

Robust Inference for Convex Pairwise Difference Estimators*

Matias D. Cattaneo[†]

Michael Jansson[‡]

Kenichi Nagasawa[§]

October 7, 2025

Abstract

This paper develops distribution theory and bootstrap-based inference methods for a broad class of convex pairwise difference estimators. These estimators minimize a kernel-weighted convex-in-parameter function over observation pairs that are similar in terms of certain covariates, where the similarity is governed by a localization (bandwidth) parameter. While classical results establish asymptotic normality under restrictive bandwidth conditions, we show that valid Gaussian and bootstrap-based inference remains possible under substantially weaker assumptions. First, we extend the theory of small bandwidth asymptotics to convex pairwise estimation settings, deriving robust Gaussian approximations even when a smaller than standard bandwidth is used. Second, we employ a debiasing procedure based on generalized jackknifing to enable inference with larger bandwidths, while preserving convexity of the objective function. Third, we construct a novel bootstrap method that adjusts for bandwidth-induced variance distortions, yielding valid inference across a wide range of bandwidth choices. Our proposed inference method enjoys demonstrable more robustness, while retaining the practical appeal of convex pairwise difference estimators.

Keywords: small bandwidth asymptotics, generalized jackknife, bootstrap, U-process, pairwise comparisons, robust distribution theory.

*This paper was prepared for the Econometric Theory Lecture delivered at the 2025 International Symposium on Econometric Theory and Applications (SETA), University of Macau (China), June 1–3, 2025. It was also presented at the Econometrics Journal Lecture of the 2024 (EC)² Conference (Amsterdam), and the 2025 Conference in Honor of Bo Honoré’s 65th Birthday (Princeton University). We thank the participants at these conferences for their feedback. Cattaneo gratefully acknowledges financial support from the National Science Foundation through grants SES-1947805, DMS-2210561, and SES-2241575. Jansson gratefully acknowledges financial support from the National Science Foundation through grant SES-1947662 and from the Aarhus Center for Econometrics (ACE) funded by the Danish National Research Foundation grant number DNRF186. Nagasawa gratefully acknowledges financial support from the British Academy through grant SRG24\241614.

[†]Department of Operations Research and Financial Engineering, Princeton University.

[‡]Department of Economics, UC Berkeley and ACE.

[§]Department of Economics, University of Warwick.

1 Introduction

Suppose $\mathbf{z}_1, \dots, \mathbf{z}_n$ is a random sample from the distribution of a random vector \mathbf{z} . This paper studies the large-sample properties of the following *convex* pairwise difference estimator:

$$\hat{\boldsymbol{\theta}}_n \in \arg \min_{\boldsymbol{\theta} \in \Theta} \binom{n}{2}^{-1} \sum_{i < j} m(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta}) K_{h_n}(\mathbf{w}_i - \mathbf{w}_j), \quad K_h(\mathbf{u}) = \frac{1}{h^d} K\left(\frac{\mathbf{u}}{h}\right), \quad (1.1)$$

where $\Theta \subseteq \mathbb{R}^k$ is a parameter space, $\sum_{i < j}$ denotes $\sum_{j=2}^n \sum_{i=1}^{j-1}$, $(\mathbf{z}_i, \mathbf{z}_j) \mapsto m(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta})$ is a permutation symmetric function, K is a symmetric, non-negative kernel, h_n is a positive bandwidth (or localization) parameter sequence, \mathbf{w} is a continuously distributed d -dimensional subvector of \mathbf{z} , and where $\boldsymbol{\theta} \mapsto m(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta})$ is a *convex* function. Pairwise difference estimation, which relies on local comparisons between observation pairs, has been used to address heterogeneity in nonlinear models. See [Powell \(1994\)](#), [Honoré and Powell \(2005\)](#), and [Aradillas-Lopez, Honoré, and Powell \(2007\)](#) for overviews, and Section 2 for three motivating examples.

In contrast to classical extremum estimators, $\hat{\boldsymbol{\theta}}_n$ is a local M -estimator that employs observation pairs (i, j) for which \mathbf{w}_i and \mathbf{w}_j are similar. The bandwidth h_n governs the degree of similarity: When $h_n \rightarrow 0$ (as $n \rightarrow \infty$), the estimator increasingly focuses on nearly identical-in- \mathbf{w} pairs. In turn, focusing on such pairs is natural in settings where identification can be based on the condition $\mathbf{w}_i \approx \mathbf{w}_j$ (combined with smoothness assumptions). The localization introduces a familiar trade-off for estimation and inference: A smaller h_n reduces bias from dissimilarity between \mathbf{w}_i and \mathbf{w}_j , but increases variance due to fewer available usable pairs. As a consequence, the large-sample behavior of $\hat{\boldsymbol{\theta}}_n$ depends critically on a delicate bias-variance trade-off determined by h_n . This paper develops novel inference methods for convex pairwise difference estimators that are demonstrably more robust to bandwidth choice than existing methods.

Under regularity conditions and assuming that

$$nh_n^d \rightarrow \infty \quad \text{and} \quad nh_n^4 \rightarrow 0,$$

the pairwise difference estimator is asymptotically linear:

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \boldsymbol{\xi}_0(\mathbf{z}_i) + o_{\mathbb{P}}(1) \rightsquigarrow \mathbf{N}(\mathbf{0}, \mathbb{E}[\boldsymbol{\xi}_0(\mathbf{z})\boldsymbol{\xi}_0(\mathbf{z})']), \quad (1.2)$$

where $\boldsymbol{\theta}_0$ is the estimand and $\boldsymbol{\xi}_0(\cdot)$ is the influence function (whose exact form is given below). The condition $nh_n^d \rightarrow \infty$ lower bounds the level of localization h_n allowed for, while the condition $nh_n^4 \rightarrow 0$ upper bounds the level of localization. The purpose of the latter condition is to control a smoothing bias term. The bias condition $nh_n^4 \rightarrow 0$ could be replaced by the weaker condition $nh_n^{2L} \rightarrow 0$ if a (higher-order) kernel of order $L > 2$ were used, but a higher-order kernel annihilates the convexity of the objective function because higher-order kernels take negative values.

The main results of this paper are obtained by combining three ideas:

1. *Small Bandwidth Asymptotics.* Utilizing the framework introduced by Cattaneo, Crump, and Jansson (2014a), we establish a more robust Gaussian distributional approximation for the pairwise difference estimator that allows for higher levels of localization by remaining valid even when the condition $nh_n^d \rightarrow \infty$ is violated. This generalized distributional approximation shows that, while the localization restriction $nh_n^d \rightarrow \infty$ is necessary for establishing asymptotic linearity, a Gaussian approximation can hold under the substantially weaker condition $n^2h_n^d \rightarrow \infty$, albeit with a convergence rate and large sample variance that depends explicitly on the level of localization used.
2. *Debiasing.* Following Honoré and Powell (2005) we debias the pairwise difference estimator using the method of *generalized jackknifing* introduced by Schucany and Sommers (1977). Doing so allows for (larger) bandwidths that violate the bias condition $nh_n^4 \rightarrow \infty$. This debiasing approach retains the convexity of the objective function, which is attractive for both theoretical (weaker regularity conditions) and practical (faster computation) reasons. The debiasing procedure combines linearly a collection of convex pairwise difference estimators constructed using different levels of localization. The resulting ensembling-based pairwise difference estimator admits a small bandwidth Gaussian approximation with an associated bias condition of the form $nh_n^{2L} \rightarrow 0$, where $L \geq 2$ denotes the order of a certain (equivalent) kernel induced by the debiasing procedure.
3. *Bootstrapping.* Building on insights in Cattaneo, Crump, and Jansson (2014b), we develop a valid bootstrap-based distributional approximation for the debiased pairwise difference estimator rescaling the localization parameter. The nonparametric bootstrap distributional approximation exhibits a mismatch in its asymptotic variance under small bandwidth asymptotics. The mismatch is characterized by a known multiplicative factor involving the localization parameter h_n . As a result, bootstrapping the (debiased) pairwise difference estimator with a different localization parameter (namely, $3^{1/d}h_n$ rather than h_n) leads to a valid bootstrap-based inference procedure also under small bandwidth asymptotics.

In combination, these three ideas enable us to offer a novel resampling-based inference method for (convex) pairwise difference estimators that are demonstrably more robust to a wider set of choices of the localization parameter h_n .

Our theoretical work is carefully developed to retain and leverage convexity of the objective function defining the pairwise difference estimator. This feature not only allows for fast implementation of the estimator and resampling-based methods, but also enables us to proceed under relatively weak conditions when obtaining theoretical results. When developing our theoretical results, we rely heavily on the foundational work of Hjort and Pollard (1993) and Pollard (1991), which we apply to the case of U -processes.

This paper is connected to several strands of the literature. Contributions to the pairwise difference estimation literature include Ahn and Powell (1993), Ahn, Ichimura, Powell, and Ruud (2018), Aradillas-Lopez (2012), Blundell and Powell (2004), Hong and Shum (2010), Honoré (1992), Honoré,

Kyriazidou, and Udry (1997), Honoré and Powell (1994), Jochmans (2013), and Kyriazidou (1997). The theoretical and practical features of small bandwidth asymptotics, and their connection with resampling methods for inference, are discussed in Cattaneo, Crump, and Jansson (2010), Cattaneo et al. (2014b), Cattaneo, Jansson, and Newey (2018), Cattaneo and Jansson (2018), Matsushita and Otsu (2021), Cattaneo and Jansson (2022), Cattaneo, Farrell, Jansson, and Masini (2025a), and references therein. The generalized jackknife has been successfully used for debiasing in density weighted average derivative estimation (Powell, Stock, and Stoker, 1989), asymptotically linear pairwise difference estimation (Honoré and Powell, 2005), nonlinear semiparametric estimation (Cattaneo, Crump, and Jansson, 2013), monotone estimation (Cattaneo, Jansson, and Nagasawa, 2024), and random forest estimation (Cattaneo, Klusowski, and Underwood, 2025b), among other settings. Shao and Tu (2012) give a textbook introduction to jackknifing, bootstrapping, and other resampling methods.

The rest of the paper proceeds as follows. Section 2 introduces the three motivating examples that are used throughout the paper to motivate our work and to illustrate the verification of the high-level assumptions imposed. Section 3 present our main theoretical distributional and bootstrap results for robust inference employing convex pairwise difference estimators. The proofs of these results are given in Section 4. Section 5 showcases how the high-level sufficient conditions imposed in our theoretical developments are verified for the three motivating examples. Section 6 gives final remarks.

2 Motivating Examples

We use three examples to motivate and illustrate our work. The first example involves an estimator that can be written in closed form (because it has a quadratic-in- θ function $m(\mathbf{z}_i, \mathbf{z}_j; \theta)$), while the other two examples do not. The second example has a smooth-in- θ function $m(\mathbf{z}_i, \mathbf{z}_j; \theta)$, while the third example does not. All three examples have convex-in- θ functions $m(\mathbf{z}_i, \mathbf{z}_j; \theta)$ and employ the following notation: $\mathbf{z}_i = (y_i, \mathbf{x}_i', \mathbf{w}_i')'$, with y_i a scalar outcome variable, \mathbf{x}_i a k -dimensional covariate, and \mathbf{w}_i a d -dimensional covariate. For more details on the examples, see Powell (1994), Honoré and Powell (2005), and Aradillas-Lopez et al. (2007).

2.1 Partially Linear Regression Model

The partially linear regression model studied here is of the form

$$y_i = \mathbf{x}_i' \theta_0 + \gamma_0(\mathbf{w}_i) + \varepsilon_i,$$

where θ_0 is the parameter of interest, $\gamma_0(\cdot)$ is an unknown function, and where $\mathbb{E}[\varepsilon_i | \mathbf{x}_i, \mathbf{w}_i] = 0$. Defining $\dot{y}_{i,j} = y_i - y_j$ and $\dot{\mathbf{x}}_{i,j} = \mathbf{x}_i - \mathbf{x}_j$, a pairwise difference estimator of θ_0 can be based on

$$m(\mathbf{z}_i, \mathbf{z}_j; \theta) = m_{\text{PLR}}(\mathbf{z}_i, \mathbf{z}_j; \theta) = \frac{1}{2}(\dot{y}_{i,j} - \dot{\mathbf{x}}_{i,j}' \theta)^2.$$

Setting $\Theta = \mathbb{R}^k$, the minimization problem defining the estimator admits a closed form solution (provided that a non-negative kernel function is used), namely

$$\hat{\theta}_n = \left(\sum_{i < j} \dot{\mathbf{x}}_{i,j} \dot{\mathbf{x}}'_{i,j} K_{h_n}(\mathbf{w}_i - \mathbf{w}_j) \right)^{-1} \sum_{i < j} \dot{\mathbf{x}}_{i,j} y_{i,j} K_{h_n}(\mathbf{w}_i - \mathbf{w}_j).$$

2.2 Partially Linear Logit Model

The partially linear logit model studied here is of the form

$$y_i = \mathbb{1}\{\mathbf{x}'_i \boldsymbol{\theta}_0 + \gamma_0(\mathbf{w}_i) + \varepsilon_i \geq 0\},$$

where $\boldsymbol{\theta}_0$ is the parameter of interest, $\gamma_0(\cdot)$ is an unknown function, and where

$$\mathbb{P}[\varepsilon_i \leq u | \mathbf{x}_i, \mathbf{w}_i] = \Lambda(u), \quad \Lambda(u) = \frac{\exp(u)}{1 + \exp(u)}.$$

The parameter $\boldsymbol{\theta}_0$ can be estimated using a pairwise difference estimator with $\Theta = \mathbb{R}^k$ and

$$m(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta}) = m_{\text{PLL}}(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta}) = -\mathbb{1}\{y_{i,j} \neq 0\} (y_i \ln \Lambda(\dot{\mathbf{x}}'_{i,j} \boldsymbol{\theta}) + y_j \ln \Lambda(-\dot{\mathbf{x}}'_{i,j} \boldsymbol{\theta})).$$

The minimization problem defining the estimator does not admit a closed form solution, but (provided that a non-negative kernel function is used) it is convex because $u \mapsto -\ln \Lambda(u)$ is.

2.3 Partially Linear Tobit Model

The partially linear censored regression model studied here is of the form

$$y_i = \max\{\mathbf{x}'_i \boldsymbol{\theta}_0 + \gamma_0(\mathbf{w}_i) + \varepsilon_i, 0\},$$

where $\boldsymbol{\theta}_0$ is the parameter of interest, $\gamma_0(\cdot)$ is an unknown function, $\mathbf{x}_i \perp \varepsilon_i | \mathbf{w}_i$, and where the conditional distribution of ε_i given \mathbf{w}_i admits a Lebesgue density. A pairwise difference estimator of $\boldsymbol{\theta}_0$ can be obtained by setting $\Theta = \mathbb{R}^k$ and employing

$$m(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta}) = m_{\text{PLT}}(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta}) = \tilde{m}_{\text{PLT}}(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta}) - \tilde{m}_{\text{PLT}}(\mathbf{z}_i, \mathbf{z}_j; \mathbf{0}),$$

where

$$\tilde{m}_{\text{PLT}}(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta}) = \begin{cases} |y_i| - (\dot{\mathbf{x}}'_{i,j} \boldsymbol{\theta} + y_j) \operatorname{sgn}(y_i) & \text{if } \dot{\mathbf{x}}'_{i,j} \boldsymbol{\theta} \leq -y_j \\ |\dot{y}_{i,j} - \dot{\mathbf{x}}'_{i,j} \boldsymbol{\theta}| & \text{if } -y_j < \dot{\mathbf{x}}'_{i,j} \boldsymbol{\theta} < y_i \\ |y_j| + (\dot{\mathbf{x}}'_{i,j} \boldsymbol{\theta} - y_i) \operatorname{sgn}(y_j) & \text{if } y_i \leq \dot{\mathbf{x}}'_{i,j} \boldsymbol{\theta} \end{cases}.$$

Because $\tilde{m}_{\text{PLT}}(\mathbf{z}_i, \mathbf{z}_j; \mathbf{0})$ does not depend on $\boldsymbol{\theta}$, the presence of $\tilde{m}_{\text{PLT}}(\mathbf{z}_i, \mathbf{z}_j; \mathbf{0})$ in $m_{\text{PLT}}(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta})$ does not affect the minimization problem defining the estimator. Nevertheless, it is theoretically attrac-

tive to work with m_{PLT} rather than \tilde{m}_{PLT} , as doing so allows for weaker regularity conditions for the existence of the expectation of the objective function.

For future reference, we note that m_{PLT} admits the alternative representation

$$m_{\text{PLT}}(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta}) = \begin{cases} |\dot{y}_{i,j} - \dot{\mathbf{x}}'_{i,j} \boldsymbol{\theta}| - |\dot{y}_{i,j}| & \text{if } y_i > 0, y_j > 0 \\ \max\{y_i - \dot{\mathbf{x}}'_{i,j} \boldsymbol{\theta}, 0\} - y_i & \text{if } y_i > 0, y_j = 0 \\ \max\{y_j + \dot{\mathbf{x}}'_{i,j} \boldsymbol{\theta}, 0\} - y_j & \text{if } y_i = 0, y_j > 0 \\ 0 & \text{if } y_i = 0, y_j = 0 \end{cases}.$$

The function $\boldsymbol{\theta} \mapsto m_{\text{PLT}}(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta})$ is convex and therefore so is the minimization problem defining the estimator (provided that a non-negative kernel function is used).

3 Distributional Approximation and Bootstrap Inference

As is standard in the literature, we generalize (1.1) slightly and define our estimator $\hat{\boldsymbol{\theta}}_n = \hat{\boldsymbol{\theta}}_n(h_n)$ to be any approximate minimizer of $\widehat{M}_n(\boldsymbol{\theta}; h_n)$, where

$$\widehat{M}_n(\boldsymbol{\theta}; h) = \binom{n}{2}^{-1} \sum_{i < j} m(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta}) K_h(\mathbf{w}_i - \mathbf{w}_j).$$

To be specific, we require

$$\widehat{M}_n(\hat{\boldsymbol{\theta}}_n(h); h) \leq \inf_{\boldsymbol{\theta} \in \Theta} \widehat{M}_n(\boldsymbol{\theta}; h) + o_{\mathbb{P}}(n^{-1}).$$

The objective function \widehat{M}_n is a sample counterpart of the function M given by

$$M(\boldsymbol{\theta}; h) = \mathbb{E}[\widehat{M}_n(\boldsymbol{\theta}; h)] = \mathbb{E}[m(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}) K_h(\mathbf{w}_1 - \mathbf{w}_2)].$$

Under regularity conditions, this function approximates, as $h \downarrow 0$, a function M_0 , which (does not depend on K and) admits a unique minimizer, namely the parameter of interest $\boldsymbol{\theta}_0$.

For the purposes of analyzing $\hat{\boldsymbol{\theta}}_n$ it is convenient to define $\boldsymbol{\theta}_n = \boldsymbol{\theta}(h_n)$, where

$$\boldsymbol{\theta}(h) \in \arg \min_{\boldsymbol{\theta} \in \Theta} M(\boldsymbol{\theta}; h)$$

is interpretable as a (fixed- h) “pseudo” parameter. With the help of $\boldsymbol{\theta}_n$ we can decompose the estimation error $\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0$ into a (non-stochastic) “bias” component $\boldsymbol{\theta}_n - \boldsymbol{\theta}_0$ and a “noise” component $\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_n$. Each component can be analyzed separately and in both cases the analysis will leverage convexity.

3.1 Regularity Conditions

The following assumption guarantees, among other things, that $\boldsymbol{\theta}_n$ is well defined for large n and that the bias component $\boldsymbol{\theta}_n - \boldsymbol{\theta}_0$ vanishes asymptotically; for details, see Lemma 1 of Section 4.1.

- Assumption 1.** (i) The kernel function K is a symmetric, bounded probability density.
- (ii) $\Theta \subseteq \mathbb{R}^k$ is convex, $(\mathbf{z}, \bar{\mathbf{z}}) \mapsto m(\mathbf{z}, \bar{\mathbf{z}}; \boldsymbol{\theta})$ is permutation symmetric, and $\boldsymbol{\theta} \mapsto m(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta})$ is convex with probability one.
- (iii) The distribution of \mathbf{w} admits a Lebesgue density $f_{\mathbf{w}}$, which is bounded and continuous on its support \mathcal{W} .
- (iv) For each $\boldsymbol{\theta} \in \Theta$,

$$\mathbb{E}[m(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta})] + \mathbb{E} \left[\sup_{\mathbf{v} \in \mathcal{W}} \mathbb{E}[m(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}) | \mathbf{w}_1 = \mathbf{w}, \mathbf{w}_2 = \mathbf{v}] f_{\mathbf{w}}(\mathbf{v}) \right] < \infty$$

and (with probability one)

$$\lim_{\mathbf{u} \rightarrow \mathbf{0}} \mathbb{E}[m(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}) | \mathbf{w}_1 = \mathbf{w}, \mathbf{w}_2 = \mathbf{w} + \mathbf{u}] = \mathbb{E}[m(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}) | \mathbf{w}_1 = \mathbf{w}, \mathbf{w}_2 = \mathbf{w}].$$

- (v) On Θ , the function M_0 given by

$$M_0(\boldsymbol{\theta}) = \int_{\mathcal{W}} \mathbb{E}[m(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}) | \mathbf{w}_1 = \mathbf{w}, \mathbf{w}_2 = \mathbf{w}] f_{\mathbf{w}}(\mathbf{w})^2 d\mathbf{w}$$

is uniquely minimized at an interior point $\boldsymbol{\theta}_0$.

The next assumption enables us to analyze the asymptotic properties of the noise component $\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_n$. To accommodate examples (such as the partially linear Tobit model) where $\boldsymbol{\theta} \mapsto m(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta})$ is not fully differentiable, we assume the existence of derivative-like functions $\mathbf{s}(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta}) \in \mathbb{R}^k$ and $\mathbf{H}(\mathbf{w}_i, \mathbf{w}_j; \boldsymbol{\theta}, \mathbf{t}) \in \mathbb{R}^{k \times k}$ such that, for any direction $\mathbf{t} \in \mathbb{R}^k$, the (remainder) terms

$$r_{\mathbf{t}}(\boldsymbol{\theta}, \tau) = \frac{m(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta} + \mathbf{t}\tau) - m(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta})}{\tau} - \mathbf{s}(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta})' \mathbf{t}$$

and

$$R_{\mathbf{t}}(\boldsymbol{\theta}, \tau) = \frac{\mathbb{E}[r_{\mathbf{t}}(\boldsymbol{\theta}, \tau) | \mathbf{w}_1, \mathbf{w}_2]}{\tau} - \frac{1}{2} \mathbf{t}' \mathbf{H}(\mathbf{w}_1, \mathbf{w}_2; \boldsymbol{\theta}, \mathbf{t}) \mathbf{t}$$

are suitably small for $\boldsymbol{\theta}$ near $\boldsymbol{\theta}_0$, $\tau > 0$ near zero, and $\mathbf{w}_1 \approx \mathbf{w}_2$. As further discussed below, functions \mathbf{s} and \mathbf{H} satisfying the following assumption exist (and are relatively easy to find) in each of our motivating examples.

Assumption 2. (i) For each $\mathbf{t} \in \mathbb{R}^k$, there is some $\delta > 0$ such that

$$\begin{aligned} \mathbb{E} \left[\sup_{\tau \in (0, \delta), \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| < \delta, \mathbf{w}_2 \in \mathcal{W}} |\mathbb{E}[r_{\mathbf{t}}(\boldsymbol{\theta}, \tau) | \mathbf{z}_1, \mathbf{w}_2]| f_{\mathbf{w}}(\mathbf{w}_2)^2 \right] &< \infty, \\ \mathbb{E} \left[\sup_{\tau \in (0, \delta), \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| < \delta, \mathbf{w}_2 \in \mathcal{W}} \mathbb{E}[r_{\mathbf{t}}(\boldsymbol{\theta}, \tau)^2 | \mathbf{w}_1, \mathbf{w}_2] f_{\mathbf{w}}(\mathbf{w}_2) \right] &< \infty, \\ \mathbb{E} \left[\sup_{\tau \in (0, \delta), \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| < \delta, \mathbf{w}_2 \in \mathcal{W}} |R_{\mathbf{t}}(\boldsymbol{\theta}, \tau)| f_{\mathbf{w}}(\mathbf{w}_2) \right] &< \infty, \end{aligned}$$

and (with probability one)

$$\begin{aligned} \lim_{\tau \downarrow 0, (\boldsymbol{\theta}, \mathbf{u}) \rightarrow (\boldsymbol{\theta}_0, \mathbf{0})} \mathbb{E}[r_{\mathbf{t}}(\boldsymbol{\theta}, \tau) | \mathbf{z}_1 = \mathbf{z}, \mathbf{w}_2 = \mathbf{w} + \mathbf{u}] &= 0, \\ \lim_{\tau \downarrow 0, (\boldsymbol{\theta}, \mathbf{u}) \rightarrow (\boldsymbol{\theta}_0, \mathbf{0})} \mathbb{E}[r_{\mathbf{t}}(\boldsymbol{\theta}, \tau)^2 | \mathbf{w}_1 = \mathbf{w}, \mathbf{w}_2 = \mathbf{w} + \mathbf{u}] &= 0, \\ \lim_{\tau \downarrow 0, (\boldsymbol{\theta}, \mathbf{u}) \rightarrow (\boldsymbol{\theta}_0, \mathbf{0})} \mathbb{E}[R_{\mathbf{t}}(\boldsymbol{\theta}, \tau) | \mathbf{w}_1 = \mathbf{w}, \mathbf{w}_2 = \mathbf{w} + \mathbf{u}] &= 0. \end{aligned}$$

(ii) There is some $\delta > 0$ and some function b with

$$\sup_{\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| < \delta} \|\mathbf{s}(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta})\| \leq b(\mathbf{z}_1)b(\mathbf{z}_2),$$

such that

$$\mathbb{E}[b(\mathbf{z})^4] + \sup_{\mathbf{w} \in \mathcal{W}} \mathbb{E}[b(\mathbf{z})^4 | \mathbf{w}] f_{\mathbf{w}}(\mathbf{w}) < \infty$$

and

$$\mathbb{E} \left[\sup_{\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| < \delta, \mathbf{w}_2 \in \mathcal{W}} \|\mathbf{H}(\mathbf{w}_1, \mathbf{w}_2; \boldsymbol{\theta}, \mathbf{t})\| f_{\mathbf{w}}(\mathbf{w}_2) \right] < \infty \quad \text{for each } \mathbf{t} \in \mathbb{R}^k.$$

(iii) There exist functions \mathbf{G}_0 , $\boldsymbol{\xi}_0$, and $\boldsymbol{\Xi}_0$ such that, for each $\mathbf{t} \in \mathbb{R}^k$ (and with probability one),

$$\begin{aligned} \lim_{(\boldsymbol{\theta}, \mathbf{u}) \rightarrow (\boldsymbol{\theta}_0, \mathbf{0})} \mathbf{H}(\mathbf{w}, \mathbf{w} + \mathbf{u}; \boldsymbol{\theta}, \mathbf{t}) f_{\mathbf{w}}(\mathbf{w}) &= \mathbf{G}_0(\mathbf{w}), \\ \lim_{(\boldsymbol{\theta}, \mathbf{u}) \rightarrow (\boldsymbol{\theta}_0, \mathbf{0})} 2\mathbb{E}[\mathbf{s}(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}) | \mathbf{z}_1 = \mathbf{z}, \mathbf{w}_2 = \mathbf{w} + \mathbf{u}] f_{\mathbf{w}}(\mathbf{w}) &= \boldsymbol{\xi}_0(\mathbf{z}), \end{aligned}$$

and

$$\lim_{(\boldsymbol{\theta}, \bar{\boldsymbol{\theta}}, \mathbf{u}) \rightarrow (\boldsymbol{\theta}_0, \boldsymbol{\theta}_0, \mathbf{0})} \mathbb{E}[\mathbf{s}(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}) \mathbf{s}(\mathbf{z}_1, \mathbf{z}_2; \bar{\boldsymbol{\theta}})' | \mathbf{w}_1 = \mathbf{w}, \mathbf{w}_2 = \mathbf{w} + \mathbf{u}] f_{\mathbf{w}}(\mathbf{w}) = \boldsymbol{\Xi}_0(\mathbf{w}).$$

(iv) $\boldsymbol{\Gamma}_0 = \mathbb{E}[\mathbf{G}_0(\mathbf{w})]$, $\boldsymbol{\Sigma}_0 = \mathbb{E}[\boldsymbol{\xi}_0(\mathbf{z})\boldsymbol{\xi}_0(\mathbf{z})']$, and $\mathbb{E}[\boldsymbol{\Xi}_0(\mathbf{w})]$ are positive definite.

3.2 Small Bandwidth Asymptotics

Defining

$$\mathbf{V}_n = \mathbf{V}_n(h_n) = \mathbf{\Gamma}_0^{-1} \left[n^{-1} \mathbf{\Sigma}_0 + \binom{n}{2}^{-1} h_n^{-d} \mathbf{\Delta}_0(K) \right] \mathbf{\Gamma}_0^{-1}, \quad \mathbf{\Delta}_0(K) = \mathbb{E}[\mathbf{\Xi}_0(\mathbf{w})] \int_{\mathbb{R}^d} K^2(\mathbf{u}) d\mathbf{u},$$

and letting Φ_k denote the distribution function of a k -dimensional standard Gaussian random vector, we have the following result.

Theorem 1. *Suppose Assumptions 1 and 2 hold. If $n^2 h_n^d \rightarrow \infty$ and if $h_n \rightarrow 0$, then*

$$\sup_{\mathbf{t} \in \mathbb{R}^k} \left| \mathbb{P} \left[\mathbf{V}_n^{-1/2} (\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_n) \leq \mathbf{t} \right] - \Phi_k(\mathbf{t}) \right| \rightarrow 0.$$

Under the assumptions of Theorem 1, the convergence rate of $\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_n$ equals

$$\rho_n = \sqrt{\min \left(n, \binom{n}{2} h_n^d \right)},$$

the magnitude of $\mathbf{V}_n^{-1/2}$. Provided that the bias is “small” in the sense that $\|\boldsymbol{\theta}_n - \boldsymbol{\theta}_0\| = o(\rho_n^{-1})$, Theorem 1 therefore encompasses the following three distinct large-sample regimes:

- *Asymptotic Linearity:* If $nh_n^d \rightarrow \infty$, then $\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)$ satisfies (1.2). In particular, $\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)$ converges in law to a mean-zero Gaussian distribution with asymptotic variance

$$\lim_{n \rightarrow \infty} n \mathbf{V}_n(h_n) = \mathbf{\Gamma}_0^{-1} \mathbf{\Sigma}_0 \mathbf{\Gamma}_0^{-1}.$$

- *Root-n Consistency without Asymptotic Linearity:* If $nh_n^d \rightarrow 2c \in (0, \infty)$, then $\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)$ is not asymptotically linear, but converges in law to a mean-zero Gaussian distribution with asymptotic variance

$$\lim_{n \rightarrow \infty} n \mathbf{V}_n(h_n) = \mathbf{\Gamma}_0^{-1} \left[\mathbf{\Sigma}_0 + \frac{1}{c} \mathbf{\Delta}_0(K) \right] \mathbf{\Gamma}_0^{-1}.$$

- *Slower than Root-n Consistency:* If $nh_n^d \rightarrow 0$ (but $n^2 h_n^d \rightarrow \infty$), then $\sqrt{n^2 h_n^d / 2}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)$ converges weakly to a mean-zero Gaussian distribution with asymptotic variance

$$\lim_{n \rightarrow \infty} \binom{n}{2} h_n^d \mathbf{V}_n(h_n) = \mathbf{\Gamma}_0^{-1} \mathbf{\Delta}_0(K) \mathbf{\Gamma}_0^{-1}.$$

The small bandwidth component (i.e., the term involving $\mathbf{\Delta}_0(K)$) in \mathbf{V}_n captures the additional uncertainty generated from increasing the localization of the observations pairs. Incorporating this component in the approximate variance is key to enabling us to replace the condition $nh_n^d \rightarrow \infty$ by the weaker condition $n^2 h_n^d \rightarrow \infty$ when obtaining a Gaussian approximation. As demonstrated by

Cattaneo et al. (2025a), incorporating the small bandwidth component can furthermore lead to a higher-order corrected distributional approximation even under asymptotic linearity.

3.3 Debiasing

In Theorem 1, we centered the estimator $\widehat{\boldsymbol{\theta}}_n = \widehat{\boldsymbol{\theta}}_n(h_n)$ at $\boldsymbol{\theta}_n = \boldsymbol{\theta}(h_n)$ to circumvent bias issues. This section focuses on the bias term $\boldsymbol{\theta}_n - \boldsymbol{\theta}_0$ and introduces an automatic debiasing approach under the assumption that $\boldsymbol{\theta}_n - \boldsymbol{\theta}_0$ can be expanded in even powers of h_n . To be specific, we follow Honoré and Powell (2005, Section 3.3) and discuss debiasing under the following high-level condition.

Assumption 3. For some even $L \geq 0$, $\boldsymbol{\theta}(\cdot)$ admits $\mathbf{b}_{2l} \in \mathbb{R}^k$ (for $l = 1, \dots, L/2$) such that

$$\boldsymbol{\theta}(h) - \boldsymbol{\theta}_0 = \sum_{l=1}^{L/2} \mathbf{b}_{2l} h^{2l} + o(h^L) \quad \text{as } h \downarrow 0.$$

The ease with which Assumption 3 can be verified depends on the magnitude of L . For instance, Assumption 1 implies that Assumption 3 holds with $L = 0$. Under additional smoothness conditions and using symmetry of K , the following result gives conditions under which Assumption 3 holds with $L = 2$. When stating the result, we employ standard multi-index notation: for $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_d)' \in \mathbb{Z}_+^d$, $\mathbf{v} = (v_1, \dots, v_d)' \in \mathbb{R}^d$, and a sufficiently smooth-in- \mathbf{v} function $f(\mathbf{w}, \mathbf{v})$,

$$\partial_{\mathbf{v}}^{\boldsymbol{\alpha}} f(\mathbf{w}, \mathbf{v}) = \frac{\partial^{|\boldsymbol{\alpha}|}}{\partial v_1^{\alpha_1} \dots \partial v_d^{\alpha_d}} f(\mathbf{w}, \mathbf{v}), \quad |\boldsymbol{\alpha}| = \sum_{j=1}^d \alpha_j.$$

Proposition 1. Suppose Assumptions 1-2 hold and that

- (i) $\int_{\mathbb{R}^d} \|\mathbf{u}\|^2 K(\mathbf{u}) d\mathbf{u} < \infty$, and
- (ii) With probability one, $\mathbf{v} \mapsto \psi(\mathbf{w}, \mathbf{v}) = \mathbb{E}[\mathbf{s}(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}_0) | \mathbf{w}_1 = \mathbf{w}, \mathbf{w}_2 = \mathbf{v}] f_{\mathbf{w}}(\mathbf{v})$ is twice continuously differentiable with $\mathbb{E}[\sup_{\mathbf{v} \in \mathcal{W}} \|\partial_{\mathbf{v}}^{\boldsymbol{\alpha}} \psi(\mathbf{w}, \mathbf{v})\|] < \infty$ for all $\boldsymbol{\alpha} \in \mathbb{Z}_+^d$ with $|\boldsymbol{\alpha}| \leq 2$.

Then $\boldsymbol{\theta}(\cdot)$ admits a $\mathbf{b}_2 \in \mathbb{R}^k$ such that

$$\boldsymbol{\theta}(h) - \boldsymbol{\theta}_0 = \mathbf{b}_2 h^2 + o(h^2) \quad \text{as } h \downarrow 0.$$

The proof of Proposition 1 leverages convexity and may therefore be of independent interest. The convexity argument in question can furthermore be adapted to form the basis of a verification by induction of Assumption 3 with $L > 2$. Details are provided in Section 4.4, which describes the induction step for general L and states explicit (smoothness) conditions under which Assumption 3 holds with $L = 4$.

To describe the debiasing procedure based on generalized jackknifing, we maintain Assumption 3, define $c_0 = 1$, and let $\mathbf{c} = (c_0, \dots, c_{L/2})'$ be a vector of (distinct) positive constants such that the

vector

$$\begin{pmatrix} \lambda_0(\mathbf{c}) \\ \lambda_1(\mathbf{c}) \\ \vdots \\ \lambda_{L/2}(\mathbf{c}) \end{pmatrix} = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ 1 & c_1^2 & \cdots & c_{L/2}^2 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & c_1^L & \cdots & c_{L/2}^L \end{pmatrix}^{-1} \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

is well defined. The debiased estimator is

$$\tilde{\boldsymbol{\theta}}_n = \tilde{\boldsymbol{\theta}}_n(\mathbf{c}, h_n) = \sum_{l=0}^{L/2} \lambda_l(\mathbf{c}) \hat{\boldsymbol{\theta}}_n(h_{n,l}), \quad h_{n,l} = c_l h_n,$$

the construction of which involves solving $L/2 + 1$ convex optimization problems. As defined, the debiased estimator is a generalization of the original pairwise difference estimator because if $L = 0$, then $\mathbf{c} = 1 = \lambda_0$ and therefore $\tilde{\boldsymbol{\theta}}_n = \hat{\boldsymbol{\theta}}_n$.

The next theorem generalizes Theorem 1 by establishing the small bandwidth Gaussian approximation for $\tilde{\boldsymbol{\theta}}_n$. To state the theorem, let

$$\bar{\boldsymbol{\theta}}_n = \bar{\boldsymbol{\theta}}_n(\mathbf{c}, h_n) = \sum_{l=0}^{L/2} \lambda_l(\mathbf{c}) \boldsymbol{\theta}(h_{n,l})$$

and

$$\bar{\mathbf{V}}_n = \bar{\mathbf{V}}_n(\mathbf{c}, h_n) = \boldsymbol{\Gamma}_0^{-1} \left[n^{-1} \boldsymbol{\Sigma}_0 + \binom{n}{2}^{-1} h_n^{-d} \boldsymbol{\Delta}_0(\bar{K}) \right] \boldsymbol{\Gamma}_0^{-1}, \quad \bar{K}(\mathbf{u}) = \bar{K}(\mathbf{u}; \mathbf{c}) = \sum_{l=0}^{L/2} \lambda_l(\mathbf{c}) K_{c_l}(\mathbf{u}).$$

As they should, the expressions have the feature that if $L = 0$, then $\bar{\boldsymbol{\theta}}_n = \boldsymbol{\theta}_n$ and $\bar{\mathbf{V}}_n = \mathbf{V}_n$. Another noteworthy feature of the expressions is that debiasing via generalized jackknifing affects the variance $\bar{\mathbf{V}}_n$ only through the kernel shape entering its small bandwidth component.

Theorem 2. *Suppose Assumptions 1 and 2 hold. If $n^2 h_n^d \rightarrow \infty$ and if $h_n \rightarrow 0$, then*

$$\sup_{\mathbf{t} \in \mathbb{R}^k} \left| \mathbb{P} \left[\bar{\mathbf{V}}_n^{-1/2} (\tilde{\boldsymbol{\theta}}_n - \bar{\boldsymbol{\theta}}_n) \leq \mathbf{t} \right] - \Phi_k(\mathbf{t}) \right| \rightarrow 0.$$

As a consequence, if also Assumption 3 holds and if $nh_n^{2L} \rightarrow 0$, then

$$\sup_{\mathbf{t} \in \mathbb{R}^k} \left| \mathbb{P} \left[\bar{\mathbf{V}}_n^{-1/2} (\tilde{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \leq \mathbf{t} \right] - \Phi_k(\mathbf{t}) \right| \rightarrow 0$$

The magnitude of $\bar{\mathbf{V}}_n^{-1/2}$ is the same as that of $\mathbf{V}_n^{-1/2}$. As a consequence, with obvious modifications the discussion of $\hat{\boldsymbol{\theta}}_n$ following Theorem 1 applies to $\tilde{\boldsymbol{\theta}}_n$, the only noteworthy difference being that (by design), the relevant “small bias” condition is different (and typically milder) in the case of $\tilde{\boldsymbol{\theta}}_n$.

It is worth noting that the equivalent kernel \bar{K} is of higher order, even though the debiased estimator $\tilde{\boldsymbol{\theta}}_n$ only employs estimators constructed using second-order kernels, hereby retaining the

desired convexity for implementation. To be specific, if $\int_{\mathbb{R}^d} \|\mathbf{u}\|^{L+2} K(\mathbf{u}) d\mathbf{u} < \infty$, then

$$\int_{\mathbb{R}^d} \bar{K}(\mathbf{u}) d\mathbf{u} = \sum_{l=0}^{L/2} \lambda_l(\mathbf{c}) \int_{\mathbb{R}^d} K_{c_l}(\mathbf{u}) d\mathbf{u} = \sum_{l=0}^{L/2} \lambda_l(\mathbf{c}) = 1$$

and, for $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_d)' \in \mathbb{Z}_+^d$ with $0 < |\boldsymbol{\alpha}| \leq L+1$,

$$\int_{\mathbb{R}^d} \mathbf{u}^\alpha \bar{K}(\mathbf{u}) d\mathbf{u} = \sum_{l=0}^{L/2} \lambda_l(\mathbf{c}) \int_{\mathbb{R}^d} \mathbf{u}^\alpha K_{c_l}(\mathbf{u}) d\mathbf{u} = \sum_{l=0}^{L/2} \lambda_l(\mathbf{c}) c_l^{|\alpha|} \int_{\mathbb{R}^d} \mathbf{v}^\alpha K(\mathbf{v}) d\mathbf{v} = 0$$

where the last equality uses the defining property of $\{\lambda_l(\mathbf{c})\}$ and symmetry of K , and where \mathbf{u}^α denotes $\prod_{j=1}^d u_j^{\alpha_j}$ for $\mathbf{u} = (u_1, \dots, u_d)' \in \mathbb{R}^d$. In other words, \bar{K} is of order $L+2$.

3.4 Bootstrapping

To develop feasible inference procedures that do not require (explicit) estimation of $\bar{\mathbf{V}}_n$, we consider nonparametric bootstrap-based approximations to the distribution of $\tilde{\boldsymbol{\theta}}_n$. Since $\tilde{\boldsymbol{\theta}}_n = \hat{\boldsymbol{\theta}}_n$ when $L=0$, results for $\hat{\boldsymbol{\theta}}_n$ can be extracted by setting $L=0$ in what follows.

Letting $\mathbf{z}_{1,n}^*, \dots, \mathbf{z}_{n,n}^*$ denote a random sample from the empirical distribution of $\mathbf{z}_1, \dots, \mathbf{z}_n$, the defining property of $\hat{\boldsymbol{\theta}}_n^*(h)$, the nonparametric bootstrap analogue of $\hat{\boldsymbol{\theta}}_n(h)$, is the following:

$$\widehat{M}_n^*(\hat{\boldsymbol{\theta}}_n^*(h); h) \leq \inf_{\boldsymbol{\theta} \in \Theta} \widehat{M}_n^*(\boldsymbol{\theta}; h) + o_{\mathbb{P}}(n^{-1}),$$

where

$$\widehat{M}_n^*(\boldsymbol{\theta}; h) = \binom{n}{2}^{-1} \sum_{i < j} m(\mathbf{z}_{i,n}^*, \mathbf{z}_{j,n}^*; \boldsymbol{\theta}) K_h(\mathbf{w}_{i,n}^* - \mathbf{w}_{j,n}^*).$$

Similarly, the nonparametric bootstrap analogue of $\tilde{\boldsymbol{\theta}}_n$ is

$$\tilde{\boldsymbol{\theta}}_n^* = \tilde{\boldsymbol{\theta}}_n^*(\mathbf{c}) = \sum_{l=0}^{L/2} \lambda_l(\mathbf{c}) \hat{\boldsymbol{\theta}}_n^*(h_{n,l}).$$

The following theorem characterizes the large sample properties of $\tilde{\boldsymbol{\theta}}_n^* - \hat{\boldsymbol{\theta}}_n^*$, the bootstrap counterpart of $\hat{\boldsymbol{\theta}}_n^* - \boldsymbol{\theta}_n$. In perfect analogy with the results in Cattaneo et al. (2014b), we find that the bootstrap distribution estimator is consistent only when $nh_n^d \rightarrow \infty$, but otherwise exhibits a variance inflation making the distributional approximation inconsistent. To state the result, let $\mathbb{P}_n^*[\cdot]$ denote $\mathbb{P}[\cdot | \mathbf{z}_1, \dots, \mathbf{z}_n]$, let $\rightarrow_{\mathbb{P}}$ denote convergence in probability, and define

$$\bar{\mathbf{V}}_n^* = \bar{\mathbf{V}}_n^*(\mathbf{c}, h_n) = \boldsymbol{\Gamma}_0^{-1} \left[n^{-1} \boldsymbol{\Sigma}_0 + 3 \binom{n}{2}^{-1} h_n^{-d} \boldsymbol{\Delta}_0(\bar{K}) \right] \boldsymbol{\Gamma}_0^{-1}.$$

Theorem 3. Suppose Assumptions 1-2 hold and that, for $\boldsymbol{\theta}$ near $\boldsymbol{\theta}_0$, $m(\mathbf{z}, \mathbf{z}; \boldsymbol{\theta}) = 0$ and $\mathbf{s}(\mathbf{z}, \mathbf{z}; \boldsymbol{\theta}) =$

$\mathbf{0}$ (with probability one). If $n^2 h_n^d \rightarrow \infty$ and if $h_n \rightarrow 0$, then

$$\sup_{\mathbf{t} \in \mathbb{R}^k} \left| \mathbb{P}^* \left[\bar{\mathbf{V}}_n^{*-1/2} (\tilde{\boldsymbol{\theta}}_n^* - \tilde{\boldsymbol{\theta}}_n) \leq \mathbf{t} \right] - \Phi_k(\mathbf{t}) \right| \rightarrow_{\mathbb{P}} 0.$$

Because $\bar{\mathbf{V}}_n^{-1} \bar{\mathbf{V}}_n^* \rightarrow \mathbf{I}_k$ if and only if $nh_n^d \rightarrow \infty$ (where \mathbf{I}_k denotes the k -dimensional identity matrix), under the assumptions of Theorem 3

$$\sup_{\mathbf{t} \in \mathbb{R}^k} \left| \mathbb{P}_n^* \left[\tilde{\boldsymbol{\theta}}_n^* - \tilde{\boldsymbol{\theta}}_n \leq \mathbf{t} \right] - \mathbb{P} \left[\tilde{\boldsymbol{\theta}}_n - \bar{\boldsymbol{\theta}}_n \leq \mathbf{t} \right] \right| \rightarrow_{\mathbb{P}} 0$$

if and only if $nh_n^d \rightarrow \infty$. In particular, if $\liminf_{n \rightarrow \infty} nh_n^d < \infty$, then the nonparametric bootstrap is inconsistent, albeit conservative in the sense that the (approximate) variance under the bootstrap distribution is larger than the (approximate) variance of the asymptotic distribution: $\bar{\mathbf{V}}_n^* > \bar{\mathbf{V}}_n$ in a positive definite sense.

The variance inflation problem associated with the nonparametric bootstrap under the small bandwidth regime can be easily fixed by appropriately rescaling the bandwidth used for the bootstrap implementation of the pairwise estimator: employing

$$\check{\boldsymbol{\theta}}_n^* = \check{\boldsymbol{\theta}}_n^*(\mathbf{c}, h_n) = \sum_{l=0}^{L/2} \lambda_l(\mathbf{c}) \hat{\boldsymbol{\theta}}_n^*(3^{1/d} h_{n,l}).$$

and centering its distribution at

$$\check{\boldsymbol{\theta}}_n = \check{\boldsymbol{\theta}}_n(\mathbf{c}, h_n) = \sum_{l=0}^{L/2} \lambda_l(\mathbf{c}) \hat{\boldsymbol{\theta}}_n(3^{1/d} h_{n,l}).$$

automatically adjusts the bootstrap variance, leading to a consistent distributional approximation. Indeed, the following result is an immediate consequence of Theorems 2 and 3.

Corollary 1. *If the assumptions of Theorem 3 hold, then*

$$\sup_{\mathbf{t} \in \mathbb{R}^k} \left| \mathbb{P}_n^* \left[\check{\boldsymbol{\theta}}_n^* - \check{\boldsymbol{\theta}}_n \leq \mathbf{t} \right] - \mathbb{P} \left[\tilde{\boldsymbol{\theta}}_n - \bar{\boldsymbol{\theta}}_n \leq \mathbf{t} \right] \right| \rightarrow_{\mathbb{P}} 0.$$

As a consequence, if also Assumption 3 holds and if $nh_n^{2L} \rightarrow 0$, then

$$\sup_{\mathbf{t} \in \mathbb{R}^k} \left| \mathbb{P}_n^* \left[\check{\boldsymbol{\theta}}_n^* - \check{\boldsymbol{\theta}}_n \leq \mathbf{t} \right] - \mathbb{P} \left[\tilde{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0 \leq \mathbf{t} \right] \right| \rightarrow_{\mathbb{P}} 0.$$

The statement of Corollary 1 emphasizes the rate-adaptive nature of the consistency property enjoyed by the bootstrap distributional approximation. The result has immediate implications for robust inference. For example, letting $\alpha \in (0, 1)$, $\mathbf{a} \in \mathbb{R}^k$ be a fixed vector, and using the “percentile method” (in the terminology of [van der Vaart, 1998](#)), the (nominal) level $1 - \alpha$ bootstrap confidence

interval for $\mathbf{a}'\boldsymbol{\theta}_0$ is

$$\check{C}_n^*(1 - \alpha) = \left[\mathbf{a}'\tilde{\boldsymbol{\theta}}_n - \check{q}_{1-\alpha/2,n}^*, \mathbf{a}'\tilde{\boldsymbol{\theta}}_n - \check{q}_{\alpha/2,n}^* \right], \quad \check{q}_{t,n}^* = \inf \left\{ q \in \mathbb{R} : \mathbb{P}_n^*[\mathbf{a}'\tilde{\boldsymbol{\theta}}_n^* - \mathbf{a}'\tilde{\boldsymbol{\theta}}_n \leq q] \geq t \right\}.$$

If Assumptions 1-3 hold and if $n^2 h_n^d \rightarrow \infty$ and $n h_n^{2L} \rightarrow 0$, then

$$\lim_{n \rightarrow \infty} \mathbb{P} \left[\mathbf{a}'\boldsymbol{\theta}_0 \in \check{C}_n^*(1 - \alpha) \right] = 1 - \alpha.$$

4 Proofs and Other Technical Results

4.1 A Useful Lemma

The following lemma is used in the proofs of Theorems 1-3.

Lemma 1. *Suppose that Assumption 1 holds. Then $\arg \min_{\boldsymbol{\theta} \in \Theta} M(\boldsymbol{\theta}; h)$ is non-empty for $h > 0$ near zero and*

$$\boldsymbol{\theta}(h) - \boldsymbol{\theta}_0 = o(1) \quad \text{as } h \downarrow 0.$$

If also Assumption 2 holds, then $\mathbb{E}[\mathbf{s}(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}(h)) K_h(\mathbf{w}_1 - \mathbf{w}_2)] = \mathbf{0}$ for $h > 0$ near zero.

Proof. For every $\boldsymbol{\theta} \in \Theta$,

$$\begin{aligned} M(\boldsymbol{\theta}; h) &= \int_{\mathcal{W}} \int_{\mathbb{R}^d} \mathbb{E}[m(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}) | \mathbf{w}_1 = \mathbf{w}, \mathbf{w}_2 = \mathbf{w} - \mathbf{u}h] f_{\mathbf{w}}(\mathbf{w}) f_{\mathbf{w}}(\mathbf{w} - \mathbf{u}h) K(\mathbf{u}) d\mathbf{u} d\mathbf{w} \\ &\rightarrow \int_{\mathcal{W}} \mathbb{E}[m(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}) | \mathbf{w}_1 = \mathbf{w}_2 = \mathbf{w}] f_{\mathbf{w}}(\mathbf{w})^2 d\mathbf{w} = M_0(\boldsymbol{\theta}) \quad \text{as } h \downarrow 0, \end{aligned}$$

the convergence being uniform on compact subsets of Θ because $\boldsymbol{\theta} \mapsto M(\boldsymbol{\theta}; h)$ is convex (e.g., Hjort and Pollard, 1993, Lemma 1).

Take any $\epsilon > 0$ with $\Theta_0^\epsilon = \{\boldsymbol{\theta} \in \mathbb{R}^k : \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| \leq \epsilon\} \subseteq \Theta$. By the preceding paragraph,

$$\sup_{\boldsymbol{\theta} \in \Theta_0^\epsilon} |M(\boldsymbol{\theta}; h) - M_0(\boldsymbol{\theta})| \leq \frac{1}{2} \left(\inf_{\boldsymbol{\theta}: \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| = \epsilon} M_0(\boldsymbol{\theta}) - M_0(\boldsymbol{\theta}_0) \right)$$

for $h > 0$ near zero. For any such h and any $\boldsymbol{\theta} \in \Theta \setminus \Theta_0^\epsilon$, we have

$$\eta = \frac{\epsilon}{\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|} \in (0, 1),$$

and therefore, by convexity of $\boldsymbol{\theta} \mapsto M(\boldsymbol{\theta}; h)$,

$$M(\eta\boldsymbol{\theta} + (1 - \eta)\boldsymbol{\theta}_0; h) \leq \eta M(\boldsymbol{\theta}; h) + (1 - \eta) M(\boldsymbol{\theta}_0; h),$$

which rearranges as

$$M(\boldsymbol{\theta}; h) - M(\boldsymbol{\theta}_0; h) \geq \frac{1}{\eta} [M(\eta\boldsymbol{\theta} + (1 - \eta)\boldsymbol{\theta}_0; h) - M(\boldsymbol{\theta}_0; h)] \geq 0.$$

As a consequence,

$$\inf_{\boldsymbol{\theta} \in \Theta} M(\boldsymbol{\theta}; h) = \inf_{\boldsymbol{\theta} \in \Theta_0^\epsilon} M(\boldsymbol{\theta}; h) = \min_{\boldsymbol{\theta} \in \Theta_0^\epsilon} M(\boldsymbol{\theta}; h),$$

where the last equality uses continuity of $\boldsymbol{\theta} \mapsto M(\boldsymbol{\theta}; h)$ and compactness of Θ_0^ϵ .

The above argument shows in particular that $\boldsymbol{\theta}(h) \in \Theta_0^\epsilon$ for $h > 0$ near zero.

If also Assumption 2 holds, then, for $\boldsymbol{\theta}$ near $\boldsymbol{\theta}_0$, $h > 0$ near zero, and for any $\mathbf{t} \in \mathbb{R}^k$,

$$\left| \frac{M(\boldsymbol{\theta} + \mathbf{t}\tau; h) - M(\boldsymbol{\theta}; h)}{\tau} - \mathbb{E}[\mathbf{s}(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}) K_h(\mathbf{w}_1 - \mathbf{w}_2)]' \mathbf{t} \right| \leq |\mathbb{E}[r_{\mathbf{t}}(\boldsymbol{\theta}, \tau) K_h(\mathbf{w}_1 - \mathbf{w}_2)]|$$

$$\rightarrow 0 \quad \text{as } \tau \downarrow 0,$$

implying that for $h > 0$ near zero, $\boldsymbol{\theta} \mapsto M(\boldsymbol{\theta}; h)$ is (directionally) differentiable near $\boldsymbol{\theta}_0$, the directional derivative $\mathbb{E}[\mathbf{s}(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}) K_h(\mathbf{w}_1 - \mathbf{w}_2)]' \mathbf{t}$ being zero when $\boldsymbol{\theta} = \boldsymbol{\theta}(h)$ because $\boldsymbol{\theta}(h)$ minimizes $M(\boldsymbol{\theta}; h)$. \square

4.2 Proof of Theorems 1 and 2

Theorem 1 can be obtained from Theorem 2 by setting $L = 0$ in the latter, so it suffices to prove Theorem 2. To do so, for $l \in \{0, \dots, L/2\}$, let $\widehat{\boldsymbol{\theta}}_{n,l} = \widehat{\boldsymbol{\theta}}_n(h_{n,l})$, $\boldsymbol{\theta}_{n,l} = \boldsymbol{\theta}(h_{n,l})$, and

$$\widehat{\mathbf{U}}_{n,l} = \binom{n}{2}^{-1} \sum_{i < j} \mathbf{s}_{n,l}^\mu(\mathbf{z}_i, \mathbf{z}_j), \quad \mathbf{s}_{n,l}^\mu(\mathbf{z}_i, \mathbf{z}_j) = \mathbf{s}_{n,l}(\mathbf{z}_i, \mathbf{z}_j) - \mathbb{E}[\mathbf{s}_{n,l}(\mathbf{z}_1, \mathbf{z}_2)],$$

where

$$\mathbf{s}_{n,l}(\mathbf{z}_i, \mathbf{z}_j) = \mathbf{s}(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta}_{n,l}) K_{h_{n,l}}(\mathbf{w}_i - \mathbf{w}_j).$$

By Lemma 1, $\lim_{n \rightarrow \infty} \boldsymbol{\theta}_{n,l} = \boldsymbol{\theta}_0$ and $\mathbb{E}[\mathbf{s}_{n,l}(\mathbf{z}_i, \mathbf{z}_j)] = 0$ for large n .

Suppose that

$$\widehat{\boldsymbol{\theta}}_{n,l} - \boldsymbol{\theta}_{n,l} = -\boldsymbol{\Gamma}_0^{-1} \widehat{\mathbf{U}}_{n,l} + o_{\mathbb{P}}(\rho_n^{-1}) \quad \text{for } l \in \{0, \dots, L/2\}. \quad (4.1)$$

Then

$$\widetilde{\boldsymbol{\theta}}_n - \bar{\boldsymbol{\theta}}_n = -\boldsymbol{\Gamma}_0^{-1} \widetilde{\mathbf{U}}_n + o_{\mathbb{P}}(\rho_n^{-1}),$$

where

$$\widetilde{\mathbf{U}}_n = \sum_{l=0}^{L/2} \lambda_l(\mathbf{c}) \widehat{\mathbf{U}}_{n,l}$$

satisfies

$$\left[n^{-1} \boldsymbol{\Sigma}_0 + \binom{n}{2}^{-1} h_n^{-d} \boldsymbol{\Delta}_0(\bar{K}) \right]^{-1/2} \widetilde{\mathbf{U}}_n \rightsquigarrow \mathcal{N}(\mathbf{0}_{k \times 1}, \mathbf{I}_k) \quad (4.2)$$

because, letting \sum_i denote $\sum_{i=1}^n$,

$$\widetilde{\mathbf{U}}_n = \widetilde{\mathbf{L}}_n + \widetilde{\mathbf{W}}_n,$$

where

$$\tilde{\mathbf{L}}_n = n^{-1} \sum_i \tilde{\ell}_n(\mathbf{z}_i), \quad \tilde{\ell}_n(\mathbf{z}_i) = \sum_{l=0}^{L/2} \lambda_l(\mathbf{c}) \mathbb{E}[\mathbf{s}_{n,l}^\mu(\mathbf{z}_i, \mathbf{z}_j) | \mathbf{z}_i] \quad (j \neq i),$$

and

$$\tilde{\mathbf{W}}_n = \binom{n}{2}^{-1} \sum_{i < j} \tilde{\omega}_n(\mathbf{z}_i, \mathbf{z}_j), \quad \tilde{\omega}_n(\mathbf{z}_i, \mathbf{z}_j) = \sum_{l=0}^{L/2} \lambda_l(\mathbf{c}) [\mathbf{s}_{n,l}^\mu(\mathbf{z}_i, \mathbf{z}_j) - \tilde{\ell}_n(\mathbf{z}_i) - \tilde{\ell}_n(\mathbf{z}_j)],$$

satisfy

$$\left(\begin{array}{c} \sqrt{n} \tilde{\mathbf{L}}_n \\ \sqrt{\binom{n}{2}} h_n^d \tilde{\mathbf{W}}_n \end{array} \right) \rightsquigarrow \mathcal{N} \left(\begin{bmatrix} \mathbf{0}_{k \times 1} \\ \mathbf{0}_{k \times 1} \end{bmatrix}, \begin{bmatrix} \Sigma_0 & \mathbf{0}_{k \times k} \\ \mathbf{0}_{k \times k} & \Delta_0(\bar{K}) \end{bmatrix} \right),$$

as can be shown by means of the Cramér-Wold device and the central limit theorem of [Heyde and Brown \(1970\)](#), the latter being applicable because it follows from routine calculations that for every $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2 \in \mathbb{R}^k$, we have

$$\varsigma_n^2 = \sum_i \mathbb{V}[g_{i,n}] = \boldsymbol{\mu}_1' \Sigma_0 \boldsymbol{\mu}_1 + \boldsymbol{\mu}_2' \Delta_0(\bar{K}) \boldsymbol{\mu}_2 + o(1),$$

$$\sum_i \mathbb{E}[g_{i,n}^4] = o(1),$$

and

$$\mathbb{V} \left[\sum_i \varsigma_{i,n}^2 - \varsigma_n^2 \right] = o(1), \quad \varsigma_{i,n}^2 = \mathbb{V}[g_{i,n} | \mathbf{z}_1, \dots, \mathbf{z}_{i-1}],$$

where

$$g_{i,n} = g_{i,n}(\boldsymbol{\mu}) = \frac{2}{\sqrt{n}} \boldsymbol{\mu}_1' \tilde{\ell}_n(\mathbf{z}_i) + \sqrt{\binom{n}{2}^{-1}} h_n^d \sum_{j=1}^{i-1} \boldsymbol{\mu}_2' \tilde{\omega}_n(\mathbf{z}_i, \mathbf{z}_j).$$

The proof of Theorem 2 can therefore be completed by verifying (4.1).

To do so, we leverage convexity. For any $l \in \{0, \dots, L/2\}$ and any $\mathbf{t} \in \mathbb{R}^k$, it can be shown that

$$\lim_{\tau \downarrow 0, h \downarrow 0, \boldsymbol{\theta} \rightarrow \boldsymbol{\theta}_0} \mathbb{E}[\mathbb{E}[r_{\mathbf{t}}(\boldsymbol{\theta}, \tau) K_h(\mathbf{w}_1 - \mathbf{w}_2) | \mathbf{z}_1]^2] = 0$$

and

$$\lim_{\tau \downarrow 0, h \downarrow 0, \boldsymbol{\theta} \rightarrow \boldsymbol{\theta}_0} h^d \mathbb{E}[r_{\mathbf{t}}(\boldsymbol{\theta}, \tau)^2 K_h(\mathbf{w}_1 - \mathbf{w}_2)^2] = 0,$$

and it therefore follows from a Hoeffding decomposition that

$$\begin{aligned} & \rho_n^2 \left[\widehat{M}_n(\boldsymbol{\theta}_{n,l} + \mathbf{t} \rho_n^{-1}; h_{n,l}) - \widehat{M}_n(\boldsymbol{\theta}_{n,l}; h_{n,l}) \right] \\ &= \rho_n^2 [M(\boldsymbol{\theta}_{n,l} + \mathbf{t} \rho_n^{-1}; h_{n,l}) - M(\boldsymbol{\theta}_{n,l}; h_{n,l})] + \mathbf{t}' \rho_n \widehat{\mathbf{U}}_{n,l} + o_{\mathbb{P}}(1). \end{aligned}$$

Moreover,

$$\rho_n^2[M(\boldsymbol{\theta}_{n,l} + \mathbf{t}\rho_n^{-1}; h_{n,l}) - M(\boldsymbol{\theta}_{n,l}; h_{n,l})] \rightarrow \frac{1}{2}\mathbf{t}'\boldsymbol{\Gamma}_0\mathbf{t},$$

and proceeding as in the proof of (4.2) it can be shown that $\rho_n\widehat{\mathbf{U}}_{n,l} = O_{\mathbb{P}}(1)$. Because $\boldsymbol{\Gamma}_0$ is positive definite and because $\mathbf{t} \mapsto \widehat{M}_n(\boldsymbol{\theta}_{n,l} + \mathbf{t}\rho_n^{-1}; h_{n,l})$ is convex (almost surely), the corollary following Hjort and Pollard (1993, Lemma 2) implies that (4.1) holds.

4.3 Proof of Theorem 3

The proof of Theorem 3 is a natural bootstrap analog of the proof of Theorem 2.

For $l \in \{0, \dots, L/2\}$, let $\widehat{\boldsymbol{\theta}}_{n,l}^* = \widehat{\boldsymbol{\theta}}_n^*(h_{n,l})$ and

$$\widehat{\mathbf{U}}_{n,l}^* = \binom{n}{2}^{-1} \sum_{i < j} \mathbf{s}_{n,l}^{\mu,*}(\mathbf{z}_{i,n}^*, \mathbf{z}_{j,n}^*), \quad \mathbf{s}_{n,l}^{\mu,*}(\mathbf{z}_{i,n}^*, \mathbf{z}_{j,n}^*) = \mathbf{s}_{n,l}(\mathbf{z}_{i,n}^*, \mathbf{z}_{j,n}^*) - \mathbb{E}_n^*[\mathbf{s}_{n,l}(\mathbf{z}_{1,n}^*, \mathbf{z}_{2,n}^*)],$$

where $\mathbb{E}_n^*[\cdot]$ denotes $\mathbb{E}[\cdot | \mathbf{z}_1, \dots, \mathbf{z}_n]$.

It suffices to show that

$$\widehat{\boldsymbol{\theta}}_{n,l}^* - \boldsymbol{\theta}_{n,l} = -\boldsymbol{\Gamma}_0^{-1}(\widehat{\mathbf{U}}_{n,l}^* + \widehat{\mathbf{U}}_{n,l}) + o_{\mathbb{P}}(\rho_n^{-1}) \quad \text{for } l \in \{0, \dots, L/2\} \quad (4.3)$$

and that

$$\left[n^{-1}\boldsymbol{\Sigma}_0 + 3\binom{n}{2}^{-1} h_n^{-d}\boldsymbol{\Delta}_0(K) \right]^{-1/2} \widetilde{\mathbf{U}}_n^* \rightsquigarrow_{\mathbb{P}} \mathcal{N}(\mathbf{0}_{k \times 1}, \mathbf{I}_k), \quad (4.4)$$

where

$$\widetilde{\mathbf{U}}_n^* = \sum_{l=0}^{L/2} \lambda_l(\mathbf{c}) \widehat{\mathbf{U}}_{n,l}^*,$$

and where $\rightsquigarrow_{\mathbb{P}}$ denotes weak convergence in probability.

For every $\mathbf{t} \in \mathbb{R}^k$, using a Hoeffding decomposition and the fact that $m(\mathbf{z}, \mathbf{z}; \boldsymbol{\theta}_{n,l}) = 0$ and $\mathbf{s}(\mathbf{z}, \mathbf{z}; \boldsymbol{\theta}_{n,l}) = \mathbf{0}$ for large n , we have

$$\begin{aligned} & \rho_n^2 \left[\widehat{M}_n^*(\boldsymbol{\theta}_{n,l} + \mathbf{t}\rho_n^{-1}; h_{n,l}) - \widehat{M}_n^*(\boldsymbol{\theta}_{n,l}; h_{n,l}) \right] \\ &= \rho_n^2(1 + o(1)) \left[\widehat{M}_n(\boldsymbol{\theta}_{n,l} + \mathbf{t}\rho_n^{-1}; h_{n,l}) - \widehat{M}_n(\boldsymbol{\theta}_{n,l}; h_{n,l}) \right] + \mathbf{t}'\rho_n\widehat{\mathbf{U}}_{n,l}^* + o_{\mathbb{P}}(1), \end{aligned}$$

where it can be shown that $\rho_n\widehat{\mathbf{U}}_{n,l}^* = O_{\mathbb{P}}(1)$ and where it follows from the proof of (4.1) that

$$\rho_n^2 \left[\widehat{M}_n(\boldsymbol{\theta}_{n,l} + \mathbf{t}\rho_n^{-1}; h_{n,l}) - \widehat{M}_n(\boldsymbol{\theta}_{n,l}; h_{n,l}) \right] = \frac{1}{2}\mathbf{t}'\boldsymbol{\Gamma}_0\mathbf{t} + \mathbf{t}'\rho_n\widehat{\mathbf{U}}_{n,l} + o_{\mathbb{P}}(1),$$

where $\rho_n\widehat{\mathbf{U}}_{n,l} = O_{\mathbb{P}}(1)$. In other words,

$$\rho_n^2 \left[\widehat{M}_n^*(\boldsymbol{\theta}_{n,l} + \mathbf{t}\rho_n^{-1}; h_{n,l}) - \widehat{M}_n^*(\boldsymbol{\theta}_{n,l}; h_{n,l}) \right]$$

$$= \frac{1}{2} \mathbf{t}' \mathbf{\Gamma}_0 \mathbf{t} + \mathbf{t}' \rho_n (\widehat{\mathbf{U}}_{n,l}^* + \widehat{\mathbf{U}}_{n,l}) + o_{\mathbb{P}}(1) \quad \text{for every } \mathbf{t} \in \mathbb{R}^k.$$

Because $\mathbf{\Gamma}_0$ is positive definite and because $\mathbf{t} \mapsto \widehat{M}_n^*(\boldsymbol{\theta}_{n,l} + \mathbf{t} \rho_n^{-1}; h_{n,l})$ is convex (almost surely), the corollary following Hjort and Pollard (1993, Lemma 2) implies that (4.3) holds.

To prove (4.4), we begin by decomposing $\widetilde{\mathbf{U}}_n^*$ as

$$\widetilde{\mathbf{U}}_n^* = \widetilde{\mathbf{I}}_n^* + \widetilde{\mathbf{W}}_n^*,$$

where

$$\widetilde{\mathbf{I}}_n^* = n^{-1} \sum_i \widetilde{\boldsymbol{\ell}}_n^*(\mathbf{z}_{i,n}^*), \quad \widetilde{\boldsymbol{\ell}}_n^*(\mathbf{z}_{i,n}^*) = \sum_{l=0}^{L/2} \lambda_l(\mathbf{c}) \mathbb{E}_n^*[\mathbf{s}_{n,l}^{\mu,*}(\mathbf{z}_{i,n}^*, \mathbf{z}_{j,n}^*) | \mathbf{z}_{i,n}^*] \quad (j \neq i),$$

and

$$\widetilde{\mathbf{W}}_n^* = \binom{n}{2}^{-1} \sum_{i < j} \widetilde{\boldsymbol{\omega}}_n^*(\mathbf{z}_{i,n}^*, \mathbf{z}_{j,n}^*), \quad \widetilde{\boldsymbol{\omega}}_n^*(\mathbf{z}_{i,n}^*, \mathbf{z}_{j,n}^*) = \sum_{l=0}^{L/2} \lambda_l(\mathbf{c}) [\mathbf{s}_{n,l}^{\mu,*}(\mathbf{z}_{i,n}^*, \mathbf{z}_{j,n}^*) - \widetilde{\boldsymbol{\ell}}_n^*(\mathbf{z}_{i,n}^*) - \widetilde{\boldsymbol{\ell}}_n^*(\mathbf{z}_{j,n}^*)].$$

Defining

$$\pi_n = \frac{\sqrt{nh_n^d}}{1 + \sqrt{nh_n^d}},$$

routine calculations can be used to show that for every $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2 \in \mathbb{R}^k$, we have

$$\widehat{\boldsymbol{\zeta}}_n^2 = \sum_i \mathbb{V}_n^*[g_{i,n}^*] = \boldsymbol{\mu}_1' [\pi_n^2 \boldsymbol{\Sigma}_0 + 4(1 - \pi_n)^2 \boldsymbol{\Delta}_0(\bar{K})] \boldsymbol{\mu}_1 + \boldsymbol{\mu}_2' \boldsymbol{\Delta}_0(\bar{K}) \boldsymbol{\mu}_2 + o_{\mathbb{P}}(1),$$

$$\sum_i \mathbb{E}_n^*[g_{i,n}^{*4}] = o_{\mathbb{P}}(1),$$

and

$$\mathbb{V} \left[\sum_i \widehat{\boldsymbol{\zeta}}_{i,n}^2 - \widehat{\boldsymbol{\zeta}}_n^2 \right] = o_{\mathbb{P}}(1), \quad \widehat{\boldsymbol{\zeta}}_{i,n}^2 = \mathbb{V}_n^*[g_{i,n}^* | \mathbf{z}_{1,n}^*, \dots, \mathbf{z}_{i-1,n}^*],$$

where $\mathbb{V}_n^*[\cdot]$ denotes $\mathbb{V}[\cdot | \mathbf{z}_1, \dots, \mathbf{z}_n]$ and where

$$g_{i,n}^* = g_{i,n}^*(\boldsymbol{\mu}) = \frac{\pi_n}{\sqrt{n}} 2 \boldsymbol{\mu}_1' \widetilde{\boldsymbol{\ell}}_n^*(\mathbf{z}_{i,n}^*) + \sqrt{\binom{n}{2}^{-1} h_n^d} \sum_{j=1}^{i-1} \boldsymbol{\mu}_2' \widetilde{\boldsymbol{\omega}}_n^*(\mathbf{z}_{i,n}^*, \mathbf{z}_{j,n}^*).$$

The Cramér-Wold device and the central limit theorem of Heyde and Brown (1970) therefore imply that if $\pi_n \rightarrow \pi_0 \in [0, 1]$, then

$$\left(\frac{\sqrt{n} \pi_n \widetilde{\mathbf{I}}_n^*}{\sqrt{\binom{n}{2} h_n^d} \widetilde{\mathbf{W}}_n^*} \right) \rightsquigarrow_{\mathbb{P}} \mathcal{N} \left(\begin{bmatrix} \mathbf{0}_{k \times 1} \\ \mathbf{0}_{k \times 1} \end{bmatrix}, \begin{bmatrix} \pi_0^2 \boldsymbol{\Sigma}_0 + 4(1 - \pi_0)^2 \boldsymbol{\Delta}_0(\bar{K}) & \mathbf{0}_{k \times k} \\ \mathbf{0}_{k \times k} & \boldsymbol{\Delta}_0(\bar{K}) \end{bmatrix} \right).$$

Whether or not π_n is convergent, the result (4.4) can be obtained from the preceding display by arguing along subsequences (if necessary).

4.4 Verifying Assumption 3

It follows from Lemma 1 that if Assumption 1 holds, then so does Assumption 3 with $L = 0$. This observation provides the base case for an induction argument. To describe the induction step, suppose that for some even $L \geq 0$, we have

$$\boldsymbol{\theta}(h) - \boldsymbol{\theta}_0 = \sum_{l=1}^{L/2} \mathbf{b}_{2l} h^{2l} + o(h^L) \quad \text{as } h \downarrow 0.$$

Suppose also that, for every $\mathbf{t} \in \mathbb{R}^k$ and some $\boldsymbol{\beta}_{L+2} \in \mathbb{R}^k$, we have

$$\begin{aligned} h^{-(L+2)} \left[M \left(\boldsymbol{\theta}_0 + \sum_{l=1}^{L/2} \mathbf{b}_{2l} h^{2l} + \mathbf{t} h^{L+2}; h \right) - M \left(\boldsymbol{\theta}_0 + \sum_{l=1}^{L/2} \mathbf{b}_{2l} h^{2l}; h \right) \right] \\ = \mathbf{t}' \boldsymbol{\beta}_{L+2} + \frac{1}{2} \mathbf{t}' \boldsymbol{\Gamma}_0 \mathbf{t} + o(1) \quad \text{as } h \downarrow 0. \end{aligned}$$

Then, the corollary following Hjort and Pollard (1993, Lemma 2) implies that

$$\begin{aligned} h^{-(L+2)} \left(\boldsymbol{\theta}(h) - \boldsymbol{\theta}_0 - \sum_{l=1}^{L/2} \mathbf{b}_{2l} h^{2l} \right) &= \arg \min_{\mathbf{t} \in \mathbb{R}^k} M \left(\boldsymbol{\theta}_0 + \sum_{l=1}^{L/2} \mathbf{b}_{2l} h^{2l} + \mathbf{t} h^{L+2}; h \right) \\ &= -\boldsymbol{\Gamma}_0^{-1} \boldsymbol{\beta}_{L+2} + o(1) \quad \text{as } h \downarrow 0; \end{aligned}$$

that is, defining $\mathbf{b}_{L+2} = -\boldsymbol{\Gamma}_0^{-1} \boldsymbol{\beta}_{L+2}$, we have

$$\boldsymbol{\theta}(h) - \boldsymbol{\theta}_0 = \sum_{l=1}^{(L+2)/2} \mathbf{b}_{2l} h^{2l} + o(h^{L+2}) \quad \text{as } h \downarrow 0.$$

To complete the proof of Proposition 1, it therefore suffices to note that (for every $\mathbf{t} \in \mathbb{R}^k$ and) under the assumptions of the proposition, we have

$$h^{-4} [M(\boldsymbol{\theta}_0 + \mathbf{t} h^2; h) - M(\boldsymbol{\theta}_0; h)] = \mathbf{t}' \boldsymbol{\beta}_2 + \frac{1}{2} \mathbf{t}' \boldsymbol{\Gamma}_0 \mathbf{t} + o(1) \quad \text{as } h \downarrow 0,$$

where

$$\boldsymbol{\beta}_2 = \frac{1}{2} \sum_{i=1}^d \int_{\mathcal{W}} \frac{\partial^2 \psi(\mathbf{w}, \mathbf{v})}{\partial v_i^2} \Big|_{\mathbf{v}=\mathbf{w}} f_{\mathbf{w}}(\mathbf{w}) d\mathbf{w} \int_{\mathbb{R}^k} u_i^2 K(\mathbf{u}) d\mathbf{u}.$$

Similarly, if in addition to the assumptions of Proposition 1 it is assumed that for every $\mathbf{t} \in \mathbb{R}^k$

and for some $\beta_4 \in \mathbb{R}^k$, we have

$$h^{-8} [M(\theta_0 + \mathbf{b}_2 h^2 + \mathbf{t} h^4; h) - M(\theta_0 + \mathbf{b}_2 h^2; h)] = \mathbf{t}' \beta_4 + \frac{1}{2} \mathbf{t}' \mathbf{T}_0 \mathbf{t} + o(1) \quad \text{as } h \downarrow 0,$$

then Assumption 3 holds with $L = 4$. One set of sufficient conditions for this to occur is that Assumptions 1-2 hold and that, for every $\mathbf{t} \in \mathbb{R}^k$, the following are satisfied (with probability one):

- (i) $\int_{\mathbb{R}^d} \|\mathbf{u}\|^4 K(\mathbf{u}) d\mathbf{u} < \infty$.
- (ii) $\mathbf{v} \mapsto \psi(\mathbf{w}, \mathbf{v}) = \mathbb{E}[\mathbf{s}(\mathbf{z}_1, \mathbf{z}_2; \theta_0) | \mathbf{w}_1 = \mathbf{w}, \mathbf{w}_2 = \mathbf{v}] f_{\mathbf{w}}(\mathbf{v})$ is four times continuously differentiable with $\mathbb{E}[\sup_{\mathbf{v} \in \mathcal{W}} \|\partial_{\mathbf{v}}^{\alpha} \psi(\mathbf{w}, \mathbf{v})\|] < \infty$ for all $\alpha \in \mathbb{Z}_+^d$ with $|\alpha| \leq 4$.
- (iii) $f_{\mathbf{w}}$ is twice continuously differentiable and $\mathbf{v} \mapsto \mathbf{H}(\mathbf{w}, \mathbf{v}; \theta_0, \mathbf{t})$ is twice continuously differentiable with $\mathbb{E}[\sup_{\mathbf{v} \in \mathcal{W}} \|\partial_{\mathbf{v}}^{\alpha} \mathbf{H}(\mathbf{w}, \mathbf{v}; \theta_0, \mathbf{t}) f_{\mathbf{w}}(\mathbf{v})\|] < \infty$ for all $\alpha \in \mathbb{Z}_+^d$ with $|\alpha| \leq 2$.
- (iv) For some function $\dot{\mathbf{H}}(\mathbf{w}, \mathbf{v}; \theta, \mathbf{t}) \in \mathbb{R}^{k \times k}$, $\mathbf{v} \mapsto \dot{\mathbf{H}}(\mathbf{w}, \mathbf{v}; \theta, \mathbf{t})$ is continuous,

$$\mathbb{E} \left[\sup_{\mathbf{v} \in \mathcal{W}} \|\dot{\mathbf{H}}(\mathbf{w}, \mathbf{v}; \theta_0, \mathbf{t}) f_{\mathbf{w}}(\mathbf{v})\| \right] < \infty,$$

$$\lim_{\tau \downarrow 0, (\theta, \mathbf{u}) \rightarrow (\theta_0, \mathbf{0})} \left\| \frac{\mathbf{H}(\mathbf{w}, \mathbf{w} + \mathbf{u}; \theta + \tau \mathbf{t}, \mathbf{t}) - \mathbf{H}(\mathbf{w}, \mathbf{w} + \mathbf{u}; \theta, \mathbf{t})}{\tau} - \dot{\mathbf{H}}(\mathbf{w}, \mathbf{w} + \mathbf{u}; \theta, \mathbf{t}) \right\| = 0,$$

and, for some $\delta > 0$,

$$\mathbb{E} \left[\sup_{\tau \in (0, \delta), \|\theta - \theta_0\| < \delta, \mathbf{w}_2 \in \mathcal{W}} \left\| \frac{\mathbf{H}(\mathbf{w}_1, \mathbf{w}_2; \theta + \tau \mathbf{t}, \mathbf{t}) - \mathbf{H}(\mathbf{w}_1, \mathbf{w}_2; \theta, \mathbf{t})}{\tau} - \dot{\mathbf{H}}(\mathbf{w}_1, \mathbf{w}_2; \theta, \mathbf{t}) \right\| \right] < \infty.$$

5 Sufficient Conditions for Motivating Examples

To demonstrate the plausibility of Assumptions 1 and 2, we revisit the examples of Section 2. In each example, Assumptions 1(ii) holds and Assumption 1(iii) is fairly primitive, so we focus on giving primitive sufficient conditions for Assumptions 1(iv)-(v) and 2.

5.1 Partially Linear Regression Model

We take $\mathbf{s} = \mathbf{s}_{\text{PLR}}$ and $\mathbf{H} = \mathbf{H}_{\text{PLR}}$, where

$$\mathbf{s}_{\text{PLR}}(\mathbf{z}_i, \mathbf{z}_j; \theta) = \frac{\partial}{\partial \theta} m_{\text{PLR}}(\mathbf{z}_i, \mathbf{z}_j; \theta) = -\dot{\mathbf{x}}_{i,j} (\dot{y}_{i,j} - \dot{\mathbf{x}}'_{i,j} \theta)$$

and

$$\begin{aligned} \mathbf{H}_{\text{PLR}}(\mathbf{w}_i, \mathbf{w}_j) &= \frac{\partial^2}{\partial \theta \partial \theta'} \mathbb{E}[m_{\text{PLR}}(\mathbf{z}_i, \mathbf{z}_j; \theta) | \mathbf{w}_i, \mathbf{w}_j] = \frac{\partial}{\partial \theta'} \mathbb{E}[\mathbf{s}_{\text{PLR}}(\mathbf{z}_i, \mathbf{z}_j; \theta) | \mathbf{w}_i, \mathbf{w}_j] \\ &= \mathbb{E}[\dot{\mathbf{x}}_{i,j} \dot{\mathbf{x}}'_{i,j} | \mathbf{w}_i, \mathbf{w}_j], \end{aligned}$$

the latter depending on neither $\boldsymbol{\theta}$ nor \mathbf{t} (because $m_{\text{PLR}}(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta})$ is quadratic in $\boldsymbol{\theta}$).

Under mild conditions, Assumptions 1(iv)-(v) and 2 hold with

$$\boldsymbol{\xi}_0(\mathbf{z}) = 2\mathbb{E}[\mathbf{s}_{\text{PLR}}(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}_0)|\mathbf{z}_1 = \mathbf{z}, \mathbf{w}_2 = \mathbf{w}]f_{\mathbf{w}}(\mathbf{w}),$$

$$\boldsymbol{\Xi}_0(\mathbf{w}) = \mathbb{E}[\mathbf{s}_{\text{PLR}}(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}_0)\mathbf{s}_{\text{PLR}}(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}_0)'|\mathbf{w}_1 = \mathbf{w}_2 = \mathbf{w}]f_{\mathbf{w}}(\mathbf{w}),$$

and

$$\mathbf{G}_0(\mathbf{w}) = \mathbf{H}_{\text{PLR}}(\mathbf{w}, \mathbf{w})f_{\mathbf{w}}(\mathbf{w}).$$

For instance, it suffices to set $b(\mathbf{z}) = (1 + \|\mathbf{x}\|)(1 + |\varepsilon| + |\gamma_0(\mathbf{w})| + \|\mathbf{x}\|)$ and to assume that

(i) The functions $\mathbf{w} \mapsto \gamma_0(\mathbf{w})$, $\mathbf{w} \mapsto \mathbb{E}[\mathbf{x}|\mathbf{w}]$, $\mathbf{w} \mapsto \mathbb{E}[\mathbf{x}\mathbf{x}'|\mathbf{w}]$, $\mathbf{w} \mapsto \mathbb{E}[\varepsilon^2|\mathbf{w}]$, $\mathbf{w} \mapsto \mathbb{E}[\mathbf{x}\varepsilon^2|\mathbf{w}]$, and $\mathbf{w} \mapsto \mathbb{E}[\mathbf{x}\mathbf{x}'\varepsilon^2|\mathbf{w}]$ are continuous on \mathcal{W} .

(ii) $\mathbb{E}[(1 + \|\mathbf{x}\|^4)\varepsilon^4] + \sup_{\mathbf{w} \in \mathcal{W}} \mathbb{E}[(1 + \|\mathbf{x}\|^4)\varepsilon^4|\mathbf{w}]f_{\mathbf{w}}(\mathbf{w}) < \infty$ and

$$\mathbb{E}[(1 + \|\mathbf{x}\|^4)\gamma_0(\mathbf{w})^4 + \|\mathbf{x}\|^8] + \sup_{\mathbf{w} \in \mathcal{W}} \mathbb{E}[(1 + \|\mathbf{x}\|^4)\gamma_0(\mathbf{w})^4 + \|\mathbf{x}\|^8|\mathbf{w}]f_{\mathbf{w}}(\mathbf{w}) < \infty.$$

(iii) With probability one, $\mathbb{V}[\mathbf{x}|\mathbf{w}]$ is positive definite and $\mathbb{V}[\varepsilon|\mathbf{x}, \mathbf{w}] > 0$.

5.2 Partially Linear Logit Model

We take $\mathbf{s} = \mathbf{s}_{\text{PLL}}$ and $\mathbf{H} = \mathbf{H}_{\text{PLL}}$, where

$$\mathbf{s}_{\text{PLL}}(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta}) = \frac{\partial}{\partial \boldsymbol{\theta}} m_{\text{PLL}}(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta}) = -\dot{\mathbf{x}}_{i,j}(y_i - \Lambda(\dot{\mathbf{x}}'_{i,j}\boldsymbol{\theta}))\mathbb{1}\{y_{i,j} \neq 0\}$$

and

$$\begin{aligned} \mathbf{H}_{\text{PLL}}(\mathbf{w}_i, \mathbf{w}_j; \boldsymbol{\theta}) &= \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \mathbb{E}[m_{\text{PLL}}(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta})|\mathbf{w}_i, \mathbf{w}_j] = \frac{\partial}{\partial \boldsymbol{\theta}'} \mathbb{E}[\mathbf{s}_{\text{PLL}}(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta})|\mathbf{w}_i, \mathbf{w}_j] \\ &= \mathbb{E}[\dot{\mathbf{x}}_{i,j}\dot{\mathbf{x}}'_{i,j}\lambda(\dot{\mathbf{x}}'_{i,j}\boldsymbol{\theta})\mathbb{1}\{y_{i,j} \neq 0\}|\mathbf{w}_i, \mathbf{w}_j], \end{aligned}$$

where the latter does not depend on \mathbf{t} (because $m_{\text{PLL}}(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta})$ is twice differentiable in $\boldsymbol{\theta}$).

Under mild conditions, Assumptions 1(iv)-(v) and 2 hold with

$$\boldsymbol{\xi}_0(\mathbf{z}) = 2\mathbb{E}[\mathbf{s}_{\text{PLL}}(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}_0)|\mathbf{z}_1 = \mathbf{z}, \mathbf{w}_2 = \mathbf{w}]f_{\mathbf{w}}(\mathbf{w}),$$

$$\boldsymbol{\Xi}_0(\mathbf{w}) = \mathbb{E}[\mathbf{s}_{\text{PLL}}(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}_0)\mathbf{s}_{\text{PLL}}(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}_0)'|\mathbf{w}_1 = \mathbf{w}_2 = \mathbf{w}]f_{\mathbf{w}}(\mathbf{w}),$$

and

$$\mathbf{G}_0(\mathbf{w}) = \mathbf{H}_{\text{PLL}}(\mathbf{w}, \mathbf{w}; \boldsymbol{\theta}_0)f_{\mathbf{w}}(\mathbf{w}).$$

For instance, it suffices to set $b(\mathbf{z}) = 1 + \|\mathbf{x}\|$ and to assume that, for some $\delta > 0$,

- (i) The function $\mathbf{w} \mapsto \gamma_0(\mathbf{w})$ is continuous on \mathcal{W} . Also, the conditional distribution of \mathbf{x} given \mathbf{w} admits a density $f_{\mathbf{x}|\mathbf{w}}$ with respect to some measure ρ such that $\mathbf{w} \mapsto f_{\mathbf{x}|\mathbf{w}}(\mathbf{x}|\mathbf{w})$ is continuous on \mathcal{W} (with probability one) and

$$\int_{\mathbb{R}^k} (1 + \|\mathbf{x}\|^2) \sup_{\|\mathbf{u}\| \leq \delta} f_{\mathbf{x}|\mathbf{w}}(\mathbf{x}|\mathbf{w} + \mathbf{u}) d\rho(\mathbf{x}) < \infty \quad \text{for every } \mathbf{w} \in \mathcal{W}.$$

- (ii) $\mathbb{E}[\|\mathbf{x}\|^4] + \sup_{\mathbf{w} \in \mathcal{W}} \mathbb{E}[\|\mathbf{x}\|^4 | \mathbf{w}] f_{\mathbf{w}}(\mathbf{w}) < \infty$.
 (iii) With probability one, $\mathbb{V}[\mathbf{x}|\mathbf{w}]$ is positive definite.

5.3 Partially Linear Tobit Model

We take $\mathbf{s} = \mathbf{s}_{\text{PLT}}$ and $\mathbf{H} = \mathbf{H}_{\text{PLT}}$, where

$$\mathbf{s}_{\text{PLT}}(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta}) = \dot{\mathbf{x}}_{i,j} (\mathbb{1}\{y_j > \max(y_i - \dot{\mathbf{x}}'_{i,j} \boldsymbol{\theta}, 0)\} - \mathbb{1}\{y_i > \max(y_j + \dot{\mathbf{x}}'_{i,j} \boldsymbol{\theta}, 0)\})$$

and

$$\begin{aligned} \mathbf{H}_{\text{PLT}}(\mathbf{w}_i, \mathbf{w}_j; \boldsymbol{\theta}, \mathbf{t}) &= \mathbb{E}[\dot{\mathbf{x}}_{i,j} \dot{\mathbf{x}}'_{i,j} (\mathbb{1}\{\dot{\mathbf{x}}'_{i,j} \boldsymbol{\theta} > 0\} + \mathbb{1}\{\dot{\mathbf{x}}'_{i,j} \boldsymbol{\theta} = 0, \dot{\mathbf{x}}'_{i,j} \mathbf{t} \geq 0\}) \eta_{\text{PLT}}(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta}) | \mathbf{w}_i, \mathbf{w}_j] \\ &\quad + \mathbb{E}[\dot{\mathbf{x}}_{i,j} \dot{\mathbf{x}}'_{i,j} (\mathbb{1}\{\dot{\mathbf{x}}'_{i,j} \boldsymbol{\theta} < 0\} + \mathbb{1}\{\dot{\mathbf{x}}'_{i,j} \boldsymbol{\theta} = 0, \dot{\mathbf{x}}'_{i,j} \mathbf{t} < 0\}) \eta_{\text{PLT}}(\mathbf{z}_j, \mathbf{z}_i; \boldsymbol{\theta}) | \mathbf{w}_i, \mathbf{w}_j], \end{aligned}$$

with

$$\begin{aligned} \eta_{\text{PLT}}(\mathbf{z}_i, \mathbf{z}_j; \boldsymbol{\theta}) &= 2 \int_0^\infty f_{\varepsilon|\mathbf{w}}(\varepsilon - \mathbf{x}'_i \boldsymbol{\theta}_0 - \gamma_0(\mathbf{w}_i) + \dot{\mathbf{x}}'_{i,j} \boldsymbol{\theta} | \mathbf{w}_i) f_{\varepsilon|\mathbf{w}}(\varepsilon - \mathbf{x}'_j \boldsymbol{\theta}_0 - \gamma_0(\mathbf{w}_j) | \mathbf{w}_j) d\varepsilon \\ &\quad + f_{\varepsilon|\mathbf{w}}(-\mathbf{x}'_i \boldsymbol{\theta}_0 - \gamma_0(\mathbf{w}_i) + \dot{\mathbf{x}}'_{i,j} \boldsymbol{\theta} | \mathbf{w}_i) \int_{-\infty}^0 f_{\varepsilon|\mathbf{w}}(\varepsilon - \mathbf{x}'_j \boldsymbol{\theta}_0 - \gamma_0(\mathbf{w}_j) | \mathbf{w}_j) d\varepsilon. \end{aligned}$$

Under mild conditions, Assumptions 1(iv)-(v) and 2 hold with

$$\boldsymbol{\xi}_0(\mathbf{z}) = 2\mathbb{E}[\mathbf{s}_{\text{PLT}}(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}_0) | \mathbf{z}_1 = \mathbf{z}, \mathbf{w}_2 = \mathbf{w}] f_{\mathbf{w}}(\mathbf{w}),$$

$$\boldsymbol{\Xi}_0(\mathbf{w}) = \mathbb{E}[\mathbf{s}_{\text{PLT}}(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}_0) \mathbf{s}_{\text{PLT}}(\mathbf{z}_1, \mathbf{z}_2; \boldsymbol{\theta}_0)' | \mathbf{w}_1 = \mathbf{w}_2 = \mathbf{w}] f_{\mathbf{w}}(\mathbf{w}),$$

and

$$\mathbf{G}_0(\mathbf{w}) = \mathbf{H}_{\text{PLT}}(\mathbf{w}, \mathbf{w}; \boldsymbol{\theta}_0, \boldsymbol{\theta}_0) f_{\mathbf{w}}(\mathbf{w}).$$

For instance, it suffices to set $b(\mathbf{z}) = 1 + \|\mathbf{x}\|$ and to assume that, for some $\delta > 0$,

- (i) The function $\mathbf{w} \mapsto \gamma_0(\mathbf{w})$ is continuous on \mathcal{W} . Also, the conditional distribution of \mathbf{x} given \mathbf{w} admits a density $f_{\mathbf{x}|\mathbf{w}}$ with respect to some measure ρ such that $\mathbf{w} \mapsto f_{\mathbf{x}|\mathbf{w}}(\mathbf{x}|\mathbf{w})$ is continuous on \mathcal{W} (with probability one) and

$$\int_{\mathbb{R}^k} (1 + \|\mathbf{x}\|^2) \sup_{\|\mathbf{u}\| \leq \delta} f_{\mathbf{x}|\mathbf{w}}(\mathbf{x}|\mathbf{w} + \mathbf{u}) d\rho(\mathbf{x}) < \infty \quad \text{for every } \mathbf{w} \in \mathcal{W}.$$

In addition, the function $(\varepsilon, \mathbf{w}) \mapsto f_{\varepsilon|\mathbf{w}}(\varepsilon|\mathbf{w})$ is continuous and bounded and the function

$$(\mathbf{x}, \mathbf{w}) \mapsto \int_{\mathbb{R}} \sup_{|u| + \|\mathbf{u}\| \leq \delta} f_{\varepsilon|\mathbf{w}}(\varepsilon - \mathbf{x}'\boldsymbol{\theta}_0 - \gamma_0(\mathbf{w}) + u|\mathbf{w} + \mathbf{u})d\varepsilon$$

is bounded.

- (ii) $\mathbb{E}[\|\mathbf{x}\|^4] + \sup_{\mathbf{w} \in \mathcal{W}} (1 + \mathbb{E}[\|\mathbf{x}\|^4|\mathbf{w}])f_{\mathbf{w}}(\mathbf{w}) < \infty$.
- (iii) With probability one, $\mathbb{V}[\mathbf{x}|\mathbf{w}]$ is positive definite.

6 Conclusion

This paper has developed bandwidth robust distribution theory and bootstrap-based inference procedures for a broad class of convex pairwise difference estimators. Our theoretical work is based on small bandwidth asymptotics and carefully leverages convexity. The theory is illustrated by means of three prominent examples. In addition to expanding the scope of small bandwidth asymptotics, our results lay the groundwork for several promising avenues of future research. First, our methods could be generalized to develop bandwidth selection based on higher-order stochastic expansions. Second, they could be expanded to allow for pairwise difference estimators based on generated regressors, a class of estimators that sometimes arises in the context of control function and related econometric methods. Third, when the objective function is smooth, plug-in variance estimation could be developed as an alternative to bootstrap inference. Finally, our current results do not cover settings where the objective function is sufficiently non-smooth to result in non-Gaussian distributional approximations. We plan to investigate these research directions in upcoming work.

References

- AHN, H., H. ICHIMURA, J. L. POWELL, AND P. A. RUUD (2018): “Simple Estimators for Invertible Index Models,” *Journal of Business & Economic Statistics*, 36, 1–10.
- AHN, H. AND J. L. POWELL (1993): “Semiparametric Estimation of Censored Selection Models with a Nonparametric Selection Mechanism,” *Journal of Econometrics*, 58, 3–29.
- ARADILLAS-LOPEZ, A. (2012): “Pairwise-Difference Estimation of Incomplete Information Games,” *Journal of Econometrics*, 168, 120–140.
- ARADILLAS-LOPEZ, A., B. E. HONORÉ, AND J. L. POWELL (2007): “Pairwise Difference Estimation with Nonparametric Control Variables,” *International Economic Review*, 48, 1119–1158.
- BLUNDELL, R. W. AND J. L. POWELL (2004): “Endogeneity in Semiparametric Binary Response Models,” *Review of Economic Studies*, 71, 655–679.

- CATTANEO, M. D., R. K. CRUMP, AND M. JANSSON (2010): “Robust Data-Driven Inference for Density-Weighted Average Derivatives,” *Journal of the American Statistical Association*, 105, 1070–1083.
- (2013): “Generalized Jackknife Estimators of Weighted Average Derivatives (with Discussions and Rejoinder),” *Journal of the American Statistical Association*, 108, 1243–1268.
- (2014a): “Small Bandwidth Asymptotics for Density-Weighted Average Derivatives,” *Econometric Theory*, 30, 176–200.
- (2014b): “Bootstrapping Density-Weighted Average Derivatives,” *Econometric Theory*, 30, 1135–1164.
- CATTANEO, M. D., M. H. FARRELL, M. JANSSON, AND R. P. MASINI (2025a): “Higher-Order Refinements of Small Bandwidth Asymptotics for Density-Weighted Average Derivative Estimators,” *Journal of Econometrics*, forthcoming.
- CATTANEO, M. D. AND M. JANSSON (2018): “Kernel-Based Semiparametric Estimators: Small Bandwidth Asymptotics and Bootstrap Consistency,” *Econometrica*, 86, 955–995.
- (2022): “Average Density Estimators: Efficiency and Bootstrap Consistency,” *Econometric Theory*, 38, 1140–1174.
- CATTANEO, M. D., M. JANSSON, AND K. NAGASAWA (2024): “Bootstrap-Assisted Inference for Generalized Grenander-type Estimators,” *Annals of Statistics*, 52, 1509–1533.
- CATTANEO, M. D., M. JANSSON, AND W. K. NEWEY (2018): “Alternative Asymptotics and the Partially Linear Model with Many Regressors,” *Econometric Theory*, 34, 277–301.
- CATTANEO, M. D., J. M. KLUSOWSKI, AND W. G. UNDERWOOD (2025b): “Inference with Mondrian Random Forests,” *arXiv preprint arXiv:2310.09702*.
- HEYDE, C. C. AND B. M. BROWN (1970): “On the Departure from Normality of a Certain Class of Martingales,” *Annals of Mathematical Statistics*, 41, 2161–2165.
- HJORT, N. L. AND D. POLLARD (1993): “Asymptotics for Minimisers of Convex Processes,” *arXiv preprint arXiv:1107.3806*.
- HONG, H. AND M. SHUM (2010): “Pairwise-Difference Estimation of a Dynamic Optimization Model,” *Review of Economic Studies*, 77, 273–304.
- HONORÉ, B. E. (1992): “Trimmed LAD and Least Squares Estimation of Truncated and Censored Regression Models with Fixed Effects,” *Econometrica*, 60, 533–565.
- HONORÉ, B. E., E. KYRIAZIDOU, AND C. UDRY (1997): “Estimation of Type 3 Tobit Models using Symmetric Rimming and Pairwise Comparisons,” *Journal of Econometrics*, 76, 107–128.

- HONORÉ, B. E. AND J. L. POWELL (1994): “Pairwise Difference Estimators of Censored and Truncated Regression Models,” *Journal of Econometrics*, 64, 241–278.
- (2005): “Pairwise Difference Estimators for Nonlinear Models,” in *Identification and Inference for Econometric Models: Essays in Honor of Thomas Rothenberg*, ed. by D. W. K. Andrews and J. H. Stock, Cambridge University Press, 520–553.
- JOCHMANS, K. (2013): “Pairwise-Comparison Estimation with Non-parametric Controls,” *Econometrics Journal*, 16, 340–372.
- KYRIAZIDOU, E. (1997): “Estimation of a Panel Data Sample Selection Model,” *Econometrica*, 65, 1335–1364.
- MATSUSHITA, Y. AND T. OTSU (2021): “Jackknife Empirical Likelihood: Small Bandwidth, Sparse Network and High-Dimensional Asymptotics,” *Biometrika*, 108, 661–674.
- POLLARD, D. (1991): “Asymptotics for Least Absolute Deviation Regression Estimators,” *Econometric Theory*, 7, 186–199.
- POWELL, J. L. (1994): “Estimation of Semiparametric Models,” in *Handbook of Econometrics, Volume IV*, ed. by R. F. Engle and D. L. McFadden, Elsevier, 2443–2521.
- POWELL, J. L., J. H. STOCK, AND T. M. STOKER (1989): “Semiparametric Estimation of Index Coefficients,” *Econometrica*, 57, 1403–1430.
- SCHUCANY, W. R. AND J. P. SOMMERS (1977): “Improvement of Kernel Type Density Estimators,” *Journal of the American Statistical Association*, 72, 420–423.
- SHAO, J. AND D. TU (2012): *The Jackknife and Bootstrap*, Springer.
- VAN DER VAART, A. W. (1998): *Asymptotic Statistics*, Cambridge University Press.