\mathbb{E}[v_2(Y, X; \theta, \eta)u(Y, X; \theta, \eta)] = \frac{1}{2\sigma^4} L_m (X^\top \otimes X^\top) (I_{n^2} + K_n) \text{vec}(I_n).$$

This shows Lemma 5.

Proof of Lemma 6. Using results from the proof of Lemma 4 we have

$$\begin{aligned}
& [L_m(X^\top \otimes X^\top)[I_{n^2} + K_n](X \otimes X)L_m^\top]^{-1} \\
&= [L_m(X^\top \otimes X^\top)(X \otimes X)[I_{m^2} + K_m]L_m^\top]^{-1} \\
&= \frac{1}{2}D_m^\top((X^\top X)^{-1} \otimes (X^\top X)^{-1})D_m.
\end{aligned}$$

Hence

$$\begin{aligned}
u_2^*(Y, X; \theta, \eta) &= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4}(Y - X\eta)^\top(Y - X\eta) \\
&\quad - \frac{1}{4}\text{vec}(I_n)^\top(I_{n^2} + K_n)(X \otimes X)L_m^\top D_m^\top((X^\top X)^{-1} \otimes (X^\top X)^{-1})D_m \\
&\quad \times L_m \text{vec} \left(-\frac{1}{\sigma^2}X^\top X + \frac{1}{\sigma^4}X^\top(Y - X\eta)(Y - X\eta)^\top X \right) \\
&= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4}(Y - X\eta)^\top(Y - X\eta) \\
&\quad - \frac{1}{2}\text{vec}(I_n)^\top(X \otimes X)L_m^\top D_m^\top((X^\top X)^{-1} \otimes (X^\top X)^{-1})D_m \\
&\quad \times L_m(X^\top \otimes X^\top)\text{vec} \left(-\frac{1}{\sigma^2}I_n + \frac{1}{\sigma^4}(Y - X\eta)(Y - X\eta)^\top \right),
\end{aligned}$$

where we have used that

$$\text{vec}(I_n)^\top(I_{n^2} + K_n) = 2\text{vec}(I_n)^\top.$$

Now, for any symmetric matrix A , $D_m L_m \text{vec}(A) = \text{vec}(A)$. Hence we have

$$\begin{aligned}
u_2^*(Y, X; \theta, \eta) &= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4}(Y - X\eta)^\top(Y - X\eta) \\
&\quad - \frac{1}{2}\text{vec}(I_n)^\top(X \otimes X)((X^\top X)^{-1} \otimes (X^\top X)^{-1}) \\
&\quad \times (X^\top \otimes X^\top)\text{vec} \left(-\frac{1}{\sigma^2}I_n + \frac{1}{\sigma^4}(Y - X\eta)(Y - X\eta)^\top \right) \\
&= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4}(Y - X\eta)^\top(Y - X\eta) \\
&\quad + \frac{1}{2\sigma^2}\text{Trace}(X(X^\top X)^{-1}X^\top) - \frac{1}{2\sigma^4}(Y - X\eta)^\top X(X^\top X)^{-1}X^\top(Y - X\eta) \\
&= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4}(Y - X\eta)^\top(Y - X\eta) \\
&\quad + \frac{1}{2\sigma^2}\text{Trace}(X(X^\top X)^{-1}X^\top) - \frac{1}{2\sigma^4}(Y - X\eta)^\top X(X^\top X)^{-1}X^\top(Y - X\eta),
\end{aligned}$$

where we have used that $X(X^\top X)^{-1}X^\top$ is symmetric and idempotent. This shows Lemma 6.

C Implementation in nonlinear regression

C.1 Expression for M

Following Constantine and Savits (1996), consider a multivariate function

$$h(x_1, \dots, x_d) = f[g^{(1)}(x_1, \dots, x_d), \dots, g^{(m)}(x_1, \dots, x_d)],$$

$h_\nu = D_{\mathbf{x}}^\nu h(\mathbf{x}^0)$, $f_\lambda = D_{\mathbf{y}}^\lambda f(\mathbf{y}^0)$, $g_\mu^{(i)} = D_{\mathbf{x}}^\mu g^{(i)}(\mathbf{x}^0)$, $\mathbf{g}_\mu = (g_\mu^{(1)}, \dots, g_\mu^{(m)})$. The Faà di Bruno formula is (Theorem 2.1 in Constantine and Savits, 1996):

$$h_\nu = \sum_{1 \leq |\lambda| \leq n} f_\lambda \underbrace{\sum_{s=1}^n \sum_{p_s(\nu, \lambda)} (\nu!) \prod_{j=1}^s \frac{[\mathbf{g}_{\ell_j}]^{\mathbf{k}_j}}{(\mathbf{k}_j!) [\ell_j!]^{|\mathbf{k}_j|}}}_{\text{elements of } M},$$

where $n = |\nu|$, and

$$p_s(\nu, \lambda) = \left\{ (\mathbf{k}_1, \dots, \mathbf{k}_s; \ell_1, \dots, \ell_s) : |\mathbf{k}_i| > 0, \right. \\ \left. \mathbf{0} \prec \ell_1 \prec \dots \prec \ell_s, \sum_{i=1}^s \mathbf{k}_i = \lambda \text{ and } \sum_{i=1}^s |\mathbf{k}_i| \ell_i = \nu \right\}.$$

C.2 Useful properties of the normal distribution

We first consider the univariate normal case.

Lemma 7. For $m \in \mathbb{R}$ and $\sigma \in [0, \infty)$, let $Y \sim \mathcal{N}(m, \sigma^2)$, with corresponding likelihood function

$$\ell(y | m, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y-m)^2}{2\sigma^2}\right).$$

Let $j, k \in \{0, 1, 2, \dots\}$, and define

$$\kappa_{jk} := \mathbb{E}_{m, \sigma} \left[\frac{1}{\ell(Y | m, \sigma)} \frac{\partial^j \ell(Y | m, \sigma)}{(\partial m)^j} \frac{1}{\ell(Y | m, \sigma)} \frac{\partial^k \ell(Y | m, \sigma)}{(\partial m)^k} \right], \\ \rho_j := \mathbb{E}_{m, \sigma} \left[\frac{1}{\ell(Y | m, \sigma)} \frac{\partial^j \ell(Y | m, \sigma)}{(\partial m)^j} \frac{\partial \log \ell(Y | m, \sigma)}{\partial \sigma} \right].$$

Then,

$$\kappa_{jk} = \mathbb{1}\{j = k\} \frac{j!}{\sigma^{2j}}, \quad \rho_j = \mathbb{1}\{j = 2\} \frac{2}{\sigma^3}.$$

Let $\phi(y) = \frac{1}{\sqrt{2\pi}} \exp(-y^2/2)$ and $\phi^{(j)}(y) = \frac{d^j \phi(y)}{dy^j}$. Hermite polynomials are defined by $h_j(y) = (-1)^j [\phi(y)]^{-1} \phi^{(j)}(y)$. The proof of Lemma 7 is given in Subsection C.4. It crucially relies on the following orthogonality property of Hermite polynomials:

$$\int_{-\infty}^{\infty} h_j(y) h_k(y) \phi(y) dy = \mathbb{1}\{j = k\} j!. \quad (3.1)$$

The result in Lemma 7 is sufficient for our purposes, but more general results can be derived.⁸

Next, we consider a vector of independent normal variables with heteroscedastic means and variances.

Lemma 8. *Let $d \in \{1, 2, 3, \dots\}$. For $m \in \mathbb{R}^d$ and $\sigma \in [0, \infty)^d$, let $\Sigma(\sigma)$ be the $d \times d$ diagonal matrix with diagonal entries σ_i^2 , and let $Y \sim \mathcal{N}(m, \Sigma(\sigma))$. The corresponding likelihood function reads*

$$\ell(y | m, \sigma) = \prod_{i=1}^d \ell(y_i | m_i, \sigma_i), \quad \ell(y_i | m_i, \sigma_i) = \left[\frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{(y_i - m_i)^2}{2\sigma_i^2}\right) \right].$$

Let $j, k \in \{0, 1, 2, \dots\}^d$, $j^* = \sum_{i=1}^d j_i$, $k^* = \sum_{i=1}^d k_i$, and define

$$\begin{aligned} \kappa(j, k) &:= \mathbb{E}_{m, \sigma} \left[\frac{1}{\ell(Y | m, \sigma)} \frac{\partial^{j^*} \ell(Y | m, \sigma)}{\prod_{i=1}^d (\partial m_i)^{j_i}} \frac{1}{\ell(Y | m, \sigma)} \frac{\partial^{k^*} \ell(Y | m, \sigma)}{\prod_{i=1}^d (\partial m_i)^{k_i}} \right], \\ \rho(j, i') &:= \mathbb{E}_{m, \sigma} \left[\frac{1}{\ell(Y | m, \sigma)} \frac{\partial^{j^*} \ell(Y | m, \sigma)}{\prod_{i=1}^d (\partial m_i)^{j_i}} \frac{\partial \log \ell(Y | m, \sigma)}{\partial \sigma_{i'}} \right], \end{aligned}$$

⁸More generally, for $k_1, k_2 \in \{0, 1, 2, \dots\}$ and $j_1, j_2 \in \{0, 1\}$, let

$$\begin{aligned} \kappa_{k_1, k_2, j_1, j_2} &:= \mathbb{E}_{m, \sigma} \left[\frac{1}{\ell(Y | m, \sigma)} \frac{\partial^{k_1+j_1} \ell(Y | m, \sigma)}{(\partial m)^{k_1} (\partial \sigma)^{j_1}} \frac{1}{\ell(Y | m, \sigma)} \frac{\partial^{k_2+j_2} \ell(Y | m, \sigma)}{(\partial m)^{k_2} (\partial \sigma)^{j_2}} \right] \\ &= \int_{-\infty}^{\infty} \frac{1}{\ell(y | m, \sigma)} \frac{\partial^{k_1+j_1} \ell(y | m, \sigma)}{(\partial m)^{k_1} (\partial \sigma)^{j_1}} \frac{\partial^{k_2+j_2} \ell(y | m, \sigma)}{(\partial m)^{k_2} (\partial \sigma)^{j_2}} dy. \end{aligned}$$

One then finds

$$\kappa_{k_1, k_2, j_1, j_2} = \mathbb{1}\{k_1 + 2j_1 = k_2 + 2j_2\} (k_1 + 2j_1)! \sigma^{-[2(k_1+2j_1)-j_1-j_2]}.$$

where $i' \in \{1, \dots, d\}$ in the last line. Then,

$$\kappa(j, k) = \mathbb{1}\{j = k\} \prod_{i=1}^d \frac{j_i!}{\sigma_i^{2j_i}}, \quad \rho(j, i) = \begin{cases} \frac{2}{\sigma_i^3} & \text{if } j_i = 2, \text{ and all other entries of } j \text{ are zero,} \\ 0 & \text{otherwise.} \end{cases}$$

Lemma 8 is an immediate corollary of Lemma 7. Using the independence of the components of Y we find

$$\frac{\partial^{j^*} \ell(y | m, \sigma)}{\prod_{i=1}^d (\partial m)^{j_i}} = \prod_{i=1}^d \frac{\partial^{j_i} \ell(y_i | m_i, \sigma_i)}{(\partial m)^{j_i}},$$

and

$$\kappa(j, k) = \prod_{i=1}^d \kappa_{j_i, k_i}.$$

Plugging in the result for κ_{jk} in Lemma 7 then gives the result for $\kappa(j, k)$ in Lemma 7. Analogously for $\rho(j, i)$.

C.3 Nonlinear regression with normal errors

Model:

$$Y_i = m(X_i; \theta, \eta) + \sigma(X_i; \theta) U_i, \quad U_i \sim iid\mathcal{N}(0, 1), \quad i = 1, \dots, d,$$

where $m(\cdot; \cdot, \cdot)$ and $\sigma(\cdot; \cdot)$ are known functions, and θ and η are unknown parameters. Ignore θ for the moment and write

$$Y_i = m_i(\eta) + \sigma_i U_i.$$

Let $m = (m_1, \dots, m_d)$ and $\sigma = (\sigma_1, \dots, \sigma_d)$. The likelihood for $y = (y_1, \dots, y_d)$ is then given by

$$\ell(y | \eta) = \ell(y | m(\eta), \sigma),$$

where $\ell(y | m, \sigma)$ is given in Lemma 8. Let $\nabla_\eta^{(p)}$ be the vector operator that collects all unique derivatives with respect to η up to order p . Let $\nabla_m^{(p)}$ be the vector operator that

collects all unique derivatives with respect to m up to order p . Then, there exists a matrix valued function $M(\eta)$, which only depends η and on the function $m(\eta)$, such that

$$\nabla_p^{(\eta)} \ell(y | \eta) = M(\eta) \nabla_p^{(m)} \ell(y | m(\eta), \sigma). \quad (3.2)$$

We want to calculate

$$\mathbb{E}_\eta \left[\frac{\nabla_\eta^{(p)} \ell(Y | \eta)}{\ell(Y | \eta)} \frac{\nabla_\eta^{(p)} \ell(Y | \eta)^\top}{\ell(Y | \eta)} \right].$$

Lemma 8 gives us explicit expressions for all the components of

$$\mathbb{E}_{m,\sigma} \left[\frac{\nabla_m^{(p)} \ell(Y | m, \sigma)}{\ell(Y | m, \sigma)} \frac{\nabla_m^{(p)} \ell(Y | m, \sigma)^\top}{\ell(Y | m, \sigma)} \right].$$

Using (3.2) we have

$$\begin{aligned} & \mathbb{E}_\eta \left[\frac{\nabla_\eta^{(p)} \ell(Y | \eta)}{\ell(Y | \eta)} \frac{\nabla_\eta^{(p)} \ell(Y | \eta)^\top}{\ell(Y | \eta)} \right] \\ &= M(\eta) \mathbb{E}_{m(\eta),\sigma} \left[\frac{\nabla_m^{(p)} \ell(Y | m(\eta), \sigma)}{\ell(Y | m(\eta), \sigma)} \frac{\nabla_m^{(p)} \ell(Y | m(\eta), \sigma)^\top}{\ell(Y | m(\eta), \sigma)} \right] M(\eta)^\top. \end{aligned}$$

Thus, by combining Lemma 8 with the multidimensional Faà di Bruno's formula we get explicit expressions for all the matrices we need.

C.4 Proof of Lemma 7

We already introduced $\phi(y) = \frac{1}{\sqrt{2\pi}} \exp(-y^2/2)$ and $\phi^{(j)}(y) = \frac{d^j \phi(y)}{dy^j}$ above. Let $j, k \in \{0, 1, 2, 3, \dots\}$. The well-known orthogonality property of Hermite polynomials in (3.1) can be rewritten as

$$\int_{-\infty}^{\infty} \frac{\phi^{(j)}(y) \phi^{(k)}(y)}{\phi(y)} dy = \mathbb{1}\{j = k\} j!. \quad (3.3)$$

Another well-known property of Hermite polynomials is the recurrence relation $h_{j+1}(y) = yh_j(y) - \frac{d}{dy} h_j(y)$. Using this, it is easy to show that for $j > k$ we have

$$\int_{-\infty}^{\infty} \frac{y \phi^{(j)}(y) \phi^{(k)}(y)}{\phi(y)} dy = -\mathbb{1}\{j = k + 1\} j!. \quad (3.4)$$

Next, we have

$$\ell(y | m, \sigma) = \frac{1}{\sigma} \phi\left(\frac{y-m}{\sigma}\right), \quad \frac{\partial^j \ell(y | m, \sigma)}{(\partial m)^j} = \frac{(-1)^j}{\sigma^{j+1}} \phi^{(j)}\left(\frac{y-m}{\sigma}\right).$$

Using this we obtain

$$\begin{aligned} \kappa_{jk} &:= \mathbb{E}_{m,\sigma} \left[\frac{1}{\ell(Y | m, \sigma)} \frac{\partial^j \ell(Y | m, \sigma)}{(\partial m)^j} \frac{1}{\ell(Y | m, \sigma)} \frac{\partial^k \ell(Y | m, \sigma)}{(\partial m)^k} \right] \\ &= \int_{-\infty}^{\infty} \frac{1}{\ell(y | m, \sigma)} \frac{\partial^j \ell(y | m, \sigma)}{(\partial m)^j} \frac{\partial^k \ell(y | m, \sigma)}{(\partial m)^k} dy \\ &= \frac{(-1)^{j+k}}{\sigma^{j+k+1}} \int_{-\infty}^{\infty} \frac{1}{\phi\left(\frac{y-m}{\sigma}\right)} \phi^{(j)}\left(\frac{y-m}{\sigma}\right) \phi^{(k)}\left(\frac{y-m}{\sigma}\right) dy \\ &= \frac{(-1)^{j+k}}{\sigma^{j+k}} \int_{-\infty}^{\infty} \frac{\phi^{(j)}(y) \phi^{(k)}(y)}{\phi(y)} dy \\ &= \mathbb{1}\{j = k\} \frac{j!}{\sigma^{2j}}, \end{aligned}$$

where the second to last step employs a change of variables in the integral ($\frac{y-m}{\sigma} \mapsto y$), and the last step uses (3.3). Similarly, for

$$\frac{\partial \ell(y | m, \sigma)}{\partial \sigma} = -\frac{1}{\sigma^2} \phi\left(\frac{y-m}{\sigma}\right) - \left(\frac{y-m}{\sigma^3}\right) \phi^{(1)}\left(\frac{y-m}{\sigma}\right),$$

one finds

$$\begin{aligned} \rho_j &:= \mathbb{E}_{m,\sigma} \left[\frac{1}{\ell(Y | m, \sigma)} \frac{\partial^j \ell(Y | m, \sigma)}{(\partial m)^j} \frac{\partial \log \ell(Y | m, \sigma)}{\partial \sigma} \right] \\ &= \int_{-\infty}^{\infty} \frac{1}{\ell(y | m, \sigma)} \frac{\partial^j \ell(y | m, \sigma)}{(\partial m)^j} \frac{\partial \ell(y | m, \sigma)}{\partial \sigma} dy \\ &= \frac{(-1)^{1+j}}{\sigma^{2+j}} \int_{-\infty}^{\infty} \frac{1}{\phi\left(\frac{y-m}{\sigma}\right)} \phi^{(j)}\left(\frac{y-m}{\sigma}\right) \left[\phi\left(\frac{y-m}{\sigma}\right) + \left(\frac{y-m}{\sigma}\right) \phi^{(1)}\left(\frac{y-m}{\sigma}\right) \right] dy \\ &= \frac{(-1)^{1+j}}{\sigma^{1+j}} \int_{-\infty}^{\infty} \frac{1}{\phi(y)} \phi^{(j)}(y) [\phi(y) + y \phi^{(1)}(y)] dy \\ &= \mathbb{1}\{j = 2\} \frac{(-1)}{\sigma^3} \int_{-\infty}^{\infty} \frac{y \phi^{(1)}(y) \phi^{(2)}(y)}{\phi(y)} dy \\ &= \mathbb{1}\{j = 2\} \frac{2}{\sigma^3}. \end{aligned}$$

where we again employed the same change of variables in the integration and also use (3.3) and (3.4).

D Monte Carlo simulation

In this section of the appendix we report on the results of a Monte Carlo experiment. We specify a CES model of team production with log-normal errors, where we take the network structure (i.e., the set \mathcal{K} in (5.3)) as given from the empirical data. We fix the true value of the substitution parameter to $\gamma_0 = 1$, the team size parameter to $\beta_0 = 1$, the log-error variance in teams of size 2 to $\sigma_0^2(2) = 1$, and the variance in teams of size 1 to $\sigma_0^2(1) = 1$. This data generating process is designed to approximate what we found on the empirical data.

We report results based on 300 simulations. In each simulated sample, we estimate the parameters using plug-in method-of-moments and the Neyman-orthogonalized method-of-moments estimates of degree $q = 1$ to $q = 6$. As we did in our empirical study, we compute sample-split preliminary estimates of the author fixed-effects based on all their sole-authored publications except for one, selected at random. However, in the simulation exercise we do not cross-fit the estimators, and simply choose a random selection of sole-authored publications for each author in each Monte Carlo run.

In Tables 3 and 4 we show the median, mean, 2.5% quantile, and 97.5% quantile of each estimate across simulations. Starting with the substitution parameter γ , we see that the plug-in estimator is severely biased, with median and mean biases of -50% (expressed in proportion of the true value). For this parameter, all Neyman-orthogonalized estimators are substantially less biased, with a median bias ranging between 1% and 6%, with the lowest bias achieved by the estimates orthogonalized to order 5 and 6. However, in some replications the orthogonalized estimates tend to have large values, which is reflected in a somewhat larger mean bias, close to 3%, and quantile bands that are not symmetric around the true value.

Turning next to the team size parameter β , we see that both the plug-in and first-order Neyman-orthogonalized estimators are biased, with a median and mean bias of 5%–6%. All orthogonalized estimates of order $q \geq 2$ are virtually unbiased, both for the mean and the median. Moreover, in this case the quantile bands are symmetric around the true

parameter value.

Shifting attention to the variance in teams of size 2, $\sigma^2(2)$, we see that both the plug-in and first-order Neyman-orthogonalized estimators are severely biased, with a median and mean bias of 16%–18%. All orthogonalized estimates of order $q \geq 2$ are virtually unbiased, both for the mean and the median, and the quantile bands are centered around the true parameter value.

Lastly, turning to the variance in teams of size 1, $\sigma^2(1)$, the plug-in estimator exhibits a large bias of 33%. First-order orthogonalization only decreases the bias slightly, to 27%. In contrast, the Neyman-orthogonalized estimators continue to show good performance. In particular, when $q \geq 4$ the estimates are virtually unbiased, and the quantile bands are symmetric around the true value.

E Restrictions independent of individual effects

Model (5.3) implies restrictions on parameters $\gamma_0, \beta_0, \sigma_0^2(1), \sigma_0^2(2)$ that do not depend on the author-specific effects η_{i0} .⁹ As an example, the model implies the following alternative expression for the team size parameter β_0 :

$$\beta_0 = \left(\frac{\mathbb{E}[Y_j^{\gamma_0} | s_j = 2]}{\mathbb{E}[Y_j^{\gamma_0} | s_j = 1]} \right)^{\frac{1}{\gamma_0}} \exp \left(\frac{1}{2} \gamma_0 [\sigma_0^2(1) - \sigma_0^2(2)] \right), \quad (5.1)$$

which does not involve the fixed-effects η_{i0} . Note that, if $\gamma_0 = 0$ and output is additive in worker inputs, then $\log \beta_0$ is simply the difference between average log-outputs in teams of size 2 and 1, respectively. As a check, in Figure 2 we report estimates of the left-hand side of (5.1), against the estimates of β_0 shown in Table 1, for various orders of orthogonalization. We see that the estimates of the two sides of (5.1) tend to agree with each other well irrespective of the orthogonalization order, with slightly closer alignment for estimates of order $q \geq 2$.

Model (5.3) also implies restrictions on γ_0 alone. To see this, let us write (5.3), within

⁹The analysis in this section was inspired by discussions with Bo Honoré.

Table 3: Monte Carlo simulation

Substitution γ (true value=1)				
	Median	Mean	2.5%	97.5%
Plug-in	0.5084	0.5100	0.4213	0.6013
$q = 1$	0.9895	0.9956	0.7317	1.3041
$q = 2$	1.0562	1.0813	0.7487	1.5680
$q = 3$	1.0353	1.0571	0.7265	1.5369
$q = 4$	1.0148	1.0364	0.7132	1.4881
$q = 5$	1.0091	1.0303	0.7109	1.4743
$q = 6$	1.0091	1.0287	0.7124	1.4841
Team size β (true value=1)				
	Median	Mean	2.5%	97.5%
Plug-in	1.0610	1.0614	1.0217	1.0962
$q = 1$	1.0457	1.0474	0.9883	1.1008
$q = 2$	1.0016	1.0012	0.9318	1.0616
$q = 3$	1.0007	0.9985	0.9247	1.0615
$q = 4$	1.0010	0.9990	0.9245	1.0605
$q = 5$	1.0014	0.9993	0.9217	1.0602
$q = 6$	1.0014	0.9994	0.9245	1.0600

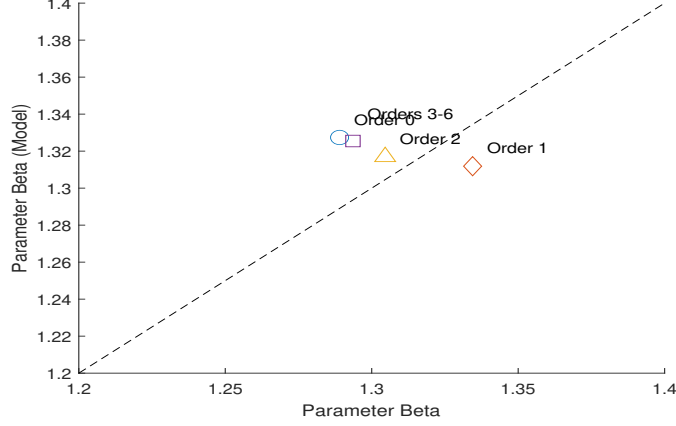
Notes: 300 simulations.

Table 4: Monte Carlo simulation (continued)

Variance $\sigma^2(2)$ (true value=1)				
	Median	Mean	2.5%	97.5%
Plug-in	1.1624	1.1621	1.1277	1.1960
$q = 1$	1.1838	1.1833	1.1456	1.2223
$q = 2$	0.9985	0.9977	0.9594	1.0371
$q = 3$	1.0017	1.0013	0.9610	1.0369
$q = 4$	1.0014	1.0009	0.9635	1.0382
$q = 5$	1.0016	1.0005	0.9622	1.0373
$q = 6$	1.0016	1.0005	0.9615	1.0371
Variance $\sigma^2(1)$ (true value=1)				
	Median	Mean	2.5%	97.5%
Plug-in	1.3337	1.3332	1.2862	1.3805
$q = 1$	1.2682	1.2702	1.2129	1.3285
$q = 2$	1.0212	1.0180	0.9429	1.0899
$q = 3$	1.0154	1.0128	0.9409	1.0910
$q = 4$	1.0052	1.0049	0.9272	1.0790
$q = 5$	1.0018	1.0022	0.9259	1.0759
$q = 6$	1.0014	1.0016	0.9248	1.0752

Notes: 300 simulations.

Figure 2: Comparing two estimates of β_0



Notes: Estimate of β_0 on the x-axis, model-based estimate of β_0 based on the right-hand side of (5.1) on the y-axis. Each point corresponds to an order of orthogonalization.

teams of size 2 only, as

$$Y_j^{\gamma_0} = \frac{1}{2} \beta_0^{\gamma_0} \left(\eta_{k(j,1)0}^{\gamma_0} + \eta_{k(j,2)0}^{\gamma_0} \right) \varepsilon_j^{\gamma_0 \sigma_0(2)},$$

which we write in vector form as

$$Y(\gamma_0) = A\tilde{\eta}_0 + \tilde{\varepsilon}, \quad (5.2)$$

where $Y(\gamma_0)$ has elements $Y_j^{\gamma_0}$, A is a matrix of zeros and ones, $\tilde{\eta}_{k0} = \frac{1}{2} \beta_0^{\gamma_0} \eta_{k0}^{\gamma_0} \exp\left(\frac{1}{2} \gamma_0^2 \sigma_0^2(2)\right)$, and $\tilde{\varepsilon}_j = \frac{1}{2} \beta_0^{\gamma_0} \left(\eta_{k(j,1)0}^{\gamma_0} + \eta_{k(j,2)0}^{\gamma_0} \right) \left[\varepsilon_j^{\gamma_0 \sigma_0(2)} - \exp\left(\frac{1}{2} \gamma_0^2 \sigma_0^2(2)\right) \right]$. Since $\mathbb{E}[\tilde{\varepsilon}_j | A] = 0$, (5.2) implies the conditional moment equalities

$$\mathbb{E}[(I - AA^\dagger)Y(\gamma_0) | A] = 0, \quad (5.3)$$

which only depend on γ_0 .

To use (5.3) for estimation, we rely on a set of instruments. For this purpose, we use interacted preliminary estimates $Z_j = \hat{\eta}_{k(j,1)} \hat{\eta}_{k(j,2)}$ for $k(j,1), k(j,2)$ the co-authors of j . Since we assume the preliminary estimates are constructed from an independent sample, we have

$$\mathbb{E}[Z'(I - AA^\dagger)Y(\gamma_0)] = 0. \quad (5.4)$$

Note these restrictions remain valid when ε_j are not Gaussian or not mutually independent, provided they are independent of A . We use GMM estimation based on (5.4), that is,

$$\hat{\gamma}^{\text{GMM}} = \underset{\gamma}{\operatorname{argmin}} |Z'(I - AA^\dagger)Y(\gamma)|. \quad (5.5)$$

We implement this estimator in the same way we have implemented our Neyman-orthogonalized equations. Specifically, we construct preliminary estimates of author effects using all but one sole-authored paper for each author, where we select the held-out sole-authored paper at random.

Using the same Monte Carlo simulation design as in Section D tends to give noisy estimates. For example, when the true value is $\gamma_0 = 1$, and $\sigma_0(1) = 1/5$ and $\sigma_0(2) = 1/5$, we obtain a mean GMM estimate of 1.0381, a median estimate of 0.9950, and a standard deviation of 0.2108 across 300 simulations. Moreover, out of the 300 simulations, in 23 cases we are unable to find another minimum in (5.5) other than $\gamma_2 = 0$. Note these findings correspond to a model with error variances that are *25 times* smaller than the variances we used for our main simulation design in Section D. This suggests this estimation approach, at least for this particular choice of instruments, is considerably less precise than our likelihood-based approach.

Lastly, computing the GMM estimator on the empirical data, cross-fitting 100 times, we obtain $\hat{\gamma}^{\text{GMM}} = 0.6110$. This is of a comparable magnitude to the estimates of γ reported in Table 1, when using a sufficiently high order of orthogonalization. However, it is worth noting that, out of the 100 random splits of the sole-authored productions, in 11 cases we are unable to find another minimum in (5.5) other than $\gamma = 0$, again reflecting the instability of this method in our setting.