

Generalized AKM: Theory and Evidence*

Francesco Del Prato[†] Yaroslav Korobka[‡] Paolo Zacchia[§]

September 2025

Preliminary draft. Please do not circulate it.

Abstract

This paper introduces an estimator for quadratic forms based on the linear parameters of a semi-parametric model. The leading example is the workhorse model of wage determination by Abowd, Kramarz, and Margolis (1999, AKM): our estimator targets standard variance components while allowing for a nonparametric treatment of both worker- and firm-level observable characteristics. We propose a bias-corrected estimator robust to heteroskedasticity that controls for approximating functions of the covariates. We show that this estimator is asymptotically unbiased and consistent when the number of linear parameters (e.g. the AKM fixed effects) is proportional to the sample size. In particular, consistency hinges on a strengthened smoothness condition (which we discuss for the first time) on the nonparametric component's functional class. In an empirical application, we show that adding a rich set of controls to the standard AKM model yields implausibly large firm effects. Our method addresses this issue, yielding estimates of variance components that are more robust relative to conventional approaches. Confounding—not functional-form choice—drives the standard model's instability.

Keywords: wage decomposition, many regressors, semiparametric model, series estimator

*We are grateful to Stanislav Anatolyev, Matias D. Cattaneo, Michal Kolesár, Mikkel Plagborg-Møller, Ulrich Müller, and seminar participants at Princeton University, for fruitful comments and suggestions. This paper results from research funded under the umbrella of the ERC-CZ project No. LL2319. Funding from the European Union's Horizon 2020 research and innovation program under the Marie Skłodowska-Curie grant agreement No. 870245 is also gratefully acknowledged.

[†]Department of Economics and Business Economics, Aarhus University. e-mail: francesco.delprato@econ.au.dk.

[‡]CERGE-EI (Charles University & Czech Academy of Sciences). e-mail: Yaroslav.Korobka@cerge-ei.cz

[§]CERGE-EI (Charles University & Czech Academy of Sciences). e-mail: Paolo.Zacchia@cerge-ei.cz.

1. Introduction

Why do some firms pay more than others, and how much of wage dispersion reflects workers, firms, or sorting? With linked employer–employee data, the AKM framework (Abowd, Kramarz, and Margolis 1999) is the workhorse reduced-form application for answering this question. Yet, standard plug-in variance-component estimators rely on linear covariate effects and can be biased under heteroskedasticity or limited mobility (Andrews et al. 2008; Kline, Saggio, and Sølvsten 2020; Cattaneo, Jansson, and Newey 2018b; Jochmans 2022). Modern applications feature many covariates and high-dimensional fixed effects; theory and evidence point to nonlinear wage profiles and worker–firm complementarities (Bonhomme, Lamadon, and Manresa 2019), so linear additivity—and common low-order fixes like cubic age profiles—can be a strong and potentially misspecified restriction, especially when several covariates are nonlinear (Mincer 1974; Card et al. 2018). This raises a question: how should one flexibly partial out many observables in AKM, and—once that is done correctly—how much wage variance is attributed to workers, firms, and sorting?

This paper develops estimation and inference for variance components (quadratic forms) when covariates enter the model through an unspecified function. In particular, we extend the familiar plug-in approach by Kline, Saggio, and Sølvsten (2020) based on a leave-one-out, bias-corrected estimator that is robust to arbitrary heteroskedasticity and to designs with many fixed effects. Our estimator further controls for a rich set of basis terms meant to approximate the actual effect of observed covariates. We establish consistency as the number of basis terms and fixed effects grow, and a chi-square limit under fixed rank, under a strengthened smoothness condition. Furthermore, we show asymptotic normality in the growing rank case. Simulations show a negligible bias. Empirically, after flexibly controlling for rich worker and firm observables, the AKM variance components are stable and align with benchmark ranges in the literature (Bonhomme et al. 2023). While there is semiparametric theory for estimators of linear functionals with growing bases, we are not aware of a framework that delivers validity for variance-component estimators when covariates enter through an unspecified function. Our results fill this gap by giving conditions under which variance components are consistently estimated and by characterizing their large-sample distributions.

We develop asymptotic results tailored to modern designs. Consistency holds when both the number of basis terms used to approximate the covariate effect and the number of fixed effects grow with the sample. For fixed-rank quadratic forms (e.g., testing a few linear contrasts of fixed effects), the statistic has a chi-square limit; for growing rank (e.g., analysis of variance across many groups), a Gaussian approximation applies. A key message is that quadratic targets are more demanding than linear ones: to control the approximation bias from the flexible covariate effect, we require a stronger smoothness condition than is standard for slope functionals. Moreover, the estimator scales to large datasets. A random-projection implementation approximates the leverage and influence quantities that enter the correction with negligible loss of accuracy, making high-order

interactions and large panels computationally feasible.

Our work connects two strands of literature. On the “many regressors” side, we complement results on inference with many covariates and robust variance estimation (Cattaneo, Jansson, and Newey 2018b; Anatolyev 2012; Anatolyev and Sølvsten 2023), and connect to leave-one-out corrections and variance-component estimation in linear models (Kline, Saggio, and Sølvsten 2020; Jochmans 2022). On the semiparametric side, we build on series estimators with growing bases (Donald and Newey 1994; Cattaneo, Jansson, and Newey 2018a), showing that quadratic targets require stronger smoothness than slope functionals, and relate our asymptotics to many-instruments and small-bandwidth regimes (Hansen, Hausman, and Newey 2008; Cattaneo, Crump, and Jansson 2014). We also discuss an alternative kernel-based construction that fits naturally within the same leave-one-out logic (Cattaneo, Crump, and Jansson 2014).

We show the finite-sample behavior of our estimation through simulations. When the true covariate effect is nonlinear, the plug-in estimator exhibits non-negligible bias that can be amplified by heteroskedasticity and many fixed effects. The leave-one-out correction removes this bias across a range of designs, with accurate coverage for the proposed standard errors.

Finally, we revisit the AKM wage decomposition using linked employer–employee data with rich worker and firm observables. The allocation of variance matters because firm premia correlate with productivity, capital intensity, and downsizing or outsourcing policies (Card, Heining, and Kline 2013; Goldschmidt and Schmieder 2017; Bertheau et al. 2023). Benchmarks in the literature place the firm component around 15–25%, while recent bias-corrected methods suggest 5–16% (Bonhomme et al. 2023). In our data, once covariates are flexibly partialled out and the quadratic-form plug-in bias is corrected, the variance components are stable across specifications and lie within the bias-corrected range. By contrast, a linear AKM with the same rich controls inflates the firm-effect variance. This contrast points to confounding from linear additivity—rather than economically meaningful nonlinearities in pay setting—as the source of instability. Allowing firm-specific, potentially nonlinear pay schedules in observed traits does not reallocate variance away from the additive worker and firm premia, indicating that such heterogeneity is second-order at the variance-decomposition margin in this sample.

Outline. The remainder of the paper is organized as follows. Section 2 introduces the semiparametric model and formalizes the variance component targets. Section 3 derives the estimator and its finite-sample representation. Section 4 establishes consistency and limiting distributions under fixed and growing ranks, and discusses the strengthened smoothness condition. Section 5 presents the random projection approximation. Section 6 derives the limiting distribution of our proposed estimator. Section 7 reports our simulation exercise, while Section 8 illustrates an empirical application on linked employer–employee data. Section 9 concludes.

2. Setup

We consider the following semiparametric model,

$$y_i = x_i' \beta + f(z_i) + e_i, \quad i = 1, \dots, n, \quad (1)$$

where the regressors $x_i \in \mathbb{R}^p$, $z_i \in \mathbb{R}^d$ are non-random. The unknown function $f(z)$ belongs to the class of smooth and real-valued functions, $f \in \mathcal{F}$. The unobserved errors $\{e_i\}_{i=1}^n$ are mutually independent and obey $\mathbb{E}[e_i] = 0$, but may possess observation-specific variances $\mathbb{E}[e_i^2] = \sigma_i^2$.

Our object of interest is a quadratic form $\theta := \beta' A \beta$ for some known non-random symmetric matrix $A \in \mathbb{R}^{p \times p}$ of rank r . This is a quantity of interest in a wide range of economic applications: two key examples are summarized next.

Example 1 (Generalized analysis of variance). We discuss a particular example of the generalized analysis of covariance model where our results allow for unspecified functional relationship between the exogenous covariates and the outcome variable. Consider having observations arranged into N groups with T_g observations in the g -th group. Since Fisher (1928) originally introduced the analysis of variance methodology, it has been common since to assume the following model:

$$y_{gt} = \alpha_g + z_{gt}' \delta + \varepsilon_{gt}, \quad g = 1, \dots, N, \quad t = 1, \dots, T_g,$$

where α_g are group-specific fixed effects, and z_{gt} is a vector of exogenous covariates. The focus here is the variability in the outcome variable attributable to groups, or

$$\sigma_\alpha^2 = \frac{1}{n} \sum_{g=1}^N T_g (\alpha_g - \bar{\alpha})^2$$

with $n := \sum_{g=1}^N T_g$, and $\bar{\alpha} := n^{-1} \sum_{g=1}^N T_g \alpha_g$. Our approach relaxes the linearity assumption on covariates and instead assumes any form of unspecified functional dependence, that is

$$y_{gt} = \alpha_g + f(z_{gt}) + \varepsilon_{gt}, \quad f \in \mathcal{F}, \quad g = 1, \dots, N, \quad t = 1, \dots, T_g.$$

We can represent it as in (1), defining $i := i(g, t)$ with $i(\cdot, \cdot)$ being a bijective function, $y_i := y_{gt}$, $z_i := z_{gt}$ and $e_i := \varepsilon_{gt}$,

$$x_i := d_i, \quad \beta := (\alpha_1, \dots, \alpha_N)', \quad d_i := (\mathbb{1}\{g=1\}, \dots, \mathbb{1}\{g=N\})'.$$

Our object of interest σ_α^2 can be represented here as $\beta' A \beta$, with

$$A := \begin{pmatrix} A_d' A_d & 0 \\ 0 & 0 \end{pmatrix}, \quad A_d := \frac{1}{\sqrt{n}}(d_1 - \bar{d}, \dots, d_n - \bar{d}), \quad \bar{d} := \frac{1}{n} \sum_{i=1}^n d_i.$$

Example 2 (Generalized AKM). Our second and leading example is a classic wage

decomposition model proposed in Abowd, Kramarz, and Margolis (1999). It models the log wage determination as an additive function of worker fixed effects, firm fixed effects, and a linear function of strictly exogenous covariates. Specifically,

$$y_{gt} = \alpha_g + \psi_{j(g,t)} + z'_{gt}\delta + \varepsilon_{gt}, \quad g = 1, \dots, N, \quad t = 1, \dots, T_g.$$

Here, α_g and $\psi_{j(g,t)}$ capture the g th worker and j th firm unobserved heterogeneity component respectively, and z_{gt} is a vector of exogenous regressors. The bijective function $j(\cdot, \cdot) : \{1, \dots, N\} \times \{1, \dots, \max_g T_g\} \rightarrow \{0, \dots, J\}$ maps $n = \sum_{g=1}^N T_g$ person-year observations to one of $J + 1$ firms. One of the model's objectives is to quantify how much of the variability in log wages is determined by firms,

$$\sigma_\psi^2 = \frac{1}{n} \sum_{g=1}^N \sum_{t=1}^{T_g} (\psi_{j(g,t)} - \bar{\psi})^2,$$

where $\bar{\psi} = \frac{1}{n} \sum_{g=1}^N \sum_{t=1}^{T_g} \psi_{j(g,t)}$. Given the empirical evidence that the relationship between the explanatory worker characteristics and the log wage is nonlinear (Mincer 1974; Card et al. 2018), the dominance of worker and firm fixed effects might partially stem from the functional linear restriction on the exogenous covariates. Our methodology allows for more plausible unspecified functional dependence between the outcome variable and regressors, that is,

$$y_{gt} = \alpha_g + \psi_{j(g,t)} + f(z_{gt}) + \varepsilon_{gt}, \quad f \in \mathcal{F}, \quad g = 1, \dots, N, \quad t = 1, \dots, T_g.$$

This formulation allows, for example, to capture rich patterns of non-linear interaction between both worker-level and firm-level characteristics, which possibly vary across observed groups in the population.

Given that the common covariates z_{gt} and the firm assignments $j(\cdot, \cdot)$ obey a strict exogeneity condition, we can rewrite the equation above as in (1) with

$$x_i := (d'_i, h'_i)', \quad \beta := (\alpha', \psi')', \quad \alpha := (\alpha_1, \dots, \alpha_N)' + \mathbb{1}'_N \psi_0, \quad \psi = (\psi_1, \dots, \psi_J)' - \mathbb{1}'_J \psi_0,$$

defining y_i , z_i , and e_i as in Example 1, and $h_i := (\mathbb{1}\{j(g, t) = 1\}, \dots, \mathbb{1}\{j(g, t) = J\})'$. The parameter of interest σ_ψ^2 then can be rewritten as $\beta' A_\psi \beta$ with

$$A_\psi := \begin{pmatrix} 0 & 0 & 0 \\ 0 & A'_h A_h & 0 \\ 0 & 0 & 0 \end{pmatrix}, \quad A_h := \frac{1}{\sqrt{n}}(h_1 - \bar{h}, \dots, h_n - \bar{h}), \quad \bar{h} := \frac{1}{n} \sum_{i=1}^n h_i.$$

3. Finite-sample properties

From now on, we restrict the analysis to model (1). To derive an estimator of β , one must regress y_i on x_i and functions of z_i . To this end, let $p^1(z), \dots, p^k(z)$ be some approximating

functions, and let $p_k(z) := (p^1(z), \dots, p^k(z)) \in \mathbb{R}^k$ be a vector of covariates $p_k(z_i)$ aimed at approximating $f(z_i)$. We assume that the class of functions \mathcal{F} to which f belongs can be well approximated by linear combinations of elements in $p_k(z_i)$.

Define M_{ij} as a (i, j) th element of $M := I_n - P_k(P'_k P_k)^{-1} P'_k$ with

$$P_k := (p_k(z_1), \dots, p_k(z_n)) \in \mathbb{R}^{n \times k},$$

and let the partialled-out design matrix $S_{xx} := \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i x_j'$ have a full column rank. Then, the estimator of β is defined as

$$\hat{\beta} := S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i y_j.$$

Under some regularity conditions, the bias of $\hat{\beta}$ is negligible for large k . Rewrite it as:

$$\begin{aligned} \hat{\beta} &= S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i (x_j' \beta + f(z_j) + e_j) \\ &= S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i x_j \beta' + S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i f(z_j) + S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i e_j \\ &= \beta + S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i f(z_j) + S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i e_j := \beta + B + U. \end{aligned}$$

Because of the zero-mean assumption on errors, $\mathbb{E}[U] = 0$ holds, and the bias in the estimator is reflected in B term only. Under some conditions discussed below,

$$B = o(1)$$

as the number of approximating functions goes to infinity, $k \rightarrow \infty$.

Analogously to the fully parametric setup examined by Kline, Saggio, and Sølvsten (2020), the plug-in estimator of the quadratic form $\hat{\theta}_{\text{PI}} := \hat{\beta}' A \hat{\beta}$ is biased, because

$$\mathbb{E}[\hat{\theta}_{\text{PI}} - \theta] = \text{trace}(A \text{var}[\hat{\beta}]) = \sum_{i=1}^n B_{ii} \sigma_i^2,$$

where $B_{ii} := \sum_{j=1}^n M_{ij} x_j' S_{xx}^{-1} A S_{xx}^{-1} \sum_{j=1}^n M_{ij} x_j$ measures the influence of the i th squared error e_i^2 on $\hat{\theta}_{\text{PI}}$. In cases when $p/n \approx 0$, individual elements B_{ii} are close to zero, and the bias is negligible. However, with increasing p , the bias is more pronounced and needs to be accounted for.

We introduce some notation to describe our estimator. Collect x_i in a matrix as $X := (x_1, \dots, x_n)' \in \mathbb{R}^{n \times p}$. Then, $W := (X, P_k) \in \mathbb{R}^{n \times p+k}$ is a matrix containing full set of regressors, with individual vectors w_i , and $W = (w_1, \dots, w_n)'$. Define

$$\hat{\gamma} := \left(\sum_{i=1}^n w_i w_i' \right)^{-1} \sum_{i=1}^n w_i y_i,$$

and the leave-one-out version of it,

$$\hat{\gamma}_{-i} := \left(\sum_{i=1}^n w_i w'_i - w_i w'_i \right)^{-1} \left(\sum_{i=1}^n w_i y_i - w_i y_i \right).$$

Also, note that first p elements of $\hat{\gamma}$ are $\hat{\beta}$, and the remaining k elements are basis coefficients from the function approximation. Similarly to Kline, Saggio, and Sølvsten (2020), we propose to estimate θ as

$$\hat{\theta} := \hat{\beta}' A \hat{\beta} - \sum_{i=1}^n B_{ii} \hat{\sigma}_i^2, \quad (2)$$

where $\hat{\sigma}_i^2$ is a leave-one-out estimator of the individual variance of e_i defined as

$$\hat{\sigma}_i^2 := y_i (y_i - w'_i \hat{\gamma}_{-i}). \quad (3)$$

We note that computing $\hat{\gamma}_{-i}$ for each $i = 1, \dots, n$ is computationally costly in large-scale applications. To avoid this, we represent (3) as

$$\hat{\sigma}_i^2 = \frac{y_i \hat{e}_i}{M_{W,ii}}, \quad (4)$$

where $M_{W,ii}$ is i th diagonal element of the matrix $M_W := I_n - W(W'W)^{-1}W'$ projecting onto the complement of the column space spanned by x_i and functions of z_i , and $\hat{e}_i := \sum_{j=1}^n M_{ij}(y_j - x_j \hat{\beta})$ are residuals.

Our estimator $\hat{\theta}$ is *exactly* unbiased if (i) $\hat{\beta}$ is unbiased, (ii) $\hat{\sigma}_i^2$ is unbiased for each i . For a moment, suppose that $\mathbb{E}[\hat{\beta}] = \beta$ holds and $f \in \mathcal{F}_{\text{linear}}$: that is, the estimator of the “slope” β is unbiased, and the unknown function $f(z) = z'\alpha$ is linear in z . Under these conditions, we have

$$\begin{aligned} \mathbb{E}[\hat{\sigma}_i^2] &= \mathbb{E}[y_i (y_i - x'_i \hat{\beta}_{-i} - p_k(z_i)' \hat{\alpha}_{-i})] \\ &= \mathbb{E}[(x'_i \beta + f(z_i) + e_i)(x'_i \beta + f(z_i) + e_i - x'_i \hat{\beta}_{-i} - p_k(z_i)' \hat{\alpha}_{-i})] \\ &= \mathbb{E}[(x'_i \beta + f(z_i) + e_i)(x'_i (\beta - \hat{\beta}_{-i}) + (f(z_i) - p_k(z_i)' \hat{\alpha}_{-i}) + e_i)] \\ &= (x'_i \beta x'_i + f(z_i) x'_i) \mathbb{E}[\beta - \hat{\beta}_{-i}] + (x'_i \beta + f(z_i)) \mathbb{E}[f(z_i) - p_k(z_i)' \hat{\alpha}_{-i}] + \sigma_i^2 \\ &= \sigma_i^2 \end{aligned}$$

where we use the fact that $\mathbb{E}[e_i x'_i (\beta - \hat{\beta}_{-i})] = \mathbb{E}[e_i] \mathbb{E}[x'_i (\beta - \hat{\beta}_{-i})] = 0$, because $\hat{\beta}_{-i}$ is independent of e_i , and similarly $\mathbb{E}[e_i (f(z_i) - p_k(z_i)' \hat{\alpha}_{-i})] = \mathbb{E}[e_i] \mathbb{E}[f(z_i) - p_k(z_i)' \hat{\alpha}_{-i}] = 0$, because $\hat{\alpha}_{-i}$ is independent of e_i .

However, for general classes of functions \mathcal{F} , the exact unbiasedness of $\hat{\sigma}_i^2$ does not hold. Instead, we rely on asymptotic approximations where we let the number of approximating functions to go to infinity, $k \rightarrow \infty$. Hence, the estimator of the slope $\hat{\beta}$, of the individual variances $\hat{\sigma}_i^2$ and, as a result, of the quadratic form $\hat{\theta}$ will be unbiased *asymptotically*¹, that

1. For expositional purposes, we sometimes drop the Euclidean norm on the vector $\|v\|$, and write asymptotic results as, say, $v = o_p(1)$ instead of $\|v\| = o_p(1)$.

is

$$\mathbb{E}[\hat{\beta}] - \beta = o(1), \quad \mathbb{E}[\hat{\sigma}_i^2] - \sigma_i^2 = o_p(1), \quad \mathbb{E}[\hat{\theta}] - \theta = o_p(1),$$

as $k \rightarrow \infty$.

For these purposes, we conventionally assume that the unknown function $f(z)$ belongs to the functional class \mathcal{F} of smooth and real-valued functions that can be well-approximated by the functions of their arguments.

Assumption 1. *We assume that $f \in \mathcal{F}$, where*

$$\mathcal{F} := \left\{ f : \min_{\alpha \in \mathbb{R}^k} \mathbb{E} [|f(z_i) - p_k(z_i)' \alpha|^2] \leq C k^{-2\alpha_f}, \alpha_f > 1 \right\} \quad (5)$$

for some absolute constant $C < \infty$.

This particular implies that

$$\sup_{f \in \mathcal{F}} \mathbb{E} [|f(z_i) - p_k(z_i)' \hat{\alpha}_{-i}|^2] = \mathcal{O}_p(k^{-2\alpha_f}),$$

so that all the functions in the class are bounded if z_i has a compact support.

The constant α_f plays an important role in our asymptotic analysis and highlights the point of departure from the results in the previous literature on semiparametric estimation. Specifically, to obtain the consistency result of our estimator, we need to strengthen the assumption on the constant α_f from being strictly positive (Donald and Newey 1994; Cattaneo, Jansson, and Newey 2018a) to being strictly greater than one. Intuitively, because we are interested in the estimator of the quadratic form, the bias should decrease to zero faster than in cases when we are solely interested in estimators that are linear in the outcome variable (for example, the estimator of β in fixed- p designs). This extra assumption is the cost of relaxing the linearity assumption on the covariates, and bears practical implications in the form of the number of exogenous variables, and assumptions on the degree of smoothness of $f(z)$ the researcher is willing to consider.²

For example, if the support of z_i is compact, while the basis functions (say, polynomials or splines), satisfy Assumption 1 with $\alpha_f = s_f/d$, where s_f is the number of continuous derivatives of $f(z)$. Given that $\alpha_f > 1$, it implies that such functions will satisfy the assumption if $s_f > d$, that is, if the function possesses more continuous derivatives than the dimensionality of z_i (Chen 2007). Practically speaking, it implies that if the researcher has relatively many exogenous covariates, she implicitly assumes that the underlying function is very smooth. We stress that there are studies (Armstrong and Kolesár 2018) that highlight the impossibility of the degree of smoothness estimation from the data, so our practical recommendation in empirical applications is to carefully look for an *ad hoc* appropriate balance between the number of covariates and robustness of the results.

2. We leave discussions on whether $\alpha_f > 1$ is a sufficient or a necessary condition for consistency for future research.

4. Consistency

In this section, we prove the consistency result for the proposed estimator $\hat{\theta}$. We study the asymptotic behavior of $\hat{\theta}$ assuming that x_i , z_i , and A are sequences of constants so that the only source of randomness is e_i . We adopt the conditional perspective approach as in Scheffe (1959), Searle, Casella, and McCulloch (2009), and Kline, Saggio, and Sølvsten (2020). This allows us to be agnostic about the potential dependency between x_i and z_i , and A . Here is also where our analysis is different from Cattaneo, Jansson, and Newey (2018b) where the authors consider sequences of random variables and conduct analysis conditional on z_i only. Limits are taken assuming the number of observations goes to infinity, $n \rightarrow \infty$, the number of approximating functions goes to infinity, $k \rightarrow \infty$ (so that the finite-sample bias of $\hat{\beta}$ and $\hat{\sigma}_i^2$ is asymptotically of negligible order), and the dimensionality of x_i goes to infinity, $p \rightarrow \infty$, to model the limited mobility bias. We make the following additional assumptions:

Assumption 2. (i) $\max_i(\mathbb{E}[e_i^4] + \sigma_i^{-2}) = \mathcal{O}(1)$; (ii) there exists a $c < 1$ such that $\max_i P_{W,ii} < c$ for all n , where $P_{W,ii} := I_n - M_{W,ii}$; (iii) $\max_i(x_i' \beta)^2 = \mathcal{O}(1)$.

Part (i) excludes heavy-tailed distributions of the errors as is typically assumed in the literature on OLS estimation (Cattaneo, Jansson, and Newey 2018b; Kline, Saggio, and Sølvsten 2020). Parts (ii) and (iii) coupled with the Assumption 1 imply that the leave-one-out estimator $\hat{\sigma}_i^2$ is well-defined and has bounded variance. Part (ii) also implies that $\frac{p+k}{n} \leq c < 1$ for all n . Below we formalize the consistency result for the proposed estimator.

Lemma 3. If Assumption 1 and 2 hold, A is positive semi-definite, $\theta = \beta' A \beta = \mathcal{O}(1)$, and $\text{trace}(\tilde{A}^2) = \sum_{\ell=1}^r \lambda_{\ell}^2 = o(1)$, then

$$\hat{\theta} - \theta \xrightarrow{p} 0.$$

4.1. Alternative estimation with kernel methods

In this subsection we briefly discuss an alternative way for estimating the quadratic form θ : one based on kernel methods. By still assuming that the data-generating process takes the form in (1), we now estimate the unknown function $f(z)$ at the point z as

$$\hat{f}(z) = \frac{\sum_{i=1}^n K\left(\frac{z-z_i}{h}\right) y_i}{\sum_{i=1}^n K\left(\frac{z-z_i}{h}\right)}, \quad \hat{f}_{-i}(z) = \frac{\sum_{j=1, j \neq i}^n K\left(\frac{z-z_j}{h}\right) y_j}{\sum_{j=1, j \neq i}^n K\left(\frac{z-z_j}{h}\right)}, \quad (6)$$

for some suitable kernel function $K(\cdot)$ that fulfills regularity requirements. This form of the "leave-one-out" kernel estimator is extensively used in the nonparametric literature on cross-validation or estimation of density-weighted average derivatives (Cattaneo, Crump, and Jansson, 2014). After obtaining the estimate of the function $\hat{f}_{-i}(z)$, we construct the

estimator of the slope as:

$$\hat{\beta}_a = \left(\sum_{i=1}^n x_i x'_i \right)^{-1} \sum_{i=1}^n x_i \hat{\xi}_i, \quad \hat{\xi}_i := y_i - \hat{f}(z_i),$$

and the estimator of the quadratic form as:

$$\hat{\theta}_a := \hat{\beta}'_a A \hat{\beta}_a - \sum_{i=1}^n B_{ii} \hat{\sigma}_{i,a}^2$$

with

$$\hat{\sigma}_{i,a}^2 := y_i \left(y_i - x'_i \hat{\beta}_{-i,a} - \hat{f}_{-i}(z_i) \right), \quad \hat{\beta}_{-i,a} := \left(\sum_{i=1}^n x_i x'_i - x_i x'_i \right)^{-1} \left(\sum_{i=1}^n x_i \hat{\xi}_i - x_i \hat{\xi}_i \right).$$

The leave-one-out kernel estimator (6) eliminates the "self-observation" bias when constructing residuals $\hat{\xi}_i$ in the first step. This form of sample-splitting is for example also used (and for similar reasons at that) in the literature on double-debiased machine learning (Chernozhukov et al. 2018). We conjecture that given conventional asymptotic conditions on the bandwidth and the sample size used in nonparametrics, namely $n \rightarrow \infty$, $h \rightarrow 0$, and $nh \rightarrow \infty$, the estimator $\hat{\theta}_a$ is consistent, in line with the Lemma 3.

5. Large-scale approximation

In this section we discuss an alternative estimator that allows for fast computation in typical large-scale applications such as those based on administrative linked employer-employee data. We follow Kline, Saggio, and Sølvsten (2020) by considering the random projection method of Achlioptas (2003). To describe it, fix $m \in \mathbb{N}$ and generate matrices $R_B, R_P \in \mathbb{R}^{m \times n}$ where each (i,j) coordinate is a random draw from the following Rademacher distribution:

$$R_{\cdot,ij} = \begin{cases} +1 & \text{with probability } 1/2, \\ -1 & \text{with probability } 1/2. \end{cases}$$

Decompose $A = 1/2(A'_1 A_2 + A'_2 A_1)$ for $A_1, A_2 \in \mathbb{R}^{n \times p}$, where $A_1 = A_2$ if A is positive semi-definite. Denote

$$\hat{P}_{W,ii} := \frac{1}{m} \left\| R_P W S_{ww}^{-1} w_i \right\|^2, \quad \hat{B}_{ii} := \frac{1}{m} \left(R_B A_1 S_{xx}^{-1} \sum_{j=1}^n M_{ij} x_j \right)' \left(R_B A_2 S_{xx}^{-1} \sum_{j=1}^n M_{ij} x_j \right).$$

The proposed estimator is then:

$$\hat{\theta}_{JLA} := \hat{\beta}' A \hat{\beta} - \sum_{i=1}^n \hat{B}_{ii} \hat{\sigma}_{i,JLA}^2, \quad \hat{\sigma}_{i,JLA}^2 := \frac{y_i(y_i - w'_i \hat{\gamma})}{1 - \hat{P}_{W,ii}} \left(1 - \frac{1}{m} \frac{3\hat{P}_{W,ii}^3 + \hat{P}_{W,ii}^2}{1 - \hat{P}_{W,ii}} \right). \quad (7)$$

As in Kline, Saggio, and Sølvsten (2020), the term $\frac{1}{m} \frac{3\hat{P}_{W,ii}^3 + \hat{P}_{W,ii}^2}{1 - \hat{P}_{W,ii}}$ removes a non-linearity bias from the approximation of $P_{W,ii}$ by $\hat{P}_{W,ii}$.

Lemma 4. *If Assumptions 1 and 2 are satisfied, $n/m^4 = o(1)$, $k \rightarrow \infty$, $\text{var}[\hat{\theta}]^{-1} = \mathcal{O}(n)$, and one of the following conditions hold, then $\text{var}[\hat{\theta}]^{-1/2}(\hat{\theta}_{JLA} - \hat{\theta} - B_m) = o_p(1)$ where $|B_m| \leq \frac{1}{m} \sum_{i=1}^n P_{W,ii}^2 |B_{ii}| \sigma_i^2$:*

- (i) *A is positive semi-definite and $\mathbb{E}[\hat{\beta}' A \hat{\beta}] - \theta = \sum_{i=1}^n B_{ii} \sigma_i^2 = \mathcal{O}(1)$.*
- (ii) *$A = 1/2(A'_1 A_2 + A'_2 A_1)$ where $\theta_1 = \beta' A'_1 A_1 \beta$ and $\theta_2 = \beta' A'_2 A_2 \beta$ satisfy (i) and $\frac{\text{var}[\hat{\theta}_1] \text{var}[\hat{\theta}_2]}{n \text{var}[\hat{\theta}]^2} = \mathcal{O}(1)$.*

6. Limiting distributions

This section examines two different results about the limiting distribution of $\hat{\theta}$. The first assumes that the rank of A : r , is fixed. Instead, the second result allows for growing r .

6.1. Fixed rank

First, we derive the limiting distribution of $\hat{\theta}$ for fixed r . We represent the estimator as

$$\hat{\theta} = \sum_{\ell=1}^r \lambda_\ell \left(\hat{b}_\ell^2 - \widehat{\text{var}}[\hat{b}_\ell] \right),$$

where $\hat{b} := \sum_{i=1}^n v_i y_i$, $\widehat{\text{var}}[\hat{b}] := \sum_{i=1}^n v_i v_i' \hat{b}_i^2$, and $v_i := Q' S_{xx}^{-1/2} \sum_{j=1}^n M_{ij} x_j$.

Below we state a formal result stating that the distribution of $\hat{\theta}$ is a sum of potentially dependent chi-squared random variables with non-centralities $b = (b_1, \dots, b_r)'$, and that the estimator of the variance of \hat{b} is consistent.

Theorem 5. *If Assumptions 1 and 2 hold, r is fixed, and $\max_i v_i' v_i = o(1)$, then*

- (i) $\text{var}[\hat{b}]^{-1/2}(\hat{b} - b) \xrightarrow{d} \mathcal{N}(0, I_r)$, where $b := Q' S_{xx}^{1/2} \beta$,
- (ii) $\text{var}[\hat{b}]^{-1} \widehat{\text{var}}[\hat{b}] \xrightarrow{p} I_r$,
- (iii) $\hat{\theta} = \sum_{\ell=1}^r \lambda_\ell \left(\hat{b}_\ell^2 - \text{var}[\hat{b}_\ell] \right) + o_p(\text{var}[\hat{\theta}]^{1/2})$.

This result is standard, though arguably of limited interest in empirical settings like the AKM model, where the rank of A grows with the sample size.

6.2. Growing rank

We can also motivate $\hat{\theta}$ using a generalized "full set" of regressors $\tilde{w}_i := \check{A} S_{ww}^{-1} w_i$, where $\check{A} := \begin{pmatrix} A & 0 \\ 0 & 0 \end{pmatrix}$ and $\gamma := (\beta \quad \alpha)'$. Indeed,

$$\theta = \gamma' \check{A} \gamma = \gamma' S_{ww} S_{ww}^{-1} \check{A} \gamma = \sum_{i=1}^n \gamma' w_i \tilde{w}_i' \gamma$$

and the asymptotically unbiased leave-one-out estimator is then

$$\hat{\theta} = \sum_{i=1}^n y_i \tilde{w}'_i \hat{\gamma}_{-i},$$

where the difference between $\gamma' w_i$ and its proxy y_i goes in expectation squared to zero as the number of approximating functions becomes large. Now we can represent

$$\hat{\theta} = \sum_{i=1}^n \sum_{\ell \neq i} C_{i\ell} y_i y_\ell,$$

where $C_{i\ell} := B_{W,i\ell} - 2^{-1} M_{W,i\ell} \left(B_{W,ii} M_{W,ii}^{-1} + B_{W,\ell\ell} M_{W,\ell\ell}^{-1} \right)$. Hence, the estimator of the quadratic form is a second-order U -statistic with the kernel $C_{i\ell}$. Below we show, using the supporting Lemma 7 provided in the Appendix, that the kernel varies with the sample size n in a way that the individual contributions are small enough leading to the Gaussian approximation. The analogical result is shown in Kline, Saggio, and Sølvsten (2020) in the context of quadratic form estimation under the linear specification, and Cattaneo, Jansson, and Newey (2018a) in the context of slope estimation in the semiparametric model.

To show normality of the proposed estimator, we decompose $\hat{\theta} - \theta$ into components that satisfy conditions of Lemma 7, and then rely on results of the Lemma to imply the joint normality of individual components. Compared to Kline, Saggio, and Sølvsten (2020), the difference between the estimator and the estimand now includes the error from the functional approximation, which is, however, of negligible order under the Assumption 1 and the assumption that the number of approximating functions becomes large.

Theorem 6. *If Assumption 1 and 2 hold, and the following conditions are satisfied,*

$$(i) \text{ var}[\hat{\theta}]^{-1} \max_i ((\tilde{w}'_i \gamma)^2 + (\check{w}'_i \gamma)^2) = o(1), \quad (ii) \frac{\lambda_1^2}{\sum_{\ell=1}^r \lambda_\ell^2} = o(1),$$

then $\text{var}[\hat{\theta}]^{-1/2}(\hat{\theta} - \theta) \xrightarrow{d} \mathcal{N}(0, 1)$.

7. Simulation study

We conduct a simulation study to verify the unbiasedness of the proposed estimator in different contexts. Our data-generating process follows (1) with the unknown function defined as $f(z_i) = \|z_i\|^7$. We generate $n = 500$ observations with $z_{\ell i} \sim \text{i.i.d. } \mathcal{U}(-1, 1)$, and $x_{\ell i} = \exp(0.3 \|z_i\|^2 + 0.7 \mathcal{N}(0, 1))$, fixing dimensionality of z_i to $d = 4$ and varying dimensionality of x_i as $p = \{20, 90, 180, 300\}$. The slope coefficients are generated as $\beta_j \sim \mathcal{U}(-0.5, 0.5)$ for $j = 1, \dots, p$, and we set the A matrix to $A = x_1 x'_1 + x_2 x'_2$. The explanatory power of x_i remains fixed to $R_x^2 = 0.7$ as we vary its dimensionality.

We evaluate the performance of the proposed estimator in (i) homoskedastic regime with errors generated as $e_i \sim \text{i.i.d. } \mathcal{N}(0, 1)$, (ii) heteroskedastic regime with errors gen-

erated as $e_i = \mathcal{N}(0, 1) + x_{1i}/12$, and different values of k for approximating $f(z_i)$, $k = \{4, 30, 50\}$. Table 1 displays the values of $|\frac{1}{s} \sum_{\ell=1}^s \hat{\theta}_\ell - \theta|$ for different settings, where $\hat{\theta}_\ell$ is an estimate of the quadratic form in ℓ th simulation as defined in (2), and the total number of simulations is $s = 1000$.

Table 1: Simulation study: results

$k \setminus p$	20	90	180	300
<i>Homoskedastic</i>				
4 (linear)	2.19	3.50	5.47	2.55
30 (3rd power + int.)	0.10	0.04	0.00	0.63
50 (5th power + int.)	0.06	0.00	0.10	0.51
<i>Heteroskedastic</i>				
4 (linear)	2.50	8.59	7.24	11.04
30 (3rd power + int.)	0.36	0.32	0.59	0.06
50 (5th power + int.)	0.08	0.01	0.22	0.02

Note: the table reports the bias of the quadratic form (in absolute value) for each experiment. The number of observations is $n = 500$, and the number of simulations is $s = 1000$. $k = 4$ corresponds to the linear model specification, $k = 30$ uses polynomials up to the third power and interactions up to the third power, and $k = 50$ uses polynomials up to the fifth power and interactions up to the fifth power.

In sum, given a high nonlinearity of the underlying function, we observe that a large number of approximating functions is required to eliminate the bias. A linear specification leads to a severe bias, while a cubic specification seems to not capture the nonlinearities fully, still leading to a bias (though of a much smaller magnitude). Approximation relying on polynomials and interactions of the fifth power appears to be the most robust solution, addressing estimation issues due to heteroskedasticity, high dimensionality of x_i , as well as the (unknown) nonparametric component of the model.

We highlight that approximations based on cubic polynomials are fairly common in practice: see e.g. Card et al. (2018). Our simulation study shows that in some situations this might not be sufficient to fully eliminate the bias, potentially leading to distorted final estimates in actual applications.

8. Empirical application

We apply our estimator to decompose wage dispersion in linked employer–employee data. The exercise illustrates that relaxing linearity in the covariates yields stable and interpretable variance components, in contrast to the standard linear AKM model once rich controls are introduced.

8.1. Data and estimation

Our empirical analysis combines two Portuguese sources of administrative data. The matched employer–employee dataset *Quadros de Pessoal* (QP) reports worker characteris-

tics (age, gender, occupation, qualification, education, contract type, hours, compensation) and firm characteristics (location, industry, revenues) for the universe of private-sector firms. The firm accounts dataset *Central de Balanços* (CB) provides annual balance-sheet and income-statement information (including total assets and sales) for non-financial corporations. We link QP and CB at the firm level to obtain a panel with rich worker- and firm-side observables.

We focus on the period 2014–2019 and on full-time workers aged 20–65. The model’s dependent variable is the logarithm of individual hourly wages. As in the typical “AKM” literature, we decompose the cross-sectional variance of wages into the variance of worker fixed effects (σ_α^2), the variance of firm fixed effects (σ_ψ^2), and their covariance $\text{cov}(\alpha, \psi)$, which captures sorting.

We compare two empirical specifications. The first is the standard linear two-way fixed effects (AKM) model with covariates entering additively and linearly. The second implements our generalized semiparametric estimator, which approximates $f(z)$ via polynomial bases of degree 1, 3, and 5 (denoted Generalized (1)/(3)/(5)). Importantly, all polynomial bases are then interacted with a fully saturated set of observed worker and firm categorical variables, allowing for a very rich set of group-specific non-linear effects. The Generalized (1) model differs from the standard AKM model as in the latter, the continuous covariates z are *not* interacted with any categorical indicators. We report estimates under a parsimonious control set (gender, age) and an extended set that adds education and qualification at the worker level and revenue and assets at the firm level.³

8.2. Results

Table 2 presents the estimates of the variance components for each model and set of controls.

Using the “parsimonious controls” set (Panel A), the standard AKM and our generalized estimator deliver similar, plausible magnitudes. The firm-effect variance lies around 0.13–0.14, consistent with ranges commonly found in the literature (Bonhomme et al. 2023). Increasing basis richness from degree 1 to degree 5 slightly reduces σ_ψ^2 and $\text{cov}(\alpha, \psi)$, in line with improved fit of $f(\cdot)$.

By contrast, introducing the extended controls (Panel B) leads the standard linear specification to break down: the estimate of σ_ψ^2 rises from roughly 0.14 to 1.116, an implausibly large value in log-variance units relative to Panel A. This result likely reflects severe confounding. Both the firm fixed effects and the firm-level observables (revenue, assets) serve as proxies for underlying firm productivity which is one of the main drivers of pay-setting policies. The linear structure of the standard AKM model cannot disentangle these collinear sources of variation, leading to a mechanical inflation of the firm effect variance. This can be interpreted as a failure of the model’s core assumption of additive separability, a restriction our method is designed to relax.

3. In practice, we compute all estimates by first partialling out covariates (which, in the “Generalized” cases, amount to the fully interacted set of polynomial bases P_k); subsequently, we compute the bias-corrected variance components of interest on the resulting residual.

Table 2: Wage variance decomposition

Model	σ_α^2	σ_ψ^2	$\text{cov}(\alpha, \psi)$
<i>Panel A: Parsimonious controls</i>			
Standard	0.3698	0.1404	0.0845
Generalized (1)	0.3655	0.1405	0.0838
Generalized (3)	0.3695	0.1351	0.0799
Generalized (5)	0.3697	0.1352	0.0799
<i>Panel B: Extended controls</i>			
Standard	0.1872	1.1157	0.0369
Generalized (1)	0.3635	0.1407	0.0843
Generalized (3)	0.3795	0.1343	0.0810
Generalized (5)	0.3660	0.1316	0.0781

Note: Estimates of the variance of worker effects, firm effects, and their covariance. The standard model imposes linearity on controls. The generalized models use our semiparametric estimator with polynomial approximations of order 1, 3, and 5. Panel A includes worker gender and age as controls. Panel B adds worker education and qualification, and firm revenue and assets.

In contrast, the generalized estimator is robust to this enrichment of controls. The variance components remain stable across panels and across degrees. For instance, in Generalized (5) the firm-effect variance is 0.135 under parsimonious controls and 0.132 with the extended set; the worker variance and the covariance display similar stability.

That these estimates are stable does not imply that the additional controls are irrelevant. Rather, it suggests that the high-dimensional worker and firm fixed effects have already absorbed most of the relevant variation featured in the covariates. For example, high-revenue firms likely already have systematically higher firm fixed effects, and workers with higher qualifications already possess higher worker fixed effects. The fixed effects capture the stable, time-invariant component of these observables, leaving little residual variance for the controls to explain once the person and firm effects are flexibly partialled out.

In summary, the empirical estimates suggest a vulnerability of the linear AKM approach in rich designs: collinearity between fixed effects and observables can severely distort the variance decomposition. Our semiparametric specification avoids rigid functional-form restrictions and yields components that are stable across specifications, pointing to person- and firm-specific premia as the primary drivers of wage dispersion, with residual nonlinearities in observed characteristics playing a second-order role once those premia are controlled for.

8.3. Implications for within-firm wage heterogeneity

Our specification includes saturated, high-order interactions between observed worker and firm characteristics. For any vector of observables z_{it} and groups that include firm identifiers, the model allows each firm to have its own (possibly nonlinear) pay schedule

$f_j(z_{it})$, in addition to worker and firm fixed effects α_i and ψ_j . In this setting, the stability of $(\sigma_\alpha^2, \sigma_\psi^2, \text{cov}(\alpha, \psi))$ across richer controls and higher-order bases shows that differences in returns to *observed* traits across firms contribute little to the variance decomposition once worker and firm fixed effects are included. Allowing firm-specific pay schedules in z does not shift variance away from the additive premia or change measured sorting in these data.

This conclusion should not be interpreted as the absence of within-firm heterogeneity *per se*. First, it speaks only to heterogeneity along observed z_{it} . Match-specific premia or complementarities outside z_{it} can still matter (e.g., a manager-worker fit that raises pay on a particular team). Second, firm-specific slopes are learned from variation in z_{it} within each firm. If a firm hires a narrow set of worker types, true slope differences are hard to identify and may appear as a higher firm effect rather than a change in the variance components. Third, our baseline interactions allow $f_j(\cdot)$ to differ across firms but not necessarily over time unless firm–time groups are included; time-varying policies can therefore exist without moving this static decomposition.

Overall, the evidence points to a wage policy that is well approximated by a firm-level shift applied broadly across workers, while firm-specific returns to observed traits—though allowed to be nonlinear and heterogeneous—are quantitatively second order in the variance decomposition for this sample.

9. Conclusions

This paper develops a semiparametric framework for estimating variance components, (as in the workhorse AKM model of wage decomposition) when covariates enter flexibly through an unknown function. We extend extant leave-one-out bias corrections to this setting and show that, under strengthened smoothness conditions, the resulting estimator is consistent and asymptotically normal. The approach remains computationally feasible in large-scale applications.

Simulations demonstrate that our proposed method eliminates the bias of plug-in estimators under heteroskedasticity and nonlinear covariate effects. In an empirical application to matched employer–employee data, we find that variance components are stable across specifications and lie within benchmark ranges reported in the recent literature. By contrast, linear AKM models with the same set of controls produce inflated firm effects. This evidence suggests that instability in standard decompositions is primarily due to functional-form restrictions rather than substantive nonlinearities in pay determination.

Relative to conventional methods, our estimator requires additional assumptions and comes with a slightly higher computational cost. In actual applications, these disadvantages must be weighed against the estimator’s ability to deliver results that are both more nuanced and different relative to standard approaches. To provide more guidance on this issue, in future work we plan on comparing the empirical implications of our method across a diverse set of administrative datasets, all of which shall ideally feature rich sets of observed covariates.

References

- Abowd, J., F. Kramarz, and D. Margolis. 1999. "High wage workers and high wage firms." *Econometrica* 67 (2): 251–333.
- Anatolyev, S. 2012. "Inference in regression models with many regressors." *Journal of Econometrics* 170 (2): 368–382.
- Anatolyev, S., and M. Sølvsten. 2023. "Testing many restrictions under heteroskedasticity." *Journal of Econometrics* 236 (1): 1–19.
- 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 (Statistics in Society)* 171 (3): 673–697.
- Armstrong, T., and M. Kolesár. 2018. "Optimal inference in a class of regression models." *Econometrica* 86 (2): 655–683.
- Bertheau, A., E. M. Acabbi, C. Barceló, A. Gulyas, S. Lombardi, and R. Saggio. 2023. "The unequal consequences of job loss across countries." *American Economic Review: Insights* 5, no. 3 (September): 393–408.
- Bonhomme, S., K. Holzheu, T. Lamadon, E. Manresa, M. Mogstad, and B. Setzler. 2023. "How much should we trust estimates of firm effects and worker sorting?" *Journal of Labor Economics* 41 (2): 291–322.
- Bonhomme, S., T. Lamadon, and E. Manresa. 2019. "A distributional framework for matched employer employee data." *Econometrica* 87 (3): 699–739.
- Card, D., A. R. Cardoso, J. Heining, and P. Kline. 2018. "Firms and labor market inequality: evidence and some theory." *Journal of Labor Economics* 36 (1): 13–70.
- Card, D., J. Heining, and P. Kline. 2013. "Workplace heterogeneity and the rise of West German wage inequality." *The Quarterly Journal of Economics* 128 (3): 967–1015.
- Cattaneo, M. D., R. K. Crump, and M. Jansson. 2014. "Small bandwidth asymptotics for density-weighted average derivatives." *Econometric Theory* 30 (1): 176–200.
- Cattaneo, M. D., M. Jansson, and W. K. Newey. 2018a. "Alternative asymptotics and the partially linear model with many regressors." *Econometric Theory* 34 (2): 277–301.
- . 2018b. "Inference in linear regression models with many covariates and heteroscedasticity." *Journal of the American Statistical Association* 113 (523): 1350–1361.
- Chen, X. 2007. "Large sample sieve estimation of semi-nonparametric models." *Handbook of Econometrics* 6:5549–5632.
- Donald, S. G., and W. K. Newey. 1994. "Series estimation of semilinear models." *Journal of Multivariate Analysis* 50 (1): 30–40.
- Fisher, Ronald. 1928. *Statistical Methods for Research Workers*. 5th ed. Oliver / Boyd.

- Goldschmidt, D., and J. Schmieder. 2017. "The rise of domestic outsourcing and the evolution of the German wage structure." *The Quarterly Journal of Economics* 132 (3): 1165–1217.
- Hansen, C., J. Hausman, and W. K. Newey. 2008. "Estimation with many instrumental variables." *Journal of Business & Economic Statistics* 26 (4): 398–422.
- Jochmans, K. 2022. "Heteroscedasticity-robust inference in linear regression models with many covariates." *Journal of the American Statistical Association* 117 (538): 887–896.
- Kline, P., R. Saggio, and M. Sølvsten. 2020. "Leave-out estimation of variance components." *Econometrica* 88 (5): 1859–1898.
- Mincer, Jacob. 1974. *Schooling, Experience, and Earnings*. National Bureau of Economic Research, Inc.

A. Proofs

Here, we provide proofs of the main results found in the paper and additional derivations.

Proof of Lemma 3. The variance of $\hat{\beta}$ is

$$\begin{aligned}\text{var}[\hat{\beta}] &= \text{var} \left[S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i f(z_j) + S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i e_j \right] \\ &= S_{xx}^{-1} \text{var} \left[\sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i f(z_j) + \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i e_j \right] S_{xx}^{-1} \\ &= S_{xx}^{-1} \left(\sum_{i=1}^n \sum_{j=1}^n M_{ij}^2 x_i x_i' \sigma_j^2 \right) S_{xx}^{-1}\end{aligned}$$

under the assumption of fixed x_i and z_i .

First, rewrite the difference between the estimator and the estimand as

$$\begin{aligned}\hat{\theta} - \theta &= \hat{\beta}' A \hat{\beta} - \beta' A \beta - \sum_{i=1}^n B_{ii} \hat{\sigma}_i^2 \\ &= \sum_{i=1}^n \sum_{j=1}^n M_{ij} y_j x_i' S_{xx}^{-1} A S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i y_j - \beta' A \beta - \sum_{i=1}^n B_{ii} \hat{\sigma}_i^2.\end{aligned}$$

Given a data-generating process for y_i , we can expand $\hat{\theta} - \theta$ further as

$$\begin{aligned}\hat{\theta} - \theta &= \sum_{i=1}^n \sum_{j=1}^n M_{ij} \left(x_j' \beta + f(z_j) + e_j \right) x_i' S_{xx}^{-1} A S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i \left(x_j' \beta + f(z_j) + e_j \right) - \beta' A \beta - \sum_{i=1}^n B_{ii} \hat{\sigma}_i^2 \\ &= \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_j' \beta x_i S_{xx}^{-1} A S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i x_j' \beta + \sum_{i=1}^n \sum_{j=1}^n M_{ij} f(z_j) x_i' S_{xx}^{-1} A S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i f(z_j) \\ &\quad + \sum_{i=1}^n \sum_{j=1}^n M_{ij} e_j x_i' S_{xx}^{-1} A S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i e_j + \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_j' \beta x_i' S_{xx}^{-1} A S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i f(z_j) \\ &\quad + \sum_{i=1}^n \sum_{j=1}^n M_{ij} f(z_j) x_i' S_{xx}^{-1} A S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i x_j' \beta + \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_j' \beta x_i' S_{xx}^{-1} A S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i e_j \\ &\quad + \sum_{i=1}^n \sum_{j=1}^n M_{ij} e_j x_i' S_{xx}^{-1} A S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i x_j' \beta + \sum_{i=1}^n \sum_{j=1}^n M_{ij} f(z_j) x_i' S_{xx}^{-1} A S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i e_j \\ &\quad + \sum_{i=1}^n \sum_{j=1}^n M_{ij} e_j x_i' S_{xx}^{-1} A S_{xx}^{-1} \sum_{i=1}^n \sum_{j=1}^n M_{ij} x_i f(z_j) - \beta' A \beta - \sum_{i=1}^n B_{ii} \hat{\sigma}_i^2.\end{aligned}$$

Next, we use the definition $B_{i\ell} := \sum_{j=1}^n M_{ij} x_j' S_{xx}^{-1} A S_{xx}^{-1} \sum_{j=1}^n M_{\ell j} x_j$, so that

$$\begin{aligned}\hat{\theta} - \theta &= \sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell} f(z_i)^2 + \sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell} e_i^2 - \sum_{i=1}^n B_{ii} \hat{\sigma}_i^2 \\ &\quad + 2 \sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell} x_\ell' \beta f(z_i) + 2 \sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell} x_\ell' \beta e_i + 2 \sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell} e_\ell f(z_i).\end{aligned}$$

Finally, rearranging, we have

$$\begin{aligned}\hat{\theta} - \theta &= \sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell} f(z_i)^2 + \sum_{i=1}^n \sum_{\ell \neq i} B_{i\ell} e_i e_\ell + \sum_{i=1}^n B_{ii} (e_i^2 - \hat{\sigma}_i^2) \\ &\quad + 2 \sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell} x'_\ell \beta f(z_i) + 2 \sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell} x'_\ell \beta e_i + 2 \sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell} e_\ell f(z_i).\end{aligned}\tag{8}$$

In general, it holds that

$$\mathbb{E}[|\hat{\theta} - \theta|^2] = |\mathbb{E}[\hat{\theta} - \theta]|^2 + \text{trace}(\text{var}[\hat{\theta} - \theta]).$$

To show that $\hat{\theta}$ is consistent for θ , we need to show that the bias and the variance of the difference in (8) goes to zero. The main idea is to compute bounds on each term and their variances, and show that these bounds are asymptotically negligible. Then, convergence in the quadratic mean would imply convergence in probability.

Applying expectations on both sides, and using independence and mean-zero properties of errors, we have

$$\mathbb{E}[\hat{\theta} - \theta] = \sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell} f(z_i)^2 + 2 \sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell} x'_\ell \beta f(z_i) + \sum_{i=1}^n B_{ii} \mathbb{E}[\sigma_i^2 - \hat{\sigma}_i^2].$$

We now prove that each term goes to zero in probability as $n \rightarrow \infty$, $k \rightarrow \infty$, and $p \rightarrow \infty$.

Denote $B := (B_{i\ell})_{i,\ell=1}^n \in \mathbb{R}^{n \times n}$, and $F := (f(z_1), \dots, f(z_n))' \in \mathbb{R}^n$. Then

$$\begin{aligned}\sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell} f(z_i)^2 &= F' B F \\ &= F' M X S_{xx}^{-1} A S_{xx}^{-1} X' M F \\ &= F' M X S_{xx}^{-1} A^{1/2} A^{1/2} S_{xx}^{-1} X' M F \\ &= (F' M X S_{xx}^{-1} A^{1/2})^2,\end{aligned}$$

where we use $A = A^{1/2} A^{1/2}$ because A is symmetric. By the Markov inequality, assumptions on \mathcal{F} , M being idempotent, and the Cauchy-Schwartz inequality,

$$\begin{aligned}\left\| \frac{1}{n} F' M X S_{xx}^{-1} A^{1/2} \right\| &\leq \text{trace} \left(\frac{1}{n} F' M F \right)^{1/2} \cdot \text{trace} \left(\frac{1}{n} A^{1/2} S_{xx}^{-1} X' M X S_{xx}^{-1} A^{1/2} \right)^{1/2} \\ &= \text{trace} \left(\frac{1}{n} F' M F \right)^{1/2} \cdot \text{trace} \left(\frac{1}{n} A^{1/2} S_{xx}^{-1} A^{1/2} \right)^{1/2} \\ &= \mathcal{O}(k^{-\alpha_f} \sqrt{p/n}).\end{aligned}$$

Thus, using the Assumption 1, we have

$$\sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell} f(z_i)^2 = \mathcal{O}(k^{-2\alpha_f} np) = \mathcal{O}(k^{-\alpha_f} n) \rightarrow 0.$$

Now, for the second term in the main decomposition,

$$\begin{aligned} \sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell} x'_\ell \beta f(z_i) &= F' BX \beta \\ &= F' MX S_{xx}^{-1} A S_{xx}^{-1} X' MX \beta \\ &= F' MX S_{xx}^{-1} A \beta. \end{aligned}$$

Given that $\beta' A S_{xx}^{-1} A \beta = \mathcal{O}(1)$, by the Markov inequality, assumption on \mathcal{F} , M being idempotent, and the Cauchy-Schwartz inequality we have

$$\begin{aligned} \left\| \frac{1}{n} F' MX S_{xx}^{-1} A \beta \right\| &\leq \text{trace} \left(\frac{1}{n} F' M F \right)^{1/2} \cdot \text{trace} \left(\frac{1}{n} \beta' A S_{xx}^{-1} X' M X S_{xx}^{-1} A \beta \right)^{1/2} \\ &= \text{trace} \left(\frac{1}{n} F' M F \right)^{1/2} \cdot \text{trace} \left(\frac{1}{n} \beta' A S_{xx}^{-1} A \beta \right)^{1/2} \\ &= \mathcal{O}(k^{-\alpha_f} / \sqrt{n}). \end{aligned}$$

So that using the Assumption 1, we have

$$\sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell} x'_\ell \beta f(z_i) = \mathcal{O}(k^{-\alpha_f} \sqrt{n}) \rightarrow 0.$$

To bound the third term, we should show that

$$\mathbb{E}[\hat{\sigma}_i^2 - \sigma_i^2] = (x'_i \beta + f(z_i)) (x'_i \mathbb{E}[\beta - \hat{\beta}_{-i}] + \mathbb{E}[f(z_i) - p_k(z_i)' \hat{\alpha}_{-i}])$$

goes in probability to zero. This is equivalent to bounding the bias of the estimator of the slope coefficient and the bias of the function approximation.

The bias of the slope estimator is

$$\mathbb{E}[\hat{\beta} - \beta] = \mathbb{E} \left[\left(\frac{1}{n} X' M X \right)^{-1} \frac{1}{n} X' M F \right].$$

By the Markov inequality, assumption on \mathcal{F} , M being idempotent, and the Cauchy-Schwartz inequality we have that:

$$\begin{aligned} \left\| \frac{1}{n} X' M F \right\| &\leq \text{trace} \left(\frac{1}{n} X' M X \right)^{1/2} \cdot \text{trace} \left(\frac{1}{n} F' M F \right)^{1/2} \\ &= \mathcal{O}(k^{-\alpha_f} \sqrt{p/n}), \end{aligned}$$

and, similarly,

$$\left\| \frac{1}{n} X' M X \right\| = \mathcal{O}(p/n),$$

so that using Assumption 1, we ultimately have:

$$\mathbb{E}[\|\hat{\beta} - \beta\|] = \mathcal{O}(\sqrt{n}/k^{\alpha_f} \sqrt{p}) \rightarrow 0.$$

Because of the assumption on the functional class \mathcal{F} , we can bound the bias of the function approximation using Jensen's inequality for convex $x \mapsto x^2$,

$$\begin{aligned}\mathbb{E}[|f(z_i) - p_k(z_i)' \hat{\alpha}_{-i}|] &\leq \mathbb{E}[|f(z_i) - p_k(z_i)' \hat{\alpha}_{-i}|^2]^{1/2} \\ &\leq (Ck^{-2\alpha_f})^{1/2} = \mathcal{O}_p(k^{-\alpha_f}),\end{aligned}$$

so that:

$$\mathbb{E}[|\hat{\sigma}_i^2 - \sigma_i^2|] = \mathcal{O}_p(\sqrt{n}/k^{\alpha_f} \sqrt{p}) \xrightarrow{p} 0.$$

From this it follows that:

$$\mathbb{E}[|\hat{\theta} - \theta|] \xrightarrow{p} 0.$$

We next turn our attention to bounding the variance. Let matrix

$$\tilde{A} := S_{xx}^{-1/2} A S_{xx}^{-1/2},$$

and let $\lambda_1, \dots, \lambda_r$ be its nonzero eigenvalues. We assume that $\lambda_1^2 \geq \dots \geq \lambda_r^2$, and that each eigenvalue appears as many times as its algebraic multiplicity. Under these assumptions, we can spectrally decompose $\tilde{A} = QDQ'$, where Q is a matrix of orthonormal vectors and $D = \text{diag}(\lambda_1, \dots, \lambda_r)$.

The variance of $2 \sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell} x_\ell' \beta e_i$ is

$$\begin{aligned}4 \sum_{i=1}^n \left(\sum_{\ell=1}^n B_{i\ell} x_\ell' \beta \right)^2 \sigma_i^2 &\leq \max_i \sigma_i^2 \beta' X' B^2 X \beta = \max_i \sigma_i^2 \beta' A S_{xx}^{-1} A \beta \\ &\leq \max_i \sigma_i^2 \lambda_1 \theta = o(1).\end{aligned}$$

To explain why the last inequality holds, define $\tilde{\beta} = S_{xx}^{1/2} \beta$ so that $\beta = S_{xx}^{-1/2} \tilde{\beta}$,

$$\beta' A S_{xx}^{-1} A \beta = \tilde{\beta}' S_{xx}^{-1/2} A S_{xx}^{-1/2} \tilde{\beta} = \tilde{\beta}' (S_{xx}^{-1/2} A S_{xx}^{-1/2})^2 \tilde{\beta} = \tilde{\beta}' \tilde{A}^2 \tilde{\beta},$$

and

$$\theta = \tilde{\beta}' S_{xx}^{-1/2} A S_{xx}^{-1/2} \tilde{\beta} = \tilde{\beta}' \tilde{A} \tilde{\beta}.$$

Then, using the Rayleigh quotient argument, it holds that

$$\frac{\tilde{\beta}' \tilde{A}^2 \tilde{\beta}}{\tilde{\beta}' \tilde{A} \tilde{\beta}} \leq \lambda_{\max}(\tilde{A}),$$

because \tilde{A} is positive semi-definite, $\tilde{\beta} \neq 0$, and $\lambda_{\max}(A)$ is the largest eigenvalue of some matrix A . From this it follows that we can bound

$$\beta' A S_{xx}^{-1} A \beta \leq \lambda_1 \theta,$$

use assumptions $\theta = \mathcal{O}(1)$, $\lambda_1 \leq \text{trace}(\tilde{A}^2)^{1/2} = o(1)$, and the fact that the variance is bounded. The conclusion then follows.

The variance of $2 \sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell} e_\ell f(z_i)$ is

$$4 \sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell}^2 \sigma_\ell^2 f(z_i)^2 \leq \max_i 4\sigma_i^2 f(z_i)^2 \sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell}^2 = \max_i 4\sigma_i^2 f(z_i)^2 \text{trace}(\tilde{A}^2) = o(1),$$

which follows since

$$\begin{aligned} \sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell}^2 &= \|B\|_F^2 = \text{trace}(B'B) \\ &= \text{trace}(MXS_{xx}^{-1}AS_{xx}^{-1}AS_{xx}^{-1}X'M) \\ &= \text{trace}(MXS_{xx}^{-1/2}S_{xx}^{-1/2}AS_{xx}^{-1/2}S_{xx}^{-1/2}AS_{xx}^{-1/2}S_{xx}^{-1/2}X'M) \\ &= \text{trace}(S_{xx}^{-1/2}AS_{xx}^{-1/2}S_{xx}^{-1/2}AS_{xx}^{-1/2}S_{xx}^{-1/2}X'MXS_{xx}^{-1/2}) \\ &= \text{trace}(S_{xx}^{-1/2}AS_{xx}^{-1/2}S_{xx}^{-1/2}AS_{xx}^{-1/2}) = \text{trace}(\tilde{A}^2), \end{aligned}$$

and we use the assumptions as above but now, instead of the bounded variance, we assume that \mathcal{F} is a class of bounded functions so that $\max |f(z_i)| < C$ for $i = 1, \dots, n$ for some absolute constant C (this is implied by the main assumption on \mathcal{F}).

Because $B_{i\ell}^2 = B_{\ell i}^2$ for any i, ℓ , the variance of $\sum_{i=1}^n \sum_{\ell \neq i} B_{i\ell} e_i e_\ell$ is

$$2 \sum_{i=1}^n \sum_{\ell \neq i} B_{i\ell}^2 \sigma_i^2 \sigma_\ell^2 \leq \max_i 2\sigma_i^4 \sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell}^2 = \max_i 2\sigma_i^4 \text{trace}(\tilde{A}^2) = o(1).$$

To compute the variance of $\sum_{i=1}^n B_{ii} (e_i^2 - \hat{\sigma}_i^2)$, we represent the leave-one-out variance estimator as

$$\begin{aligned} \hat{\sigma}_i^2 &= y_i(y_i - w'_i \hat{\gamma}_{-i}) = y_i M_{W,ii}^{-1} (y_i - w'_i \hat{\gamma}) \\ &= y_i M_{W,ii}^{-1} \hat{e}_i = y_i M_{W,ii}^{-1} \sum_{\ell=1}^n M_{W,i\ell} y_\ell \\ &= y_i M_{W,ii}^{-1} \sum_{\ell=1}^n M_{W,i\ell} (x'_\ell \beta + f(z_\ell) + e_\ell) \\ &= y_i M_{W,ii}^{-1} \sum_{\ell=1}^n M_{W,i\ell} (x'_\ell \beta + f(z_\ell)) + y_i M_{W,ii}^{-1} \sum_{\ell=1}^n M_{W,i\ell} e_\ell. \end{aligned}$$

Thus, the variance is now expressed as:

$$\begin{aligned} \sum_{i=1}^n \left(\sum_{\ell=1}^n M_{W,\ell\ell}^{-1} B_{\ell\ell} M_{W,i\ell} (x'_\ell \beta + f(z_\ell)) \right)^2 \sigma_i^2 &+ 2 \sum_{i=1}^n \sum_{\ell \neq i} M_{W,ii}^{-2} B_{ii}^2 M_{W,i\ell}^2 \sigma_i^2 \sigma_\ell^2 \\ &\leq \frac{1}{c^2} \max_i \sigma_i^2 \max_i (x'_i \beta + f(z_i))^2 \sum_{i=1}^n B_{ii}^2 + \frac{2}{c} \max_i \sigma_i^4 \sum_{i=1}^n B_{ii}^2 = o(1), \end{aligned}$$

because $\min_i M_{W,ii} \geq c > 0$, $\sum_{i=1}^n B_{ii}^2 \leq \text{trace}(\tilde{A}^2) = o(1)$, and $\max_i (x'_i \beta + f(z_i))^2 \leq 2 \max_i (x'_i \beta)^2 + 2 \max_i f(z_i)^2 = \mathcal{O}(1)$.

Because we have that

$$\mathbb{E}[\hat{\theta} - \theta] \xrightarrow{p} 0, \quad \text{var}[\hat{\theta} - \theta] \xrightarrow{p} 0,$$

the proposed estimator $\hat{\theta}$ is consistent. \square

Proof of Lemma 4. We prove the result by considering a second-order approximation of $\hat{\theta}_{\text{JLA}} - \hat{\theta}$ around $\hat{a}_i := (1 - P_{W,ii})^{-1}(\hat{P}_{W,ii} - P_{W,ii})$ as

$$(\hat{\theta}_{\text{JLA}} - \hat{\theta})_2 := \sum_{i=1}^n \hat{\sigma}_i^2 \left(B_{ii} - \hat{B}_{ii} - \hat{B}_{ii}\hat{a}_i - \hat{B}_{ii} \left(\hat{a}_i^2 - \frac{1}{m} \frac{3P_{W,ii}^3 + P_{W,ii}^2}{1 - P_{W,ii}} \right) \right),$$

and an approximation error that we show to be negligible,

$$\text{AE}_2 := \sum_{i=1}^n \hat{\sigma}_i^2 \hat{B}_{ii} \left(\frac{1}{m} \frac{3\hat{P}_{W,ii}^3 + \hat{P}_{W,ii}^2 - (3P_{W,ii}^3 + P_{W,ii}^2)(1 - \hat{a}_i)^2}{(1 - \hat{a}_i)^2(1 - P_{W,ii})} - \frac{\hat{a}_i^3}{1 - \hat{a}_i} \right).$$

To decompose $\hat{\theta}_{\text{JLA}} - \hat{\theta} = (\hat{\theta}_{\text{JLA}} - \hat{\theta})_2 + \text{AE}_2$, note that

$$\begin{aligned} \hat{\sigma}_{i,\text{JLA}}^2 &= \frac{y_i(y_i - w'_i \hat{y})}{1 - \hat{P}_{W,ii}} \left(1 - \frac{1}{m} \frac{3\hat{P}_{W,ii}^3 + \hat{P}_{W,ii}^2}{1 - \hat{P}_{W,ii}} \right) \\ &= \frac{1 - P_{W,ii}}{1 - \hat{P}_{W,ii}} \hat{\sigma}_i^2 \left(1 - \frac{1}{m} \frac{3\hat{P}_{W,ii}^3 + \hat{P}_{W,ii}^2}{1 - \hat{P}_{W,ii}} \right), \end{aligned}$$

so that

$$\begin{aligned} \hat{\theta}_{\text{JLA}} - \hat{\theta} &= \sum_{i=1}^n B_{ii} \hat{\sigma}_i^2 - \hat{B}_{ii} \hat{\sigma}_{i,\text{JLA}}^2 \\ &= \sum_{i=1}^n B_{ii} \hat{\sigma}_i^2 - \hat{B}_{ii} \hat{\sigma}_i^2 \left(1 - \frac{1}{m} \frac{3\hat{P}_{W,ii}^3 + \hat{P}_{W,ii}^2}{1 - \hat{P}_{W,ii}} \right) \\ &= \sum_{i=1}^n B_{ii} \hat{\sigma}_i^2 + \hat{B}_{ii} \hat{\sigma}_i^2 \left(\frac{1}{m} \frac{1 - P_{W,ii}}{1 - \hat{P}_{W,ii}} \frac{3\hat{P}_{W,ii}^3 + \hat{P}_{W,ii}^2}{1 - \hat{P}_{W,ii}} - \frac{1 - P_{W,ii}}{1 - \hat{P}_{W,ii}} \right). \end{aligned}$$

Add and subtract $m^{-1}\hat{\sigma}_i^2 \hat{B}_{ii}(1 - P_{W,ii})^{-1}(3P_{W,ii}^3 + P_{W,ii}^2)$ to obtain

$$\begin{aligned} \hat{\theta}_{\text{JLA}} - \hat{\theta} &= \sum_{i=1}^n B_{ii} \hat{\sigma}_i^2 + \frac{1}{m} \hat{\sigma}_i^2 \hat{B}_{ii} \frac{3P_{W,ii}^3 + P_{W,ii}^2}{1 - P_{W,ii}} + \hat{B}_{ii} \hat{\sigma}_i^2 \left(\frac{1}{m} \frac{1 - P_{W,ii}}{1 - \hat{P}_{W,ii}} \frac{3\hat{P}_{W,ii}^3 + \hat{P}_{W,ii}^2}{1 - \hat{P}_{W,ii}} - \frac{1 - P_{W,ii}}{1 - \hat{P}_{W,ii}} \right) \\ &\quad - \frac{1}{m} \hat{\sigma}_i^2 \hat{B}_{ii} \frac{3P_{W,ii}^3 + P_{W,ii}^2}{1 - P_{W,ii}} \\ &= \sum_{i=1}^n \hat{\sigma}_i^2 \left(B_{ii} + \frac{1}{m} \hat{\sigma}_i^2 \hat{B}_{ii} \frac{3P_{W,ii}^3 + P_{W,ii}^2}{1 - P_{W,ii}} \right) \\ &\quad + \sum_{i=1}^n \hat{\sigma}_i^2 \hat{B}_{ii} \left(\frac{1}{m} \frac{(1 - P_{W,ii})^2(3\hat{P}_{W,ii}^3 + \hat{P}_{W,ii}^2) - (1 - \hat{P}_{W,ii})^2(3P_{W,ii}^3 + P_{W,ii}^2)}{(1 - \hat{P}_{W,ii})^2(1 - P_{W,ii})} - \frac{1 - P_{W,ii}}{1 - \hat{P}_{W,ii}} \right), \end{aligned}$$

and using $(1 - \hat{P}_{W,ii})^{-1}(1 - P_{W,ii}) = (1 - \hat{a}_i)^{-1}$, and expanding up to the third order as

$$\begin{aligned}\frac{1 - P_{W,ii}}{1 - \hat{P}_{W,ii}} &= 1 + \frac{1 - P_{W,ii}}{1 - \hat{P}_{W,ii}} \hat{a}_i \\ &= 1 + \hat{a}_i + \hat{a}_i^2 + \frac{1 - P_{W,ii}}{1 - \hat{P}_{W,ii}} \hat{a}_i^3,\end{aligned}$$

we have that

$$\begin{aligned}\hat{\theta}_{\text{JLA}} - \hat{\theta} &= \sum_{i=1}^n \hat{\sigma}_i^2 \left(B_{ii} + \frac{1}{m} \hat{B}_{ii} \frac{3P_{W,ii}^3 + P_{W,ii}^2}{1 - P_{W,ii}} \right) \\ &\quad + \sum_{i=1}^n \hat{\sigma}_i^2 \hat{B}_{ii} \left(\frac{1}{m} \frac{3\hat{P}_{W,ii}^3 + \hat{P}_{W,ii}^2 - (1 - \hat{a}_i)^2(3P_{W,ii}^3 + P_{W,ii}^2)}{(1 - \hat{a}_i)^2(1 - P_{W,ii})} - \left(1 + \hat{a}_i + \hat{a}_i^2 + \frac{\hat{a}_i^3}{1 - \hat{a}_i} \right) \right) \\ &= (\hat{\theta}_{\text{JLA}} - \hat{\theta})_2 + \text{AE}_2.\end{aligned}$$

To describe the bias, we note that $\hat{P}_{W,ii}$, \hat{B}_{ii} , and $\hat{\sigma}_i^2$ are independent of each other, $\mathbb{E}[\hat{P}_{W,ii}] = P_{W,ii}$, $\mathbb{E}[\hat{B}_{ii}] = B_{ii}$, $\mathbb{E}[\hat{\sigma}_i^2] = \sigma_i^2 + \mathcal{O}_p(\sqrt{n}/k^{\alpha_f} \sqrt{p})$, and using properties of the Rademacher random variables,

$$\text{var}[\hat{a}_i] = \frac{2}{m} \frac{P_{W,ii} - \sum_{\ell=1}^n P_{W,i\ell}^4}{(1 - P_{W,ii})^2} = \frac{1}{m} \frac{3P_{W,ii}^3 + P_{W,ii}^2}{1 - P_{W,ii}} + \frac{P_{W,ii}(1 - P_{W,ii})^2 - 2 \sum_{\ell \neq i}^n P_{W,i\ell}^4}{m(1 - P_{W,ii})^2}.$$

Therefore, in total we have:

$$\mathbb{E}[(\hat{\theta}_{\text{JLA}} - \hat{\theta})_2] = - \sum_{i=1}^n \sigma_i^2 B_{ii} \left(\text{var}[\hat{a}_i] - \frac{1}{m} \frac{3P_{W,ii}^3 + P_{W,ii}^2}{1 - P_{W,ii}} \right) + \mathcal{O}_p(\sqrt{n}/k^{\alpha_f} \sqrt{p}),$$

or, assuming $k \rightarrow \infty$,

$$\mathbb{E}[(\hat{\theta}_{\text{JLA}} - \hat{\theta})_2] = B_m + o(1), \quad B_m := \sum_{i=1}^n B_{ii} \sigma_i^2 \left(\frac{2 \sum_{\ell \neq i}^n P_{W,i\ell}^4 - P_{W,ii}^2 (1 - P_{W,ii})^2}{m(1 - P_{W,ii})^2} \right).$$

Focusing on the variance next, denote $y := (y_1, \dots, y_n)'$, so that

$$\begin{aligned}\text{var} \left[\sum_{i=1}^n \hat{\sigma}_i^2 (B_{ii} - \hat{B}_{ii}) \right] &= \mathbb{E} \left[\text{var} \left[\sum_{i=1}^n \hat{\sigma}_i^2 \hat{B}_{ii} \right] \middle| y \right] + \text{var} \left[\mathbb{E} \left[\sum_{i=1}^n \hat{\sigma}_i^2 \hat{B}_{ii} \right] \middle| y \right] = \mathbb{E} \left[\text{var} \left[\sum_{i=1}^n \hat{\sigma}_i^2 \hat{B}_{ii} \right] \middle| y \right] \\ &\leq 2m^{-1} \sum_{i=1}^n \sum_{\ell=1}^n B_{i\ell}^2 \mathbb{E}[\hat{\sigma}_i^2 \hat{\sigma}_\ell^2] = \mathcal{O} \left(m^{-1} \text{trace}(\tilde{A}^2) \right), \\ \text{var} \left[\sum_{i=1}^n \hat{\sigma}_i^2 \hat{B}_{ii} \hat{a}_i \right] &= \mathbb{E} \left[\text{var} \left[\sum_{i=1}^n \hat{\sigma}_i^2 \hat{B}_{ii} \hat{a}_i \middle| y, R_B \right] \right] \leq 2m^{-1} \sum_{i=1}^n \sum_{\ell=1}^{\ell} P_{W,i\ell}^2 \frac{\mathbb{E}[\hat{B}_{ii} \hat{B}_{\ell\ell}] \mathbb{E}[\hat{\sigma}_i^2 \hat{\sigma}_\ell^2]}{(1 - P_{W,ii})(1 - P_{W,\ell\ell})} \\ &= \mathcal{O} \left(m^{-1} \text{trace}(\tilde{A}^2) + m^{-2} \text{trace}(\tilde{A}_1^2)^{1/2} \text{trace}(\tilde{A}_2^2)^{1/2} \right)\end{aligned}$$

for $\tilde{A}_k := S_{xx}^{-1/2} A'_k A_k S_{xx}^{-1/2}$ for $k = 1, 2$. Regarding the ensuing terms, it holds that:

$$\begin{aligned}\text{var}\left[\sum_{i=1}^n \hat{\sigma}_i^2 \hat{B}_{ii} (\hat{a}_i^2 - \text{var}[\hat{a}_i])\right] &= \sum_{i=1}^n \sum_{\ell=1}^n \mathbb{E}[\hat{B}_{ii} \hat{B}_{\ell\ell}] \mathbb{E}[\hat{\sigma}_i^2 \hat{\sigma}_\ell^2] \text{cov}[\hat{a}_i^2, \hat{a}_\ell^2] \\ &= \mathcal{O}(m^{-2} \text{trace}(\tilde{A}^2) + m^{-3} \text{trace}(\tilde{A}_1^2)^{1/2} \text{trace}(\tilde{A}_2^2)^{1/2}), \\ \text{var}\left[\sum_{i=1}^n \hat{\sigma}_i^2 (\hat{B}_{ii} - B_{ii}) \frac{2 \sum_{\ell \neq i}^n P_{W,i\ell}^4 - P_{W,ii}(1 - P_{W,ii})^2}{m(1 - P_{W,ii})^2}\right] &= \mathcal{O}(m^{-3} \text{trace}(\tilde{A}^2)), \\ \text{var}\left[\sum_{i=1}^n B_{ii} (\hat{\sigma}_i^2 - \sigma_i^2) \frac{2 \sum_{\ell \neq i}^n P_{W,i\ell}^4 - P_{W,ii}(1 - P_{W,ii})^2}{m(1 - P_{W,ii})^2}\right] &= \mathcal{O}(m^{-2} \text{var}[\hat{\theta}]).\end{aligned}$$

Because $\text{trace}(\tilde{A}^2) = \mathcal{O}(\text{var}[\hat{\theta}])$ and $m^{-4} \text{var}[\hat{\theta}]^{-2} \text{var}[\hat{\theta}_1] \text{var}[\hat{\theta}_2] = o(1)$, it can be ultimately established that $\text{var}[\hat{\theta}]^{-1/2} ((\hat{\theta}_{\text{JLA}} - \hat{\theta})_2 - B_m) = o_p(1)$.

Using that $\mathbb{E}[\hat{a}_i^3] = \mathcal{O}(m^{-2})$, $\mathbb{E}[\hat{a}_i^4] = \mathcal{O}(m^{-2})$, and $\max_i |\hat{a}_i| = o_p(\log n / \sqrt{m})$, the terms in the approximation error are as follows:

$$\begin{aligned}\sum_{i=1}^n \hat{\sigma}_i^2 \hat{B}_{ii} \hat{a}_i^3 + \sum_{i=1}^n \hat{\sigma}_i^2 \hat{B}_{ii} \hat{a}_i^4 &= m^{-2} \mathcal{O}_p (\mathbb{E}[\hat{\theta}_{1,\text{PI}} - \theta_1] + \mathbb{E}[\hat{\theta}_{2,\text{PI}} - \theta_2]), \\ \sum_{i=1}^n \hat{\sigma}_i^2 \hat{B}_{ii} \frac{\hat{a}_i^5}{1 - \hat{a}_i} &= \frac{\log n}{m^{5/4}} \mathcal{O}_p (\mathbb{E}[\hat{\theta}_{1,\text{PI}} - \theta_1] + \mathbb{E}[\hat{\theta}_{2,\text{PI}} - \theta_2]), \\ \frac{1}{m} \sum_{i=1}^n \hat{\sigma}_i^2 \hat{B}_{ii} \frac{3\hat{P}_{W,ii}^3 + \hat{P}_{W,ii}^2 - (3P_{W,ii}^3 + P_{W,ii}^2)(1 - \hat{a}_i)^2}{(1 - \hat{a}_i)^2(1 - P_{W,ii})} &\\ = \left(m^{-2} + \frac{\log n}{p^{5/4}}\right) \mathcal{O}_p (\mathbb{E}[\hat{\theta}_{1,\text{PI}} - \theta_1] + \mathbb{E}[\hat{\theta}_{2,\text{PI}} - \theta_2]).\end{aligned}$$

□

Proof of Theorem 5. Representation in the theorem holds because

$$\sum_{\ell=1}^r \lambda_\ell \hat{b}_\ell^2 = \hat{\beta}' S_{xx}^{1/2} Q D Q' S_{xx}^{1/2} \hat{\beta} = \hat{\beta}' S_{xx}^{1/2} S_{xx}^{-1/2} A S_{xx}^{-1/2} S_{xx}^{1/2} \hat{\beta} = \hat{\beta}' A \hat{\beta},$$

and

$$\sum_{i=1}^n B_{ii} \hat{\sigma}_i^2 = \text{trace}(A \widehat{\text{var}}[\hat{\beta}]) = \text{trace}(D \widehat{\text{var}}[\hat{b}]) = \sum_{\ell=1}^r \lambda_\ell \widehat{\text{var}}[\hat{b}_\ell].$$

We prove the theorem in three steps.

Approximation. Equivalently, represent $\hat{\theta}$ as

$$\hat{\theta} = \sum_{\ell=1}^r \lambda_\ell (\hat{b}_\ell^2 - \text{var}[\hat{b}_\ell]) + \sum_{i=1}^n B_{ii} (\sigma_i^2 - \hat{\sigma}_i^2),$$

and below we show that the second term is asymptotically dominated by the variance of the estimator, $\text{var}[\hat{\theta}]$. Given that, the asymptotic distribution of $\hat{\theta}$ is then driven by the joint distribution of random vector \hat{b} .

We need to show that the second term is asymptotically mean-zero and is of smaller order than variance of $\hat{\theta}$. The first claim is immediate because we have shown in Lemma 3 that

$$\mathbb{E}[\hat{\sigma}_i^2 - \sigma_i^2] = \mathcal{O}_p(\sqrt{n}/k^{\alpha_f} \sqrt{p}) \xrightarrow{p} 0,$$

so that

$$\mathbb{E} \left[\sum_{i=1}^n B_{ii} (\sigma_i^2 - \hat{\sigma}_i^2) \right] \xrightarrow{p} 0.$$

Now, rewrite

$$\begin{aligned} \sum_{i=1}^n B_{ii} (\hat{\sigma}_i^2 - \sigma_i^2) &= \sum_{i=1}^n B_{ii} M_{W,ii}^{-1} x'_i \beta \sum_{\ell=1}^n M_{W,i\ell} e_\ell + \sum_{i=1}^n (e_i^2 - \sigma_i^2) \\ &\quad + \sum_{i=1}^n B_{ii} M_{W,ii}^{-1} f(z_i) \sum_{\ell=1}^n M_{W,i\ell} e_\ell + \sum_{i=1}^n B_{ii} M_{W,ii}^{-1} \sum_{\ell \neq i} M_{W,i\ell} e_i e_\ell. \end{aligned}$$

The variances of the first and the third terms are

$$\sum_{\ell=1}^n \sigma_\ell^2 \left(\sum_{i=1}^n M_{W,i\ell} B_{ii} M_{W,ii}^{-1} x'_i \beta \right)^2 \leq \max_i \sigma_i^2 \sum_{i=1}^n B_{ii}^2 M_{W,ii}^{-2} (x'_i \beta)^2 \leq \max_i \sigma_i^2 \max_i (x'_i \beta)^2 M_{W,ii}^{-2} \sum_{i=1}^n B_{ii}^2,$$

$$\sum_{\ell=1}^n \sigma_\ell^2 \left(\sum_{i=1}^n M_{W,i\ell} B_{ii} M_{W,ii}^{-1} f(z_i) \right)^2 \leq \max_i \sigma_i^2 \sum_{i=1}^n B_{ii}^2 M_{W,ii}^{-2} f(z_i)^2 \leq \max_i \sigma_i^2 \max_i f(z_i)^2 M_{W,ii}^{-2} \sum_{i=1}^n B_{ii}^2,$$

and of the second and the forth

$$\sum_{i=1}^n B_{ii}^2 \text{var}[e_i^2] \leq \max_i \mathbb{E}[e_i^4] \sum_{i=1}^n B_{ii}^2,$$

$$\sum_{i=1}^n \sum_{\ell \neq i} (B_{ii}^2 M_{W,ii}^{-2} + B_{ii} M_{W,ii}^{-1} B_{\ell\ell} M_{W,\ell\ell}^{-1}) M_{i\ell}^2 \sigma_i^2 \sigma_\ell^2 \leq 2 \max_i \sigma_i^4 M_{W,ii}^{-2} \sum_{i=1}^n B_{ii}^2.$$

Because each variance is bounded by $C \sum_{i=1}^n B_{ii}^2$, to show that it is of smaller order than the variance of $\hat{\theta}$, we need $\text{var}[\hat{\theta}]^{-1} \sum_{i=1}^n B_{ii}^2 = o(1)$. It holds because

$$\text{var}[\hat{\theta}]^{-1} \sum_{i=1}^n B_{ii}^2 \leq \max_i v'_i v_i \text{var}[\hat{\theta}]^{-1} \sum_{\ell=1}^r \lambda_\ell^2 \leq \max_i v'_i v_i \max_i \sigma_i^{-4} = o(1).$$

Variance estimator consistency. Now we show that the variance estimator is consistent, i.e. $\text{var}[\hat{b}]^{-1} \widehat{\text{var}}[\hat{b}] \xrightarrow{p} I_r$. For it we need to show that

$$\text{var}[\vartheta' \hat{b}]^{-1} \left(\widehat{\text{var}}[\vartheta' \hat{b}] - \text{var}[\vartheta' \hat{b}] \right) = o_p(1), \quad \vartheta \in \mathbb{R}^r, \quad \vartheta' \vartheta = 1$$

for some non-random ϑ . Rewrite the expression above as

$$\delta(\vartheta) := \sum_{i=1}^n v_i(\vartheta) (\hat{\sigma}_i^2 - \sigma_i^2), \tag{9}$$

where

$$v_i(\vartheta) := \frac{(\vartheta' v_i)^2}{\sum_{i=1}^n \sigma_i^2 (\vartheta' v_i)^2}.$$

We know that $\mathbb{E}[\delta(\vartheta)] = o_p(1)$ because by the triangle inequality, the Cauchy-Schwartz inequality, and $|\mathbb{E}[\hat{\sigma}_i^2 - \sigma_i^2]| = o_p(1)$,

$$|v(\vartheta)\mathbb{E}[\hat{\sigma}_i^2 - \sigma_i^2]| \leq (v(\vartheta)^2)^{1/2} \cdot (|\mathbb{E}[\hat{\sigma}_i^2 - \sigma_i^2]|^2)^{1/2} = o_p(1)$$

for $i = 1, \dots, n$. The variance of $\delta(\vartheta)$ is

$$\begin{aligned} \sum_{i=1}^n \delta(\vartheta)^2 \text{var}[\hat{\sigma}_i^2] &\leq \sum_{i=1}^n v_i(\vartheta)^4 \\ &\leq \max_i \sigma_i^{-4} \max_i v_i v'_i \frac{\vartheta' \vartheta}{\sum_{i=1}^n v_i v'_i \vartheta' \vartheta} \\ &= \max_i \sigma_i^{-4} \max_i v_i v'_i = o(1) \end{aligned}$$

because $\max_i v_i v'_i = o(1)$ by assumption.

Asymptotic normality. Our objective is to prove that

$$\text{var}[\vartheta' \hat{b}]^{-1/2} \left(\vartheta' \hat{b} - \vartheta' b \right) \xrightarrow{d} \mathcal{N}(0, 1),$$

where $\hat{b} := Q' S_{xx}^{1/2} \hat{\beta}$, and $b := Q' S_{xx}^{1/2} \beta$. Lyapunov's condition implies that it is sufficient to show that

$$\text{var}[\vartheta' \hat{b}]^{-2} \sum_{i=1}^n \mathbb{E} \left[\left(\vartheta' (\hat{b} - b) \right)^4 \right] = o_p(1).$$

Because we have that

$$\begin{aligned} \vartheta' (\hat{b} - b) &= \vartheta' (Q' S_{xx}^{1/2} \hat{\beta} - Q' S_{xx}^{1/2} \beta) \\ &= \vartheta' \left(Q' S_{xx}^{1/2} \sum_{j=1}^n M_{ij} x_j f(z_i) + Q' S_{xx}^{1/2} \sum_{j=1}^n M_{ij} x_j e_i \right) \\ &= \vartheta' (v_i f(z_i) + v_i e_i), \end{aligned}$$

and $(a + b)^4 \leq C(a^4 + b^4)$ for some constant C , Lyapunov's condition is equivalent to

$$\text{var}[\vartheta' \hat{b}]^{-2} \sum_{i=1}^n C \left(f(z_i)^4 + \mathbb{E}[e_i^4] \right) \cdot (\vartheta' v_i)^4 = o_p(1).$$

It holds because $\max_i |f(z_i)| = \mathcal{O}(1)$, and $\max_i \mathbb{E}[e_i^4] = \mathcal{O}(1)$ by assumption, so that $\max_i f(z_i)^4 + \mathbb{E}[e_i^4] = \mathcal{O}(1)$, also $\max_i (\vartheta' v_i)^2 \leq \max_i v'_i v_i = o(1)$, $\sum_{i=1}^n (\vartheta' v_i)^2 = \sum_{i=1}^n \vartheta' v_i v'_i = 1$, and $\text{var}[\vartheta' \hat{b}]^{-2} \leq \max_i \sigma_i^{-2} = \mathcal{O}(1)$. \square

Proof of Theorem 6. We derive the limiting distribution of $\hat{\theta}$ with growing rank based on the following result regarding the joint normality of independent and not necessarily identical random variables as in Kline, Saggio, and Sølvsten (2020).

Let $\{q_{n,i}\}_{i,n}$ be a triangular array of row-wise independent random variables with $\mathbb{E}[q_{n,i}] = 0$ and $\text{var}[q_{n,i}] = \sigma_{n,i}^2$, let $\{\tilde{w}_{n,i}\}_{i,n}$ be a triangular array of non-random weights that satisfy $\sum_{i=1}^n \tilde{w}_{n,i} \sigma_{n,i}^2 = 1$ for $\forall n$, and let $(Q_n)_n$ be a sequence of symmetric non-random matrices in $\mathbb{R}^{n \times n}$ with zeros on the diagonal and having $2 \sum_{i=1}^n \sum_{\ell \neq i} Q_{n,i\ell}^2 \sigma_{n,i}^2 \sigma_{n,\ell}^2 = 1$. Define

$$\mathcal{S}_n := \sum_{i=1}^n \tilde{w}_{n,i} q_{n,i}, \quad \mathcal{U}_n := \sum_{i=1}^n \sum_{\ell \neq i} Q_{n,i\ell} q_{n,i} q_{n,\ell}.$$

Lemma 7. If $\max_i \mathbb{E}[q_{n,i}^4] + \sigma_{n,i}^{-2} = O(1)$, (i) $\max_i \tilde{w}_{n,i}^2 = o(1)$, and (ii) $\text{trace}(Q_n^4) = o(1)$, then $(\mathcal{S}_n, \mathcal{U}_n)' \xrightarrow{d} \mathcal{N}(0, I_2)$.

Proof. See Appendix B in Kline, Saggio, and Sølvsten (2020) and Appendix A2 in Solvsten (2020). \square

The U -statistic representation holds because

$$\begin{aligned} \hat{\theta} &= \sum_{i=1}^n y_i \tilde{w}'_i \hat{\gamma}_{-i} = \sum_{i=1}^n y_i \tilde{w}'_i \left(S_{ww}^{-1} - w_i w'_i \right) \sum_{\ell \neq i} w_\ell y_\ell \\ &= \sum_{i=1}^n y_i \tilde{w}'_i \left(S_{ww}^{-1} + \frac{S_{ww}^{-1} w_i w'_i S_{ww}^{-1}}{1 - w'_i S_{ww}^{-1} w_i} \right) \sum_{\ell \neq i} w_\ell y_\ell \\ &= \sum_{i=1}^n y_i \tilde{w}'_i \left(S_{ww}^{-1} + M_{W,ii}^{-1} S_{ww}^{-1} w_i w'_i S_{ww}^{-1} \right) \sum_{\ell \neq i} w_\ell y_\ell \\ &= \sum_{i=1}^n y_i w'_i S_{ww}^{-1} \check{A} S_{ww}^{-1} \sum_{\ell \neq i} w_\ell y_\ell + \sum_{i=1}^n y_i w'_i S_{ww}^{-1} \check{A} M_{W,ii}^{-1} S_{ww}^{-1} w_i w'_i S_{ww}^{-1} \sum_{\ell \neq i} w_\ell y_\ell \\ &= \sum_{i=1}^n \sum_{\ell \neq i} y_i y_\ell w'_i S_{ww}^{-1} \check{A} S_{ww}^{-1} w_\ell + \sum_{i=1}^n \sum_{\ell \neq i} y_i y_\ell M_{W,ii}^{-1} w'_i S_{ww}^{-1} \check{A} S_{ww}^{-1} w_i w'_i S_{ww}^{-1} w_\ell \\ &= \sum_{i=1}^n \sum_{\ell \neq i} y_i y_\ell B_{W,i\ell} + y_i y_\ell M_{W,ii}^{-1} B_{W,ii} (1 - M_{W,i\ell}) = \sum_{i=1}^n \sum_{\ell \neq i} C_{i\ell} y_i y_\ell. \end{aligned}$$

Now, define $\check{e}_i := f(z_i) - p_k(z_i)' \alpha$ to be an approximation error (it is implicitly indexed by the unknown function but we omit this dependence for brevity). We can write the difference

$$\begin{aligned} \hat{\theta} - \theta &= \sum_{i=1}^n \sum_{\ell \neq i} C_{i\ell} y_i y_\ell - \sum_{i=1}^n \gamma' w_i \tilde{w}'_i \gamma \\ &= \sum_{i=1}^n \sum_{\ell \neq i} C_{i\ell} (x'_i \beta + f(z_i) + e_i) (x'_\ell \beta + f(z_\ell) + e_\ell) - \sum_{i=1}^n \gamma' w_i \tilde{w}'_i \gamma \\ &= \sum_{i=1}^n \sum_{\ell \neq i} C_{i\ell} (\gamma' w_i + e_i + \check{e}_i) (\gamma' w_\ell + e_\ell + \check{e}_\ell) - \sum_{i=1}^n \gamma' w_i \tilde{w}'_i \gamma \\ &= \sum_{i=1}^n \sum_{\ell \neq i} C_{i\ell} \gamma' w_i \gamma' w_\ell + \sum_{i=1}^n \sum_{\ell \neq i} C_{i\ell} (\gamma' w_i e_\ell + \gamma' w_\ell e_i) - \sum_{i=1}^n \gamma' w_i \tilde{w}'_i \gamma \\ &\quad + \sum_{i=1}^n \sum_{\ell \neq i} C_{i\ell} (\gamma' w_i \check{e}_\ell + \gamma' w_\ell \check{e}_i + e_i \check{e}_\ell + e_\ell \check{e}_i) + \sum_{i=1}^n \sum_{\ell \neq i} C_{i\ell} e_i e_\ell + \sum_{i=1}^n \sum_{\ell \neq i} C_{i\ell} \check{e}_i \check{e}_\ell \end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^n e_i \sum_{\ell \neq i} (\gamma' w_i + \gamma' w_\ell) C_{i\ell} + \sum_{i=1}^n \check{e}_i \sum_{\ell \neq i} (\gamma' w_i + \gamma' w_\ell + e_i + e_\ell) \\
&\quad + \sum_{i=1}^n \sum_{\ell \neq i} C_{i\ell} e_i e_\ell + \sum_{i=1}^n \sum_{\ell \neq i} C_{i\ell} \check{e}_i \check{e}_\ell \\
&= \sum_{i=1}^n (2\tilde{w}'_i \gamma - \check{w}'_i \gamma) e_i + \sum_{i=1}^n \sum_{\ell \neq i} C_{i\ell} e_i e_\ell + \mathcal{O}_p(nk^{-\alpha_f}),
\end{aligned}$$

where the last equality follows from defining $\check{w}_i := \sum_{\ell=1}^n M_{W,i\ell} \frac{B_{W,\ell\ell}}{1-P_{W,\ell\ell}} w_\ell$, and the fact that by the Assumption 1

$$\check{e}_i^2 \leq (f(z_i) - p_k(z_i)' \alpha)^2 = \mathbb{E}[(f(z_i) - p_k(z_i)' \alpha)^2] \leq \min_{\alpha \in \mathbb{R}^k} \mathbb{E}[(f(z_i) - p_k(z_i)' \alpha)^2] \leq Ck^{-2\alpha_f}.$$

As $k \rightarrow \infty$, we have that

$$\hat{\theta} - \theta = \sum_{i=1}^n (2\tilde{w}'_i \gamma - \check{w}'_i \gamma) e_i + \sum_{i=1}^n \sum_{\ell \neq i} C_{i\ell} e_i e_\ell + o_p(1).$$

Having dispensed with asymptotically negligible contributions to $\hat{\theta}$, asymptotic variance is

$$\text{var}[\hat{\theta}] = \sum_{i=1}^n (2\tilde{w}'_i \gamma - \check{w}'_i \gamma)^2 \sigma_i^2 + 2 \sum_{i=1}^n \sum_{\ell \neq i} C_{i\ell} \sigma_i^2 \sigma_\ell^2,$$

that is, a sum of two components given by

$$\mathcal{V}_s := \sum_{i=1}^n (2\tilde{w}'_i \gamma - \check{w}'_i \gamma)^2 \sigma_i^2, \quad \mathcal{V}_u := 2 \sum_{i=1}^n \sum_{\ell \neq i} C_{i\ell} \sigma_i^2 \sigma_\ell^2.$$

The normalized difference is then given by

$$\text{var}[\hat{\theta}]^{-1/2}(\hat{\theta} - \theta) = \omega_1 \mathcal{S}_n + \omega_2 \mathcal{U}_n,$$

with $\omega_1 := \text{var}[\hat{\theta}]^{-1/2} \mathcal{V}_s^{1/2}$, $\omega_2 := \text{var}[\hat{\theta}]^{-1/2} \mathcal{V}_u^{1/2}$, and

$$\mathcal{S}_n := \mathcal{V}_s^{-1/2} \sum_{i=1}^n (2\tilde{w}'_i \gamma - \check{w}'_i \gamma) e_i, \quad \mathcal{U}_n := \mathcal{V}_u^{-1/2} \sum_{i=1}^n \sum_{\ell \neq i} C_{i\ell} e_i e_\ell.$$

Consider the case where the limit of ω_1 is nonzero. If it is not, then asymptotic normality of the difference $\text{var}[\hat{\theta}]^{-1/2}(\hat{\theta} - \theta)$ is implied by asymptotic normality of \mathcal{U}_n . Using notation of Lemma 7, we have that $\check{w}_i = \mathcal{V}_s^{-1/2}(2\tilde{w}'_i \gamma - \check{w}'_i \gamma)$, and $Q_{i\ell} = \mathcal{V}_u^{-1/2} C_{i\ell}$.

To verify the condition (i) of Lemma 7, note that

$$\begin{aligned}
\max_i \check{w}_i^2 &= \max_i \mathcal{V}_s (2\tilde{w}'_i \gamma - \check{w}'_i \gamma)^2 \leq \max_i 4\mathcal{V}_s ((\tilde{w}'_i \gamma)^2 + (\check{w}'_i \gamma)^2) \\
&= \max_i 4\omega_1^{-2} \frac{(\tilde{w}'_i \gamma)^2 + (\check{w}'_i \gamma)^2}{\text{var}[\hat{\theta}]} = o(1),
\end{aligned}$$

where the last equality follows from Theorem 6 (i), and the nonzero limit of ω_1 .

For the condition (ii) of Lemma 7, denote $\tilde{A}_W := S_{ww}^{-1/2} \check{A} S_{ww}^{-1/2}$, and note that the first r eigenvalues of the \tilde{A}_W matrix are equal to eigenvalues of the \check{A} matrix. Now, with constants c_U and c_L not dependent on n , we have that $\text{trace}(C^4) \leq c_U \cdot \text{trace}(B_W^4) = c_U \cdot \text{trace}(\tilde{A}_W^4) \leq c_U \lambda_1^2 \cdot \text{trace}(\tilde{A}_W^2)$, and $\mathcal{V}_u \geq c_L \min_i \sigma_i^4 \cdot \text{trace}(\tilde{A}_W)$, which implies

$$\text{trace}(Q^4) \leq \frac{c_U \lambda_1^2 \cdot \text{trace}(\tilde{A}_W^2)}{(c_L \min_i \sigma_i^4 \cdot \text{trace}(\tilde{A}_W^2))^2} = \mathcal{O}\left(\frac{\lambda_1^2}{\text{trace}(\tilde{A}_W^2)}\right) = o(1),$$

where the last equality follows from Theorem 6 (ii). \square

Sherman-Morrison representation. The representation in (4) holds due to application of the Sherman-Morrison formula to the leave-one-out estimator of γ . Denoting $S_{ww} := \sum_{i=1}^n w_i w_i'$, it holds that

$$(S_{ww} - w_i w_i')^{-1} = S_{ww}^{-1} + \frac{S_{ww}^{-1} w_i w_i' S_{ww}^{-1}}{1 - w_i' S_{ww}^{-1} w_i}.$$

Plugging the above into the definition of $\hat{\gamma}_{-i}$, we obtain

$$\begin{aligned} \hat{\gamma}_{-i} &= \left(S_{ww}^{-1} + \frac{S_{ww}^{-1} w_i w_i' S_{ww}^{-1}}{1 - w_i' S_{ww}^{-1} w_i} \right) \left(\sum_{i=1}^n w_i y_i - w_i y_i \right) \\ &= S_{ww}^{-1} \sum_{i=1}^n w_i y_i - S_{ww} w_i y_i + \frac{S_{ww}^{-1} w_i w_i' S_{ww}^{-1} \sum_{i=1}^n w_i y_i}{1 - w_i' S_{ww}^{-1} w_i} - \frac{S_{ww}^{-1} w_i w_i' S_{ww}^{-1} w_i y_i}{1 - w_i' S_{ww}^{-1} w_i} \\ &= \hat{\gamma} + S_{ww}^{-1} w_i \left(\frac{w_i' \hat{\gamma}}{1 - w_i' S_{ww}^{-1} w_i} - y_i - \frac{w_i' S_{ww}^{-1} w_i y_i}{1 - w_i' S_{ww}^{-1} w_i} \right) \\ &= \hat{\gamma} + S_{ww}^{-1} w_i \left(\frac{w_i' \hat{\gamma} - y_i}{1 - w_i' S_{ww}^{-1} w_i} \right) \\ &= \hat{\gamma} - S_{ww}^{-1} w_i \frac{(y_i - w_i' \hat{\gamma})}{1 - w_i' S_{ww}^{-1} w_i}. \end{aligned}$$

Substituting for $\hat{\gamma}_{-i}$ in (3), we have

$$\begin{aligned} \hat{\sigma}_i^2 &= y_i (y_i - w_i' \hat{\gamma}_{-i}) \\ &= y_i \left(y_i - w_i' \left(\hat{\gamma} - S_{ww}^{-1} w_i \frac{(y_i - w_i' \hat{\gamma})}{1 - w_i' S_{ww}^{-1} w_i} \right) \right) \\ &= y_i \left(y_i - w_i' \hat{\gamma} + w_i' S_{ww}^{-1} w_i \frac{(y_i - w_i' \hat{\gamma})}{1 - w_i' S_{ww}^{-1} w_i} \right) \\ &= y_i \left((y_i - w_i' \hat{\gamma}) \left(\frac{1}{1 - w_i' S_{ww}^{-1} w_i} \right) \right) \\ &= \frac{y_i \hat{e}_i}{M_{W,ii}}. \end{aligned}$$