

# AN ASSUMPTION-LEAN SKEPTICISM TEST OF INFERENCE VALIDITY FOR DOUBLY ROBUST FUNCTIONALS

LIN LIU<sup>1</sup>, RAJARSHI MUKHERJEE<sup>2</sup>, AND JAMES M. ROBINS<sup>3</sup>

A Skeptic’s Empirical Theory of Statistics (SETS) In this article we study assumption-free tests that can falsify or refute the validity of a nominal  $(1 - \alpha)$  Wald confidence interval of the so-called doubly robust functionals (DR functionals), recently characterized in [Rotnitzky et al. \(2019\)](#). The class of such DR functionals is quite broad, encompassing many parameters of interest in econometrics, such as average treatment effect. The state-of-the-art estimators for DR functionals  $\psi$  are doubly robust machine learning (DRML) estimators  $\hat{\psi}_1$  ([Chernozhukov et al., 2018a](#)), which combine the benefits of the superior predictive performance of modern machine learning algorithms, the decreased bias of doubly robust estimation, and the analytical tractability and relaxed assumptions on the model complexity of sample splitting with cross fitting. However, the associated nominal  $(1 - \alpha)$  Wald confidence interval  $\hat{\psi}_1 \pm z_{\alpha/2} \widehat{\text{s.e.}}(\hat{\psi}_1)$  may still undercover even when the sample size is large, if the bias of the DRML estimator is of the same or even greater order than its standard error.

The assumption-free procedure developed in this article (i) can have the power to detect if the bias of  $\hat{\psi}_1$  is of the same or greater order than its standard error, (ii) can provide a lower confidence limit of the degree of under coverage of the confidence interval  $\hat{\psi}_1 \pm z_{\alpha/2} \widehat{\text{s.e.}}(\hat{\psi}_1)$ , and (iii) are valid under essentially no assumptions whatsoever. Our impressive claims need to be tempered in several important ways. First no test, including ours, of the null hypothesis that the ratio of the bias to its standard error can be consistent [without making additional assumptions (e.g. smoothness or sparsity) that may be incorrect]. Furthermore the above claims, without further restrictions, apply to testing if a certificate or an upper bound of the bias of  $\hat{\psi}_1$ , based on Cauchy Schwarz inequality, is of the same or greater order than its standard error.

Finally simulation experiments are shown to demonstrate how the proposed assumption-free test can be used in data analysis.

**Key words:** Causal Inference, Machine Learning, Assumption-free Inference, Hypothesis Testing, Confidence Intervals, Doubly Robust Functionals, Higher-Order Influence Functions, U-Statistics

## 1. INTRODUCTION

In this paper we consider the problem of nearly assumption-free test of the coverage of confidence intervals for a class of low dimensional functionals of high dimensional multivariate distributions that are of great interest in

---

1. Department of Epidemiology, Harvard University, [lil490@mail.harvard.edu](mailto:lil490@mail.harvard.edu).

2. Department of Biostatistics, Harvard University, [ram521@mail.harvard.edu](mailto:ram521@mail.harvard.edu).

3. Department of Epidemiology and Biostatistics, Harvard University, [robins@hsph.harvard.edu](mailto:robins@hsph.harvard.edu).

THIS WORK IS BASED ON A TALK GIVEN BY THE LAST AUTHOR IN THE CONFERENCE OF “WHITNEY NEWKEY’S CONTRIBUTIONS TO ECONOMETRICS” HELD AT MIT IN 2019.

economics and statistics. Specifically, we consider the class of mixed bias functionals (MBFs) studied by [Rotnitzky et al. \(2019\)](#). This class includes the mean of a response missing at random, the average treatment effect under strong ignorability, the weighted average derivative, and the expected conditional covariance. In fact, the class strictly includes both (i) the class of mean-square continuous functionals that can be written as an expectation of a affine functional of a conditional expectation studied by [Chernozhukov et al. \(2018c\)](#) and the class of functionals studied by [Robins et al. \(2008\)](#).

Formally the set up is as follows. We observe  $N$  i.i.d. copies of the data vector  $O = (W, X)$  drawn from some probability distribution  $P_\theta$  belonging to a model  $\mathcal{M} = \{P_\theta; \theta \in \Theta\}$  where  $\Theta$  is a large, infinite dimensional parameter space and  $X$  is a  $d$ -dimensional random vector with compact support. We assume  $\mathcal{M}$  is locally nonparametric in the sense that the tangent space for the model at each  $\theta \in \Theta$  is equal to  $L_2(P_\theta)$ . To avoid technicalities that are orthogonal to the issues studied in this paper we assume that observed data  $O$  is bounded with probability 1 (see Condition [W](#) for further discussion). We consider functionals (i.e. parameters)  $\psi : \theta \mapsto \psi(\theta)$  that possess a (first order) influence function<sup>1</sup> ([Bickel et al., 1998; van der Vaart, 2002, 1998](#))  $IF_{1,\psi}(\theta) = if_{1,\psi}(O; \theta)$  (and thus a positive semiparametric variance bound ([Newey, 1990](#))) and are contained in the mixed bias (MB) class of functionals of [Rotnitzky et al. \(2019\)](#) defined as follows.

**Definition 1.1** (Definition of a mixed bias (MB) functional (Definition 1 of [Rotnitzky et al. \(2019\)](#))).  *$\psi(\theta)$  is an MB functional if, for each  $\theta \in \Theta$  there exists  $b : x \mapsto b(x) \in \mathcal{B}$  and  $p : x \mapsto p(x) \in \mathcal{P}$  such that (i)  $\theta = (b, p, \theta_{\setminus(b,p)})$  and  $\Theta = \mathcal{B} \times \mathcal{P} \times \Theta_{\setminus(\mathcal{B},\mathcal{P})}$  and (ii) for any  $\theta, \theta'$*

$$(1.1) \quad \psi(\theta) - \psi(\theta') + E_\theta [IF_{1,\psi}(\theta')] = E_\theta [S_{bp}(b(X) - b'(X))(p(X) - p'(X))]$$

where  $S_{bp} \equiv s_{bp}(O)$  and  $o \mapsto s_{bp}(o)$  is a known function that does not depend on  $\theta$  or  $\theta'$  satisfying either  $P_\theta(S_{bp} \geq 0) = 1$  or  $P_\theta(S_{bp} \leq 0) = 1$ .

To understand the importance of the MB property, let  $P_N$  be the empirical expectation operator and suppose  $\theta'$  were an estimate of  $\theta$  from a separate sample (which we regard as fixed, i.e. non-random). It then follows that the one step estimator  $\psi(\theta') + P_N [IF_{1,\psi}(\theta')]$  is doubly robust ([Bang and Robins, 2005; Robins and Rotnitzky, 2001; Scharfstein et al., 1999a,b](#)). That is, by [\(1.1\)](#), it is unbiased for  $\psi(\theta)$  under  $P_\theta$  if either  $b = b'$  or  $p = p'$ . Because of this fact we will use the term MB functional and doubly robust (DR) functional interchangeably, as is done in much of the current literature. [For convenience and with no effect on our final results, we have added the variation independence of  $b$  and  $p$  to the above definition of [Rotnitzky et al. \(2019\)](#).] To ease notation, we will restrict consideration to DR functionals for which  $P_\theta(S_{bp} \geq 0) = 1$ . For a DR functional  $\psi^\dagger(\theta)$  of substantive interest for which  $P_\theta(S_{bp} \leq 0) = 1$ , we will instead analyze the DR functional  $\psi(\theta) = -\psi^\dagger(\theta)$ .

Let  $W = (Y, A)$ . Below are some examples of DR functionals.

<sup>1</sup>When writing “influence function” only, we mean “first order influence function” in this paper.

- (1) The counterfactual mean  $E_\theta[Y(a=1)]$  of  $Y$  when  $\{0, 1\}$ -valued  $A$  is set to 1 (under strong ignorability) is identified by the DR functional  $\psi^\dagger(\theta) = E_\theta[b(X)]$  with  $b(x) = E_\theta[Y|X=x, A=1]$ ,  $p(x) = 1/E_\theta[A|X=x]$ , and  $S_{bp} = -A$ . Here  $p(x)$  is the inverse of the propensity score  $\pi(x) = E_\theta[A|X=x]$ . Because  $S_{bp} = -A$ , we instead analyze the DR functional

$$\psi(\theta) = -\psi^\dagger(\theta) = -E_\theta[b(X)] = -E_\theta[Y(a=1)]$$

for which  $S_{bp} = A$ . We will use  $\psi(\theta) = -E_\theta[b(X)] = -E_\theta[Y(a=1)]$  as a running example below.

- (2) The expected conditional covariance

$$\psi(\theta) = E_\theta[(Y - b(X))(A - p(X))]$$

with  $b(x) = E_\theta[Y|X=x]$ ,  $p(x) = E_\theta[A|X=x]$  is a DR functional (equivalently, a MB functional) with  $S_{bp} = 1$ .

The state-of-the-art estimators for DR functionals such as  $\psi(\theta) = -E_\theta[b(X)] = -E_\theta[Y(a=1)]$  are doubly robust machine learning (DRML) estimators  $\hat{\psi}_{cf,1}$  (Chernozhukov et al., 2018a), which combine the benefits of the superior predictive performance of modern machine learning algorithms, the decreased bias of doubly robust estimation (Bang and Robins, 2005; Robins and Rotnitzky, 2001; Scharfstein et al., 1999a,b), and the analytical tractability and the decreased bias of sample splitting with cross fitting. More specifically, the data is randomly divided into two (or more) samples - the estimation sample of size  $n$  and the training or nuisance sample of size  $n_{tr} = N - n$  with  $1 - c > n/N > c$  for some  $c \in (0, 1)$ . To simplify the exposition we will take  $c$  to be  $1/2$ . Machine learning estimators  $\hat{b}(x)$  of the outcome regression  $b(x)$  and  $\hat{p}(x)$  of the inverse propensity score  $p(x)$  are fit using the training sample data. The one step estimator of  $\psi(\theta)$  is  $\hat{\psi}_1 = \psi(\hat{\theta}) + P_n[IF_{1,\psi}(\hat{\theta})] = P_n[-\hat{b}(X) - A\hat{p}(X)(Y - \hat{b}(X))]$  of  $\psi(\theta)$  is computed from the estimation sample treating the machine learning regression estimators as fixed functions. Here  $\psi(\hat{\theta})$  is the plug in estimator of  $\psi(\theta)$ . This approach is required because the ML estimates of the regression functions generally have unknown statistical properties and, in particular, may not lie in a Donsker class – a condition often needed for valid inference when sample splitting is not employed. Under conditions given in Theorem 2.3 below, the efficiency lost due to sample splitting can be recovered by cross fitting. The cross-fit estimator  $\hat{\psi}_{cf,1}$  averages  $\hat{\psi}_1$  with its ‘twin’ obtained by exchanging the roles of the estimation and training sample. However, even using DRML estimators, the associated nominal  $(1 - \alpha)$  Wald confidence interval  $\hat{\psi}_{cf,1} \pm z_{\alpha/2}\widehat{se}(\hat{\psi}_{cf,1})$  may still undercover  $\psi(\theta)$  even when the sample size is large, if the bias of the DRML estimator  $\hat{\psi}_{cf,1}$  is of the same or even greater order than its standard error of order  $n^{-1/2}$ . In the following we will say that a nominal  $(1 - \alpha)$  confidence interval is valid if its coverage probability under repeated sampling is greater than or equal to  $(1 - \alpha)$ .

In general a data analyst who reports the Wald confidence interval  $\hat{\psi}_{cf,1} \pm z_{\alpha/2}\widehat{se}(\hat{\psi}_{cf,1})$  for a DR functional justifies its validity by (i) imposing complexity-reducing assumptions on the nuisance functions  $b$  and  $p$  (e.g. in terms of smoothness (Farrell et al., 2019; Robins et al., 2009b) or sparsity (Bradic et al., 2019; Smucler et al., 2019)) and then

(ii) appealing to theorems wherein the asymptotic validity of the Wald confidence interval has been proved under these assumptions. However, these assumptions may both be incorrect and empirically untestable about the true but unknown functions  $b$  and  $p$  generating the data. In particular, [Robins and Ritov \(1997\)](#), [Robins et al. \(2009b\)](#) and [Ritov et al. \(2014\)](#) proved that, if the parameter space  $\Theta$  does not entail any restriction on the complexity of  $b$  and  $p$  [beyond the assumption that the propensity score  $1/p(X)$  is bounded away from zero, which is needed to insure a positive semiparametric variance bound], then (i) the bias of any estimator of  $\psi(\theta) = -E_\theta[b(X)] = -E_\theta[Y(a=1)]$  may be order 1 [i.e., there does not exist an estimator of  $\psi(\theta)$  such that the bias of estimator converges to zero as  $N \rightarrow \infty$  for all  $\theta \in \Theta$ ]; and (ii) there does not exist a uniform (in  $\theta \in \Theta$ ) consistent test of the composite null hypothesis that  $b$  and  $p$  satisfy given smoothness or sparsity assumptions. From (ii) it follows that no matter the sample size  $N$ , for any  $\alpha^\dagger$ -level test of this composite null, there exists  $\theta = (b, p, \theta_{\setminus(b,p)})$  in the complement to the null such that the power under  $P_\theta$  is less than or equal to  $\alpha^\dagger$ . From (i) it follows that there exists no valid nominal  $(1 - \alpha)$  confidence interval whose median length converges to 0 for all  $\theta \in \Theta$  as the sample size  $N \rightarrow \infty$ . In this sense, in the absence of restrictive untestable assumptions on  $b$  and  $p$ , no meaningful inference concerning  $\psi(\theta)$  is possible.

Even so, in this paper, we will take a skeptic's stance and impose no restrictions on the complexity of  $b$  and  $p$ ! However, we adopt a novel *Skeptical Empiric Theory of Statistics* described below and show that, under this theory, something important remains to be said. The *Skeptic's Theory of Statistics* is motivated by the contrast between two competing social norms. The first is the social norm most of our parents taught us and the second is the skeptic's social norm.

- Parental Social Norm: If You Don't Have Anything Positive to Contribute, Don't Go Criticizing Others.
- Skeptic's Social Norm: Not Having Anything Positive to Contribute Does Not Relieve You of your Duty to Criticize What Others Say.

As we saw above, because we do not impose untestable complexity reducing assumptions on  $b$  and  $p$ , we have nothing to contribute if we follow parental social norms. However, in this paper, we adopt the *Skeptic's Theory of Statistics* that criticizes, where possible, the data analyst who reports a state of the art  $(1 - \alpha)$  Wald confidence interval  $\hat{\psi}_{cf,1} \pm z_{\alpha/2} \widehat{se}(\hat{\psi}_{cf,1})$  as valid. However, our critique will have to be stronger than simply informing the analyst that one can prove that his interval will not be valid for all  $\theta \in \Theta$  if his complexity-reducing assumptions are incorrect, as he will likely respond that he believes his assumptions to be reasonable and likely true under the law  $P_\theta$  actually generating the data.

Therefore the goal of the *Skeptic's Theory of Statistics* is to empirically demonstrate, when possible, that the bias of the analyst's estimator  $\hat{\psi}_{cf,1}$  (under the unknown true law  $P_\theta$ ) is of the same or greater order than its standard error. This would prove to him not only that his assumptions must be incorrect but, more importantly, that his Wald interval cannot be valid. However, for parameters in the DR functional (equivalently MB functional) class this is only possible under a the so-called faithfulness assumption given in Section 3.2.1. Heuristically, faithfulness is

the assumption that near perfect cancelling of the bias of two separate components of the the total bias of  $\widehat{\psi}_{\text{cf},1}$  will essentially never occur.

If we do not assume faithfulness, we can consider the less ambitious goal of demonstrating to the analyst, when possible, that his complexity reducing assumptions are incorrect [now without being able to ever empirically prove the bias of his estimator is of the order of its standard error or greater]. We shall see, this goal is often achievable for parameters in the DR functional class. If successfully achieved, the analyst would then have to admit that he can no longer justify his earlier claim of validity for his state-of-the-art confidence interval. The approach described here is one of being in dialogue with current practices and practitioners. This is not surprising, as it is the statements of the practitioners that the skeptic is critiquing.

To be concrete, suppose, as is often the case, the analyst justifies the validity of  $\widehat{\psi}_1 \pm z_{\alpha/2} \widehat{\text{s.e.}}(\widehat{\psi}_1)$  and thus its cross-fit version  $\widehat{\psi}_{\text{cf},1} \pm z_{\alpha/2} \widehat{\text{s.e.}}(\widehat{\psi}_{\text{cf},1})$  by first proving that, under his complexity reducing assumptions, the Cauchy Schwarz (CS) bias functional

$$\text{CSBias}_{\theta}(\widehat{\psi}_1) = \mathbb{E}_{\theta} \left[ A \left( \widehat{b}(X) - b(X) \right)^2 \right]^{1/2} \mathbb{E}_{\theta} \left[ A \left( \widehat{p}(X) - p(X) \right)^2 \right]^{1/2}$$

is  $o(n^{-1/2})$  conditional on the training sample data  $\mathbf{O}_{\text{tr}}$ <sup>2</sup> (and thus also on the functions  $\widehat{b}, \widehat{p}$  computed from the training sample) and then noting the CS bias upper bounds the absolute conditional bias

$$\left| \mathbb{E}_{\theta} \left[ A \left( \widehat{b}(X) - b(X) \right) \left( \widehat{p}(X) - p(X) \right) \right] \right|$$

of  $\widehat{\psi}_1$ . It then follows if we empirically show that Cauchy Schwarz bias  $\text{CSBias}_{\theta}(\widehat{\psi}_1)$  exceeds some given multiple  $\delta > 0$ , e.g.  $\delta = 3/4$ , times  $\widehat{\psi}_1$ 's conditional standard error of order  $n^{-1/2}$ , then we have falsified the basis of the analyst's claim that his Wald confidence interval is valid.

To give a glimpse of our idea, we shall construct  $\alpha^{\dagger}$ -level tests of the null hypothesis  $\text{CSBias}_{\theta}(\widehat{\psi}_1) < \text{s.e.}_{\theta}(\widehat{\psi}_1)\delta$  via higher order influence function tests and estimators introduced in [Mukherjee et al. \(2017\)](#); [Robins et al. \(2008, 2017\)](#). In particular, the null hypothesis  $\text{CSBias}_{\theta}(\widehat{\psi}_1) < \text{s.e.}_{\theta}(\widehat{\psi}_1)\delta$ , if true, guarantees that the nominal  $(1 - \alpha)$  Wald confidence interval covers  $\psi(\theta)$  with probability at least

$$\Phi(z_{\alpha/2} - \delta) - \Phi(-z_{\alpha/2} - \delta)$$

in large sample. For example, when  $\alpha = 0.10, \delta = 3/4, \Phi(z_{\alpha/2} - \delta) - \Phi(-z_{\alpha/2} - \delta) \approx 80\%$ .

In this paper, we will not review the definition of higher order influence functions. We refer the interested readers to works such as [Robins et al. \(2008, 2016\)](#); [van der Vaart \(2014\)](#). Higher order influence functions tests and estimators are particular higher order U-statistics. Our tests are constructed without: i) refitting, modifying, or even having knowledge of the ML algorithms that have been employed to compute  $\widehat{b}, \widehat{p}$  from the training sample and ii)

<sup>2</sup>In this paper, essentially all expectations and probabilities are to be understood as being conditional on the training sample. Hence we can and do omit this conditioning event in our notation. The randomness in the training sample will also be ignored.

requiring any assumptions at all (aside from a few standard, quite weak assumptions given later) – in particular, without making any assumptions about the smoothness or sparsity of the true outcome regression  $b$  or propensity score  $1/p$ .

We will show there is an unavoidable limitation to what can be achieved with our or any other method. No test, including ours, of the null hypothesis  $\text{CSBias}_\theta(\hat{\psi}_1) < \text{s.e.}_\theta(\hat{\psi}_1)\delta$  can be uniformly (in  $\theta$ ) consistent [without making additional complexity reducing assumptions]. Thus, when our  $\alpha^\dagger$ -level test rejects this null hypothesis for  $\alpha^\dagger$  small, we have strong evidence that the null hypothesis is false; nonetheless when the test does not reject, we cannot conclude that there is good evidence that  $\text{CSBias}_\theta(\hat{\psi}_1) < \text{s.e.}_\theta(\hat{\psi}_1)\delta$ , no matter how large the sample size. This reflects the fact that, because we make (essentially) no assumptions,  $\text{CSBias}_\theta(\hat{\psi}_1)$  may be order 1 and thus  $n^{1/2}$  times greater than  $\text{s.e.}_\theta(\hat{\psi}_1)$  under the law  $P_\theta$  generating the data and yet our  $\alpha^\dagger$ -level test, since not uniformly consistent, may be powerless to detect this fact!

*Overview of Our Approach.* All claims made in this subsection will be formally proved later in the paper. DRML estimators are based on the influence function of the DR functional  $\psi(\theta)$  (van der Vaart, 2002, 1998). Our proposed approach begins by computing a second order influence function estimator  $\hat{\mathbb{IF}}_{22,k}$  (defined later) of the estimable part of the conditional bias  $E_\theta[\hat{\psi}_1 - \psi(\theta)|\mathbf{O}_{\text{tr}}]$  of  $\hat{\psi}_1$  given  $\mathbf{O}_{\text{tr}}$ . Our bias corrected estimator is  $\hat{\psi}_{2,k} \equiv \hat{\psi}_1 - \hat{\mathbb{IF}}_{22,k}$ . The statistic  $\hat{\mathbb{IF}}_{22,k}$  is a second-order U-statistic that depends on a choice of  $k$  (with  $k = o(n^2)$  for reasons explained in Comment 3.4), a vector of basis functions  $\bar{Z}_k \equiv \bar{z}_k(X) \equiv (z_1(X), \dots, z_k(X))$  of the  $d$ -dimensional vector of potential confounders  $X$  and an estimator  $\hat{\Sigma}_k^{-1}$  computed from the training sample of the inverse expected weighted outer product  $\Sigma_k^{-1} := \{E_\theta[\lambda(X)\bar{z}_k(X)\bar{z}_k(X)^\top]\}^{-1}$ , with weight defined as

$$\lambda(x) \equiv \lambda(x; \theta) := E_\theta[S_{bp}|X = x].$$

Henceforth we denote  $\hat{\mathbb{IF}}_{22,k}$  and  $\hat{\psi}_{2,k}$  as  $\hat{\mathbb{IF}}_{22,k}(\hat{\Sigma}_k^{-1})$  and  $\hat{\psi}_{2,k}(\hat{\Sigma}_k^{-1})$  to make explicit the dependence on the estimator  $\hat{\Sigma}_k^{-1}$ . The oracle estimators  $\hat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})$  and  $\hat{\psi}_{2,k}(\Sigma_k^{-1})$  replace  $\hat{\Sigma}_k^{-1}$  by the true  $\Sigma_k^{-1}$ . The oracle estimators are generally infeasible. (An exception is the case of semisupervised data; that is a data set in which the number  $N$  of subjects with complete data on  $(W, X)$  is many fold less than the number of subjects for whom data on just the covariates  $X$  and the statistic  $S_{bp}$  are available. In that case, assuming the subjects with complete data are a random sample of all subjects, we can first estimate  $\Sigma_k$  by its empirical covariance matrix based on all subjects; and then treat  $\Sigma_k$ , and thus  $\Sigma_k^{-1}$ , as known and equal to the estimate in an analysis based on the  $N$  subjects with complete data.) The estimators  $\hat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})$ ,  $\hat{\psi}_{2,k}(\Sigma_k^{-1})$ ,  $\hat{\mathbb{IF}}_{22,k}(\hat{\Sigma}_k^{-1})$  and  $\hat{\psi}_{2,k}(\hat{\Sigma}_k^{-1})$  will be asymptotically normal conditional on the training sample when, as in our asymptotic set-up,  $k = k(n) \rightarrow \infty$  and  $k = o(n^2)$  as  $n \rightarrow \infty$ . Furthermore, the variance of  $\hat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})$  and  $\hat{\psi}_{2,k}(\Sigma_k^{-1})$  are at most of order  $1/n$ .

The fraction of conditional bias of  $\hat{\psi}_1$  corrected by  $\hat{\mathbb{IF}}_{22,k}(\hat{\Sigma}_k^{-1})$  depends critically on (i) the choice of  $k$ , (ii) the accuracy of the estimator  $\hat{\Sigma}_k^{-1}$  of  $\Sigma_k^{-1}$ , and (iii) the particular  $k$ -vector of (basis) functions  $\bar{Z}_k \equiv \bar{z}_k(X)$  chosen from



a much larger, possibly countably infinite, dictionary of candidate functions. Data adaptive choice of  $\bar{z}_k(x)$  is an important research topic but it is beyond the scope of the current paper.

In this paper we shall always choose  $k$  to be less than the sample size  $n = N/2$  of the estimation sample for the following three reasons:  $k < n$  is necessary (i) for  $\hat{\mathbb{I}}\mathbb{F}_{22,k}$ 's standard error of order  $k^{1/2}/n$  to be smaller than the order  $n^{-1/2}$  of the standard error of  $\hat{\psi}_1$ , thereby creating the possibility of detecting, for any given  $\delta > 0$ , that the (absolute value) of the ratio of the estimable part of the bias of  $\hat{\psi}_1$  to its standard error exceeds  $\delta$ , when the sample size  $n$  is sufficiently large, (ii) to be able to use the inverse of the sample covariance matrix  $\hat{\Sigma}_k^{-1}$  computed from the training sample to estimate  $\Sigma_k^{-1}$  accurately without imposing, possibly incorrect, additional smoothness or sparsity assumptions required for accurate estimation when  $k \geq n$ , and iii) for bias-corrected confidence intervals centered at  $\hat{\psi}_{2,k} \equiv \hat{\psi}_1 - \hat{\mathbb{I}}\mathbb{F}_{22,k}$  to have length approximately equal to  $n^{-1/2}$  centered at  $\hat{\psi}_1$ .

The importance of (i) is two-fold. First, if  $|\mathbb{E}_\theta[\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})]| \geq \text{s.e.}_\theta(\hat{\psi}_1)\delta$ , then under faithfulness, the analysts nominal  $(1 - \alpha)$  Wald confidence interval is invalid. Further  $|\mathbb{E}_\theta[\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})]| \geq \text{s.e.}_\theta(\hat{\psi}_1)\delta$  implies  $\text{CSBias}_\theta(\hat{\psi}_1) \geq \text{s.e.}_\theta(\hat{\psi}_1)\delta$ . Thus, even without the faithfulness assumption,  $|\mathbb{E}_\theta[\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})]| \geq \text{s.e.}_\theta(\hat{\psi}_1)\delta$  falsifies the analyst's justification for the validity of his nominal  $(1 - \alpha)$  Wald interval, if his justification needed  $\text{CSBias}_\theta(\hat{\psi}_1)$  be  $o(n^{-1/2})$ .

*Organization of the paper.* The remainder of the paper is organized as follows. Section 2 reviews needed results on DR functionals obtained by [Rotnitzky et al. \(2019\)](#), the state of the art DRML estimators, and the statistical properties of these estimators.

In Section 3, we study the theoretical statistical properties of the following: the oracle estimator  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})$  of the bias of the DRML estimator  $\hat{\psi}_1$ , the oracle biased-corrected estimator  $\hat{\psi}_{2,k}(\Sigma_k^{-1})$ , the  $\alpha^\dagger$ -level oracle test  $\hat{\chi}_{2,k}(\Sigma_k^{-1}; z_{\alpha^\dagger/2}, \delta)$  of the null hypothesis  $|\mathbb{E}_\theta[\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})]| \geq \text{s.e.}_\theta(\hat{\psi}_1)\delta$ .

In Section 4, we no longer assume  $\Sigma_k^{-1}$  is known and replace it by inverse of the sample covariance matrix  $\hat{\Sigma}_k^{-1}$  computed from the training sample. We propose a test statistic  $\hat{\chi}_{3,k}(\hat{\Sigma}_k^{-1}; z_{\alpha^\dagger/2}, \delta)$  based on a third order U-statistic and derive additional conditions under which this new test recovers the statistical properties of the oracle test  $\hat{\chi}_{2,k}(\Sigma_k^{-1}; z_{\alpha^\dagger/2}, \delta)$ .

In Section 5, we generalize the results in Section 4 by relaxing the additional conditions required for  $\hat{\chi}_{3,k}(\hat{\Sigma}_k^{-1}; z_{\alpha^\dagger/2}, \delta)$  to have good statistical properties via even higher order U-statistic. We also propose an "early stopping" procedure that can save computational cost while still protecting the level of the test.

In Section 6 we present results of simulation studies to evaluate the finite sample performance of our methods. Section 7 concludes with discussion of some open problems. Many of the technical details are deferred to the Appendix.

Before going into the main part of our paper, we define the observed data structure and collect some frequently used notation.

**1.1. Observation schemes and notations.** As mentioned above, we observe  $N$  i.i.d. copies of data  $O = (W, X)$  drawn from some unknown probability distribution  $P_\theta$  belonging to a locally nonparametric model

$$\mathcal{M} = \{P_\theta; \theta = (b, p, \theta_{\setminus(b,p)}) \in \Theta = \mathcal{B} \times \mathcal{P} \times \Theta_{\setminus(\mathcal{B}, \mathcal{P})}\}$$

parameterized by the parameter  $\theta = (b, p, \theta_{\setminus(b,p)})$  with  $b, p, \theta_{\setminus(b,p)}$  variation independent. Here the maps  $b : x \mapsto b(x) \in \mathcal{B}$  and  $p : x \mapsto p(x) \in \mathcal{P}$  have range contained in  $\mathbb{R}$ .

To avoid extraneous technical issues, we assume  $W$  is bounded with  $P_\theta$ -probability 1 and that  $X$  is a  $d$ -dimensional random vector with compact support whose density is bounded away from 0 and  $\infty$  on its support. We use  $E_\theta[\cdot]$ ,  $\text{var}_\theta[\cdot]$  and etc. to denote the expectation, variance, and etc. with respect to  $P_\theta$ . The data is randomly divided into an estimation sample and a training (equivalently, nuisance) sample of size  $n = N/2$ . To avoid notational clutter, all expectations, variances, and probabilities are conditional on the training sample unless otherwise stated. Hence when we write  $E_\theta[X]$  we mean  $E_\theta[X | \mathbf{O}_{\text{tr}}]$  with  $\mathbf{O}_{\text{tr}}$  the training sample data. For a (random) vector  $V$ ,  $\|V\|_\theta \equiv E_\theta[V^{\otimes 2}]^{1/2} = E_\theta[V^\top V]^{1/2}$  denotes its  $L_2(P_\theta)$  norm conditioning on the training sample,  $\|V\| \equiv (V^{\otimes 2})^{1/2} = (V^\top V)^{1/2}$  denotes its  $\ell_2$  norm and  $\|V\|_\infty$  denotes its  $L_\infty(P_\theta)$  norm. For any matrix  $A$ ,  $\|A\|$  will be reserved for its operator norm. Given an integer  $k$ , and a random vector  $\bar{Z}_k = \bar{z}_k(X)$ ,  $\Pi_\theta[\cdot | \bar{Z}_k]$  denotes the population linear projection operator onto the space spanned by  $\bar{Z}_k$  conditioning on the training sample, and  $\Pi_\theta^\perp[\cdot | \bar{Z}_k] = (I - \Pi_\theta)[\cdot | \bar{Z}_k]$  is the projection onto the orthogonal complement of  $\bar{Z}_k$  in the Hilbert space  $L_2(P_{\theta, X})$  where  $P_{\theta, X}$  denotes the law of the  $X$  under  $\theta$ . That is, for a random variable  $V$ ,

$$(1.2) \quad \Pi_\theta[V | \bar{Z}_k] = \bar{Z}_k^\top (E_\theta[\bar{Z}_k \bar{Z}_k^\top])^{-1} E_\theta[\bar{Z}_k V], \Pi_\theta^\perp[V | \bar{Z}_k] = V - \Pi_\theta[V | \bar{Z}_k].$$

We adopt the following condition, which is really a notational convention.

The following common asymptotic notations are used throughout the paper:  $x \lesssim y$  (equivalently  $x = O(y)$  or  $y = \Omega(x)$ ) denotes that there exists some constant  $C > 0$  such that  $x \leq Cy$ ,  $x \asymp y$  (equivalently  $x = \Theta(y)$ ) means there exist some constants  $c_1 > c_2 > 0$  such that  $c_2|y| \leq |x| \leq c_1|y|$ .  $x = o(y)$  or  $y = \omega(x)$  or  $y \gg x$  or  $x \ll y$  is equivalent to  $\lim_{x, y \rightarrow \infty} \frac{x}{y} = 0$ . For a random variable  $X_n$  with law  $P$  possibly depending on the sample size  $n$ ,  $X_n = O_P(a_n)$  denotes that  $X_n/a_n$  is bounded in  $P$ -probability, and  $X_n = o_P(a_n)$  means that  $\lim_{n \rightarrow \infty} P(|X_n/a_n| \geq \epsilon) = 0$  for every positive real number  $\epsilon$ .

## 2. DR FUNCTIONALS AND THE STATISTICAL PROPERTIES OF DRML ESTIMATORS

We begin this section by reviewing some useful results on DR functionals (see Definition 1.1) established in [Rotnitzky et al. \(2019\)](#). First, the theorem below characterizes the (nonparametric) influence functions of DR functionals.

**Theorem 2.1** (Theorem 1 of [Rotnitzky et al. \(2019\)](#)). *Suppose  $\psi(\theta)$  for  $\theta = (b, p, \theta_{\setminus(b,p)}) \in \Theta \equiv \mathcal{B} \times \mathcal{P} \times \Theta_{\setminus(\mathcal{B}, \mathcal{P})}$  is a DR functional (equivalently MB functional) according to Definition 1.1; further the regularity Condition 1 of [Rotnitzky et al.](#)*



(2019) holds. Then there exist a statistic  $S_0$  and maps  $h \mapsto m_1(O, h)$  for  $h \in \mathcal{B}$  and  $h \mapsto m_2(O, h)$  for  $h \in \mathcal{P}$  independent of  $\theta$  satisfying the following:

- the maps  $h \in \mathcal{B} \mapsto \mathbb{E}_\theta[m_1(O, h)]$  and  $h \in \mathcal{P} \mapsto \mathbb{E}_\theta[m_2(O, h)]$  are linear;

- 

$$\psi(\theta) = \mathbb{E}_\theta[m_1(O, b)] + \mathbb{E}_\theta[S_0] = \mathbb{E}_\theta[m_2(O, p)] + \mathbb{E}_\theta[S_0];$$

- and

$$(2.1) \quad \text{IF}_{1,\psi}(\theta) = \mathcal{H}(b, p) - \psi(\theta)$$

$$\text{where } \mathcal{H}(b, p) := S_{bp}b(X)p(X) + m_1(O, b) + m_2(O, p) + S_0.$$

Furthermore

$$(2.2) \quad \begin{cases} \mathbb{E}_\theta[S_{bp}h(X)p(X) + m_1(O, h)] = 0 \text{ for all } h \in \mathcal{B}, \\ \mathbb{E}_\theta[S_{bp}b(X)h(X) + m_2(O, h)] = 0 \text{ for all } h \in \mathcal{P}. \end{cases}$$

In addition, with  $\mathbf{P}_{\theta,X}$  denoting the marginal law of  $X$  under  $\mathbf{P}_\theta$ , suppose for all  $h \in L_2(\mathbf{P}_{\theta,X})$ ,  $m_1(O, h)$  and  $m_2(O, h)$  are in  $L_2(\mathbf{P}_\theta)$ . Then the maps  $h \in \mathcal{B} \mapsto m_1(O, h)$  and  $h \in \mathcal{P} \mapsto m_2(O, h)$  are linear.

Under very slightly stronger regularity conditions on the maps  $m_1$  and  $m_2$ , the converse of Theorem 2.1 also holds, together with several additional results.

**Theorem 2.2** (Theorem 2 of [Rotnitzky et al. \(2019\)](#)). Suppose the regularity Condition 1 of [Rotnitzky et al. \(2019\)](#) holds. Suppose a functional  $\psi(\theta)$  has influence function  $\text{IF}_{1,\psi}(\theta)$  of the form (2.1) for  $m_1$  and  $m_2$  such that all  $\theta = (b, p, \theta_{\setminus(b,p)}) \in \Theta$ , the maps  $h \mapsto \mathbb{E}_\theta[m_1(O, h)]$  and  $h \mapsto \mathbb{E}_\theta[m_2(O, h)]$  for  $h \in L_2(\mathbf{P}_{\theta,X})$  are continuous and linear with Riesz representers  $\mathcal{R}_1(X) \equiv \mathcal{R}_1(X; \theta)$  and  $\mathcal{R}_2(X) \equiv \mathcal{R}_2(X; \theta)$ <sup>3</sup> respectively. Moreover, suppose that  $b(X)$ ,  $p(X)$ ,  $\lambda(X)b(X)$  and  $\lambda(X)p(X)$  belong to  $L_2(\mathbf{P}_{\theta,X})$  for all  $\theta \in \Theta$ , where  $\lambda(X) = \mathbb{E}_\theta[S_{bp}|X]$ . Then, for all  $\theta = (b, p, \theta_{\setminus(b,p)}) \in \Theta$

- (1) the identity (1.1) holds for any  $\theta'$  and hence  $\psi(\theta)$  is an DR functional (equivalently MB functional);
- (2) for all  $h \in L_2(\mathbf{P}_{\theta,X})$

$$\mathbb{E}_\theta[S_{bp}bh + m_2(O, h)] = 0 \quad \text{and} \quad \mathbb{E}_\theta[S_{bp}hp + m_1(O, h)] = 0.$$

- (3)  $b(X) = -\mathcal{R}_2(X)/\lambda(X)$  and  $p(X) = -\mathcal{R}_1(X)/\lambda(X)$ .

- (4)  $\psi(\theta) = \mathbb{E}_\theta[m_2(O, p) + S_0] = \mathbb{E}_\theta[m_1(O, b) + S_0] = \mathbb{E}_\theta[-S_{bp}bp + S_0]$ .

- (5)  $b$  and  $p$  are the minimizers of the following minimization problems

$$b = \arg \min_{h \in L_2(\mathbf{P}_{\theta,X})} \mathbb{E}_\theta \left[ S_{bp} \frac{h^2}{2} + m_2(O, h) \right], \quad p = \arg \min_{h \in L_2(\mathbf{P}_{\theta,X})} \mathbb{E}_\theta \left[ S_{bp} \frac{h^2}{2} + m_1(O, h) \right].$$

<sup>3</sup>The Riesz representer  $\mathcal{R}(X)$  of a continuous linear functional  $h \mapsto \mathbb{E}[m(O, h)]$  for  $h \in L_2(\mathbf{P}_X)$  is, by definition, the function of  $X$  satisfying  $\mathbb{E}[m(O, h)] = \mathbb{E}[\mathcal{R}(X)h(X)]$ ; see [Chernozhukov et al. \(2018b,c\)](#); [Rotnitzky et al. \(2019\)](#).

We now give the influence functions  $\text{IF}_{1,\psi}(\theta)$  and Riesz representers for several DR functionals.

**Example 1** (Some examples of Riesz representers for DR functionals).

- (1)  $\psi(\theta) = -\mathbb{E}_\theta[b(X)] = -\mathbb{E}_\theta[Y(a=1)]$ . Then  $\text{IF}_{1,\psi}(\theta) = \mathcal{H}(b, p) - \psi(\theta)$  with  $\mathcal{H}(b, p) = -(-Ab(X)p(X) + b(X) + AYp(X))$ :  $S_{bp} = A$  and  $\lambda(X) = p(X)^{-1}$ ,  $m_1(O, b) = -b(X)$ ,  $m_2(O, p) = -AYp(X)$  and  $S_0 = 0$ . Then  $\mathcal{R}_1(X) = -1$  and  $\mathcal{R}_2(X) = -b(X)/p(X)$  since

$$\begin{cases} \mathbb{E}_\theta[m_1(O, h)] = \mathbb{E}_\theta[-h(X)], \\ \mathbb{E}_\theta[m_2(O, h)] = \mathbb{E}_\theta[-AYh(X)] = \mathbb{E}_\theta[-\{p(X)\}^{-1}b(X)h(X)]. \end{cases}$$

- (2)  $\psi(\theta) = \mathbb{E}_\theta[(Y - b(X))(A - p(X))]$  – expected conditional covariance between  $A$  and  $Y$  given  $X$  with  $b(X) = \mathbb{E}_\theta[Y|X]$ ,  $p(X) = \mathbb{E}_\theta[A|X]$ :  $\text{IF}_{1,\psi}(\theta) = \mathcal{H}(b, p) - \psi(\theta)$  with  $\mathcal{H}(b, p) = b(X)p(X) - Ab(X) - Yp(X) + AY$ ,  $S_{bp} = 1$  and  $\lambda(X) = 1$ ,  $m_1(O, b) = -Ab(X)$ ,  $m_2(O, p) = -Yp(X)$ , and  $S_0 = AY$ . Then  $\mathcal{R}_1(X) = -p(X)$  and  $\mathcal{R}_2(X) = -b(X)$  since

$$\begin{cases} \mathbb{E}_\theta[m_1(O, h)] = \mathbb{E}_\theta[-Ah(X)] = \mathbb{E}_\theta[-p(X)h(X)], \\ \mathbb{E}_\theta[m_2(O, h)] = \mathbb{E}_\theta[-Yh(X)] = \mathbb{E}_\theta[-b(X)h(X)]. \end{cases}$$

The following algorithm defines the DRML estimators  $\hat{\psi}_1$  and  $\hat{\psi}_{\text{cf},1}$  of a DR functional  $\psi(\theta)$  satisfying the suppositions of Theorem 2.2.

- (i) The  $N$  study subjects are randomly split into two parts: an estimation sample of size  $n$  and a training (nuisance) sample of size  $n_{tr} = N - n$  with  $n/N \approx 1/2$ . Without loss of generality we shall assume that  $i = 1, \dots, n$  corresponds to the estimation sample.
- (ii) Since  $b$  and  $p$  are generally unknown, their estimates  $\hat{b}$  and  $\hat{p}$  are constructed from the training sample data using ML methods, such as deep neural networks and define  $\hat{\theta} := (\hat{b}, \hat{p}, \theta_{\setminus(b,p)})$ . [Note that unlike  $(\hat{b}, \hat{p})$ ,  $\theta_{\setminus(b,p)}$  in the definition of  $\hat{\theta}$  is not estimated from data.]
- (iii) Denote  $\mathbb{IF}_1 \equiv \mathbb{IF}_1(\theta) := \frac{1}{n} \sum_{i=1}^n \text{IF}_{1,\psi}(\theta)$  and  $\hat{\mathbb{IF}}_1 \equiv \mathbb{IF}_1(\hat{\theta}) := \frac{1}{n} \sum_{i=1}^n \text{IF}_{1,\psi}(\hat{\theta})$ . Then

$$\begin{aligned} \hat{\psi}_1 &= \hat{\mathbb{IF}}_1 + \psi(\hat{\theta}) = \frac{1}{n} \sum_{i=1}^n \mathcal{H}_i(\hat{b}, \hat{p}) \\ &= \frac{1}{n} \sum_{i=1}^n \left( S_{bp,i} \hat{b}(X_i) \hat{p}(X_i) + m_1(O_i, \hat{b}) + m_2(O_i, \hat{p}) + S_{0,i} \right) \end{aligned}$$

from  $n$  subjects in the estimation sample and

$$\hat{\psi}_{\text{cf},1} = \frac{1}{2} \left( \hat{\psi}_1 + \overline{\hat{\psi}}_1 \right)$$

where  $\overline{\hat{\psi}}_1$  is  $\hat{\psi}_1$  but with the training and estimation samples reversed. [Note that  $\hat{\psi}_1$  does not depend on  $\theta_{\setminus(b,p)}$  and depends on the training sample data only through the functions  $\hat{b}$  and  $\hat{p}$ .]

Then the following theorem provides conditional and unconditional asymptotic properties of  $\widehat{\psi}_1$ .

**Theorem 2.3.** *Under the conditions of Theorem 2.2, conditional on the training sample,  $\widehat{\psi}_1$  is asymptotically normal with conditional bias*

$$(2.3) \quad \text{Bias}_\theta(\widehat{\psi}_1) := \mathbb{E}_\theta [\widehat{\psi}_1 - \psi(\theta)] = \mathbb{E}_\theta \left[ S_{bp} \left\{ b(X) - \widehat{b}(X) \right\} \{p(X) - \widehat{p}(X)\} \right].$$

If further, a)  $\text{Bias}_\theta(\widehat{\psi}_1)$  is  $o(n^{-1/2})$  and b)  $\widehat{b}(x)$  and  $\widehat{p}(x)$  converge to  $b(x)$  and  $p(x)$  in  $L_2(\mathbb{P}_\theta)$ , then

(1)

$$\begin{aligned} \widehat{\psi}_1 - \psi(\theta) &= n^{-1} \sum_{i=1}^n \text{IF}_{1,\psi,i}(\theta) + o(n^{-1/2}) \\ \widehat{\psi}_{\text{cf},1} - \psi(\theta) &= N^{-1} \sum_{i=1}^N \text{IF}_{1,\psi,i}(\theta) + o(N^{-1/2}). \end{aligned}$$

Further  $n^{1/2}(\widehat{\psi}_1 - \psi(\theta))$  converges conditionally and unconditionally to a normal distribution with mean zero;  $\widehat{\psi}_{\text{cf},1}$  is a regular, asymptotically linear estimator; i.e.,  $N^{1/2}(\widehat{\psi}_{\text{cf},1} - \psi(\theta))$  converges unconditionally to a normal distribution with mean zero and variance equal to the semiparametric variance bound  $\text{var}_\theta [\text{IF}_{1,\psi}(\theta)]$  (Bickel et al., 1998; Newey, 1990).

(2) The  $(1 - \alpha)$  nominal Wald confidence intervals (CIs)

$$(2.4) \quad \text{CI}_\alpha(\widehat{\psi}_1) := \widehat{\psi}_1 \pm z_{\alpha/2} \widehat{\text{s.e.}}[\widehat{\psi}_1]$$

$$(2.5) \quad \text{CI}_\alpha(\widehat{\psi}_{\text{cf},1}) := \widehat{\psi}_{\text{cf},1} \pm z_{\alpha/2} \widehat{\text{s.e.}}[\widehat{\psi}_{\text{cf},1}]$$

are valid  $(1 - \alpha)$  unconditional asymptotic confidence interval for  $\psi(\theta)$ . Here  $\widehat{\text{s.e.}}[\widehat{\psi}_1] = \{\widehat{\text{var}}[\widehat{\psi}_1]\}^{1/2}$  with

$$\begin{aligned} \widehat{\text{var}}[\widehat{\psi}_1] &= \frac{1}{n^2} \sum_{i=1}^n \left( \text{IF}_{1,\psi,i}(\widehat{\theta}) - \frac{1}{n} \sum_{i=1}^n \text{IF}_{1,\psi,i}(\widehat{\theta}) \right)^2 \\ \widehat{\text{var}}[\widehat{\psi}_{\text{cf},1}] &= \frac{1}{4} \{ \widehat{\text{var}}[\widehat{\psi}_1] + \widehat{\text{var}}[\widehat{\psi}_1] \}. \end{aligned}$$

The proof of Theorem 2.3 is straightforward and can be found in Chernozhukov et al. (2018a), Smucler et al. (2019) or Farrell et al. (2019, Theorem 3). However, in reality, it is often not enough just to invoke Theorem 2.3 and claim the validity of  $\text{CI}_\alpha(\widehat{\psi}_1)$ . First, owing to the fact that we make no complexity reducing assumptions on  $b$  and  $p$  it follows that because of their asymptotic nature, there is no finite sample size  $n$  at which any test could empirically reject either the hypothesis  $\text{Bias}_\theta(\widehat{\psi}_1) = o(n^{-1/2})$  or the hypothesis that  $n^{1/2}(\widehat{\psi}_1 - \psi(\theta))$  is asymptotically normal with mean zero. Rather, as discussed in Section 1, our interest, instead, lies in testing and rejecting hypotheses such as, at the observed sample size  $n$ , the actual conditional coverage of the nominal  $(1 - \alpha)$  two-sided interval  $\widehat{\psi}_1 \pm z_{\alpha/2} \widehat{\text{s.e.}}[\widehat{\psi}_1]$  under the distribution  $\mathbb{P}_\theta$  that generated the data is less than a fraction  $\varrho < 1$  of its nominal coverage. Second, since the

ML estimators  $\widehat{b}(x)$  and  $\widehat{p}(x)$  have unknown statistical properties and we are interested in developing assumption-free statistical procedure relying only on empirical evidence, our inferential statements regard the training sample as fixed rather than random. In particular, the randomness referred to in subsequent sections is that of the estimation sample. In fact, our inferences will rely on being in ‘asymptopia’, but only to be able to posit that, at our sample size of  $n$ , the quantiles of the finite sample distribution of a conditionally asymptotically normal statistic (e.g.  $\widehat{\mathbb{F}}_{22,k}(\Sigma_k^{-1})$  defined later) are close to the quantiles of a normal. Non-asymptotic version of our proposed procedure is an interesting and important open problem that requires future efforts: e.g. developing procedures that also estimate constants in concentration or deviation type inequalities.

### 3. THE ORACLE PROCEDURE WHEN $\Sigma_k^{-1}$ IS KNOWN

It is useful to first consider an oracle procedure that would only be implementable if  $\Sigma_k^{-1}$  were known. We will focus on such situation in this section. Such oracle procedure bears practical values when one has semisupervised data with sample size much greater than  $N$ , nonetheless only containing information on  $(S_{bp}, X)$  but not on the component of  $W$  excluding  $S_{bp}$ .

**3.1. Truncated parameters and their unbiased estimators.** Though there does not exist uniformly (in  $\theta$ ) consistent test of the null hypothesis that the bias  $\text{Bias}_\theta(\widehat{\psi}_1)$  of  $\widehat{\psi}_1$  for estimating  $\psi(\theta)$  is smaller than a fraction  $\delta$  of its standard error  $\text{s.e.}_\theta(\widehat{\psi}_1)$ , we can instead consider the bias of  $\widehat{\psi}_1$  as an estimator of the truncated parameter  $\widetilde{\psi}_k(\theta)$  of  $\psi(\theta)$  (Robins et al. (2008, Section 3.2)), defined as follows:

**Definition 3.1** (The truncated parameter  $\widetilde{\psi}_k(\theta)$ , truncation bias  $\text{TB}_k(\theta)$ , and the projected bias  $\text{Bias}_{\theta,k}(\widehat{\psi}_1)$ ). *Given estimators  $\widehat{b} \in \mathcal{B}$  and  $\widehat{p} \in \mathcal{P}$  of  $b$  and  $p$  computed from the training sample and some  $k$  dimensional basis functions  $\bar{z}_k(x)$  in  $L_2(\mathbb{P}_{\theta,X})$ :*

- Define the (conditional) truncated parameter of a DR functional  $\psi(\theta)$  as

$$\widetilde{\psi}_k(\theta) := \mathbb{E}_\theta \left[ \mathcal{H} \left( \widetilde{b}_{k,\theta}, \widetilde{p}_{k,\theta} \right) \right],$$

where  $\widetilde{b}_{k,\theta}$  and  $\widetilde{p}_{k,\theta}$  are defined as follows. Consider the working linear models

$$\begin{cases} b_k^*(X; \bar{\zeta}_{b,k}) = \widehat{b}(X) + \bar{\zeta}_{b,k}^\top \bar{z}_k(X), \\ p_k^*(X; \bar{\zeta}_{p,k}) = \widehat{p}(X) + \bar{\zeta}_{p,k}^\top \bar{z}_k(X). \end{cases}$$

Define  $\widetilde{\zeta}_b(\theta)$  and  $\widetilde{\zeta}_p(\theta)$  as the solutions to the following system of equations

$$\begin{cases} 0 = \mathbb{E}_\theta \left[ \frac{\partial \mathcal{H} (b_k^*(X; \bar{\zeta}_{b,k}), p^*(X; \bar{\zeta}_{p,k}))}{\partial \bar{\zeta}_{b,k}} \right] = \mathbb{E}_\theta [S_{bp} p_k^*(X; \bar{\zeta}_{p,k}) \bar{z}_k(X) + m_1(O, \bar{z}_k)] \\ 0 = \mathbb{E}_\theta \left[ \frac{\partial \mathcal{H} (b_k^*(X; \bar{\zeta}_{b,k}), p^*(X; \bar{\zeta}_{p,k}))}{\partial \bar{\zeta}_{p,k}} \right] = \mathbb{E}_\theta [S_{bp} b_k^*(X; \bar{\zeta}_{b,k}) \bar{z}_k(X) + m_2(O, \bar{z}_k)]. \end{cases}$$

where, for  $j = 1, 2$ ,  $m_j(O, \bar{z}_k) = (m_j(O, z_1), \dots, m_j(O, z_k))$ . Hence

$$(3.1) \quad \begin{cases} \tilde{\zeta}_{b,k}(\theta) = -\Sigma_k^{-1} \mathbb{E}_\theta \left[ \left( S_{bp} \widehat{b}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) \right) \right], \\ \tilde{\zeta}_{p,k}(\theta) = -\Sigma_k^{-1} \mathbb{E}_\theta \left[ \left( S_{bp} \widehat{p}(X) \bar{z}_k(X) + m_1(O, \bar{z}_k) \right) \right] \end{cases}$$

where

$$\Sigma_k := \mathbb{E}_\theta \left[ S_{bp} \bar{z}_k(X) \bar{z}_k(X)^\top \right] \equiv \mathbb{E}_\theta \left[ \lambda(X) \bar{z}_k(X) \bar{z}_k(X)^\top \right].$$

Now define

$$\begin{cases} \tilde{b}_{k,\theta}(X) := b_k^*(X; \tilde{\zeta}_{b,k}(\theta)) \equiv \widehat{b}(X) + \tilde{\zeta}_{b,k}(\theta)^\top \bar{z}_k(X), \\ \tilde{p}_{k,\theta}(X) := p_k^*(X; \tilde{\zeta}_{p,k}(\theta)) \equiv \widehat{p}(X) + \tilde{\zeta}_{p,k}(\theta)^\top \bar{z}_k(X). \end{cases}$$

- Define the difference between  $\tilde{\psi}_k(\theta)$  and  $\psi(\theta)$  as the truncation bias  $\text{TB}_k(\theta)$ :

$$(3.2) \quad \text{TB}_k(\theta) := \tilde{\psi}_k(\theta) - \psi(\theta) = \mathbb{E}_\theta \left[ S_{bp} \left( \tilde{b}_{k,\theta}(X) - b(X) \right) \left( \tilde{p}_{k,\theta}(X) - p(X) \right) \right].$$

- Define the bias of  $\widehat{\psi}_1$  as an estimator of  $\tilde{\psi}_k(\theta)$  as

$$(3.3) \quad \text{Bias}_{\theta,k}(\widehat{\psi}_1) := \mathbb{E}_\theta \left[ \widehat{\psi}_1 - \tilde{\psi}_k(\theta) \right].$$

By definition,  $\text{Bias}_{\theta,k}(\widehat{\psi}_1) = \text{Bias}_\theta(\widehat{\psi}_1) - \text{TB}_k(\theta)$ . Combining equation (3.2), we have  $\tilde{\psi}_k(\theta) - \psi(\theta) \equiv \text{Bias}_\theta(\widehat{\psi}_1) - \text{Bias}_{\theta,k}(\widehat{\psi}_1)$ . Note that  $\text{TB}_k(\theta)$  is not consistently estimable without imposing smoothness/sparsity assumptions on the functions  $b$  and  $p$ . But  $\text{Bias}_{\theta,k}(\widehat{\psi}_1)$  can be unbiasedly estimated by the following oracle (i.e.  $\Sigma_k^{-1}$  known) second-order U-statistic:

$$(3.4) \quad \widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1}) := \frac{1}{n(n-1)} \sum_{1 \leq i_1 \neq i_2 \leq n} \widehat{\mathbb{IF}}_{22,k,\bar{i}_2}(\Sigma_k^{-1})$$

where

$$(3.5) \quad \widehat{\mathbb{IF}}_{22,k,\bar{i}_2}(\Sigma_k^{-1}) := \left[ \mathcal{E}_{\widehat{b},m_2}(\bar{z}_k)(O) \right]_{i_1}^\top \Sigma_k^{-1} \left[ \mathcal{E}_{\widehat{p},m_1}(\bar{z}_k)(O) \right]_{i_2}.$$

Here  $\mathcal{E}_{1,b'}(h)$  and  $\mathcal{E}_{2,p'}(h)$  are linear functionals defined as: given any  $h \in L_2(\mathbb{P}_{\theta,X})$ ,  $b' \in \mathcal{B}$  and  $p' \in \mathcal{P}$

$$(3.6) \quad \mathcal{E}_{b',m_2}(h)(O) := S_{bp} b'(X) h(X) + m_2(O, h)$$

and

$$(3.7) \quad \mathcal{E}_{p',m_1}(h)(O) := S_{bp} p'(X) h(X) + m_1(O, h).$$

When the range of  $h$  is multidimensional (e.g.  $h = \bar{z}_k$ ), the linear functionals  $\mathcal{E}_{1,b'}$  and  $\mathcal{E}_{2,p'}$  act on  $h$  component-wise.

For the examples covered in Example 1,  $\mathcal{E}_{\widehat{b},m_2}(\bar{z}_k)(O)$  and  $\mathcal{E}_{\widehat{p},m_1}(\bar{z}_k)(O)$  are

(1) Example 1(1):

$$\begin{cases} \mathcal{E}_{\hat{b}, m_2}(\bar{z}_k)(O) = \hat{A}\hat{b}(X)\bar{z}_k(X) - AY\bar{z}_k(X), \\ \mathcal{E}_{\hat{p}, m_1}(\bar{z}_k)(O) = A\hat{p}(X)\bar{z}_k(X) - \bar{z}_k(X). \end{cases}$$

(2) Example 1(2):

$$\begin{cases} \mathcal{E}_{\hat{b}, m_2}(\bar{z}_k)(O) = \hat{b}(X)\bar{z}_k(X) - Y\bar{z}_k(X), \\ \mathcal{E}_{\hat{p}, m_1}(\bar{z}_k)(O) = \hat{p}(X)\bar{z}_k(X) - A\bar{z}_k(X). \end{cases}$$

The following lemma proves that (i)  $\hat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})$  is an unbiased estimator of  $\text{Bias}_{\theta,k}(\hat{\psi}_1)$  and (ii)  $\text{Bias}_{\theta,k}(\hat{\psi}_1)$  is the expected product of the  $L_2(P_\theta)$  projections of the weighted residuals

$$\begin{aligned} \Delta_{1,\hat{b}} &:= \lambda(X)^{1/2} \left( \hat{b}(X) - b(X) \right), \\ \Delta_{2,\hat{p}} &:= \lambda(X)^{1/2} \left( \hat{p}(X) - p(X) \right) \end{aligned}$$

onto the subspace spanned by the weighted basis functions  $\bar{v}_k(X) := \lambda(X)^{1/2}\bar{z}_k(X)$ :

**Lemma 3.1.** Define  $\beta_{\hat{b},k} := \mathbb{E}_\theta \left[ \mathcal{E}_{\hat{b}, m_2}(\bar{z}_k)(O) \right]^\top \Sigma_k^{-1}$  and  $\beta_{\hat{p},k} := \mathbb{E}_\theta \left[ \mathcal{E}_{\hat{p}, m_1}(\bar{z}_k)(O) \right]^\top \Sigma_k^{-1}$ . Under the conditions of Theorem 2.2, we have

$$\begin{aligned} \text{Bias}_{\theta,k}(\hat{\psi}_1) &= \mathbb{E}_\theta \left[ \hat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1}) \right] \\ (3.8) \quad &= \mathbb{E}_\theta \left[ \Pi_\theta \left[ \Delta_{1,\hat{b}} | \bar{v}_k(X) \right] \Pi_\theta \left[ \Delta_{2,\hat{p}} | \bar{v}_k(X) \right] \right] = \beta_{\hat{b},k}^\top \Sigma_k \beta_{\hat{p},k} \end{aligned}$$

and

$$(3.9) \quad \text{TB}_k(\theta) := \tilde{\psi}_k(\theta) - \psi(\theta) = \text{Bias}_\theta(\hat{\psi}_1) - \text{Bias}_{\theta,k}(\hat{\psi}_1) = \mathbb{E}_\theta \left[ \Pi_\theta^\perp \left[ \Delta_{1,\hat{b}} | \bar{v}_k(X) \right] \Pi_\theta^\perp \left[ \Delta_{2,\hat{p}} | \bar{v}_k(X) \right] \right].$$

Furthermore, define the bias corrected estimator  $\hat{\psi}_{2,k}(\Sigma_k^{-1}) := \hat{\psi}_1 - \hat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})$ . Then

$$\begin{aligned} \text{Bias}_\theta(\hat{\psi}_{2,k}(\Sigma_k^{-1})) &\equiv \mathbb{E}_\theta \left[ \hat{\psi}_{2,k}(\Sigma_k^{-1}) - \psi(\theta) \right] = \text{TB}_k(\theta), \\ (3.10) \quad &\text{and thus } \mathbb{E}_\theta \left[ \hat{\psi}_{2,k}(\Sigma_k^{-1}) \right] = \psi(\theta) + \text{TB}_k(\theta) = \tilde{\psi}_k(\theta). \end{aligned}$$



*Proof.* First, we show the first part of equation (3.8).

$$\begin{aligned}
\text{TB}_k(\theta) &= \mathbb{E}_\theta \left[ S_{bp} \left( \tilde{b}_{k,\theta}(X) - b(X) \right) \left( \tilde{p}_{k,\theta}(X) - p(X) \right) \right] \\
&= \mathbb{E}_\theta \left[ S_{bp} \left\{ \hat{b}(X) - b(X) - \mathbb{E}_\theta \left[ (S_{bp} \hat{b}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k)) \right]^\top \Sigma_k^{-1} \bar{z}_k(X) \right\} \right. \\
&\quad \left. \times \left\{ \hat{p}(X) - p(X) - \mathbb{E}_\theta \left[ (S_{bp} \hat{p}(X) \bar{z}_k(X) + m_1(O, \bar{z}_k)) \right]^\top \Sigma_k^{-1} \bar{z}_k(X) \right\} \right] \\
&= \mathbb{E}_\theta \left[ S_{bp} (\hat{b}(X) - b(X)) (\hat{p}(X) - p(X)) \right] \\
&\quad - \mathbb{E}_\theta \left[ S_{bp} (\hat{b}(X) - b(X)) \bar{z}_k(X) \right]^\top \Sigma_k^{-1} \mathbb{E}_\theta [S_{bp} (\hat{p}(X) - p(X)) \bar{z}_k(X)] \\
&\quad - \mathbb{E}_\theta [S_{bp} (\hat{p}(X) - p(X)) \bar{z}_k(X)]^\top \Sigma_k^{-1} \mathbb{E}_\theta [S_{bp} (\hat{b}(X) - b(X)) \bar{z}_k(X)] \\
&\quad + \mathbb{E}_\theta [S_{bp} (\hat{p}(X) - p(X)) \bar{z}_k(X)]^\top \Sigma_k^{-1} \mathbb{E}_\theta [S_{bp} (\hat{b}(X) - b(X)) \bar{z}_k(X)] \\
&= \underbrace{\mathbb{E}_\theta [S_{bp} (\hat{b}(X) - b(X)) (\hat{p}(X) - p(X))]}_{\equiv \text{Bias}_\theta(\hat{\psi}_1)} \\
&\quad - \mathbb{E}_\theta \left[ S_{bp} (\hat{b}(X) - b(X)) \bar{z}_k(X) \right]^\top \Sigma_k^{-1} \mathbb{E}_\theta [S_{bp} (\hat{p}(X) - p(X)) \bar{z}_k(X)]
\end{aligned}$$

where in the third equality we apply Theorem 2.2(2). It follows that

$$\text{Bias}_{\theta,k}(\hat{\psi}_1) = \mathbb{E}_\theta \left[ S_{bp} (\hat{b}(X) - b(X)) \bar{z}_k(X) \right]^\top \Sigma_k^{-1} \mathbb{E}_\theta [S_{bp} (\hat{p}(X) - p(X)) \bar{z}_k(X)].$$

Again by Theorem 2.2(2),

$$\begin{aligned}
\mathbb{E}_\theta [\mathcal{E}_{\hat{b},m_2}(\bar{z}_k)(O)] &= \mathbb{E}_\theta [S_{bp} \hat{b}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k)] = \mathbb{E}_\theta [S_{bp} (\hat{b}(X) - b(X)) \bar{z}_k(X)] \\
\mathbb{E}_\theta [\mathcal{E}_{\hat{p},m_1}(\bar{z}_k)(O)] &= \mathbb{E}_\theta [S_{bp} \hat{p}(X) \bar{z}_k(X) + m_1(O, \bar{z}_k)] = \mathbb{E}_\theta [S_{bp} (\hat{p}(X) - p(X)) \bar{z}_k(X)]
\end{aligned}$$

and thus

$$\mathbb{E}_\theta [\hat{\mathbb{IIF}}_{22,k}(\Sigma_k^{-1})] = \mathbb{E}_\theta [\mathcal{E}_{\hat{b},m_2}(\bar{z}_k)(O)]^\top \Sigma_k^{-1} \mathbb{E}_\theta [\mathcal{E}_{\hat{p},m_1}(\bar{z}_k)(O)] \equiv \text{Bias}_{\theta,k}(\hat{\psi}_1).$$

In terms of the second part of equation (3.8):

$$\begin{aligned}
\text{Bias}_{\theta,k}(\widehat{\psi}_1) &= \mathbb{E}_\theta \left[ S_{bp}(\widehat{b}(X) - b(X))\bar{z}_k(X) \right]^\top \Sigma_k^{-1} \mathbb{E}_\theta [S_{bp}(\widehat{p}(X) - p(X))\bar{z}_k(X)] \\
&= \mathbb{E}_\theta \left[ \lambda(X)(\widehat{b}(X) - b(X))\bar{z}_k(X) \right]^\top \Sigma_k^{-1} \mathbb{E}_\theta [\lambda(X)(\widehat{p}(X) - p(X))\bar{z}_k(X)] \\
&= \mathbb{E}_\theta \left[ \lambda(X)^{1/2}(\widehat{b}(X) - b(X))\bar{v}_k(X) \right]^\top \Sigma_k^{-1} \mathbb{E}_\theta [\lambda(X)^{1/2}(\widehat{p}(X) - p(X))\bar{v}_k(X)] \\
&= \mathbb{E}_\theta \left[ \Delta_{1,\widehat{b}}\bar{v}_k(X) \right]^\top \Sigma_k^{-1} \mathbb{E}_\theta [\Delta_{2,\widehat{p}}\bar{v}_k(X)] \\
&= \mathbb{E}_\theta \left[ \Pi_\theta \left[ \Delta_{1,\widehat{b}}|\bar{v}_k \right] (X)^\top \Pi_\theta [\Delta_{2,\widehat{p}}|\bar{v}_k] (X) \right].
\end{aligned}$$

where in the last line we use the definition of the projection operator  $\Pi_\theta [\cdot|\bar{v}_k(X)]$  (see equation (1.2)).

Then equation (3.9) immediately follows:

$$\begin{aligned}
\text{TB}_k(\theta) &= \text{Bias}_\theta(\widehat{\psi}_1) - \text{Bias}_{\theta,k}(\widehat{\psi}_1) \\
&= \mathbb{E}_\theta \left[ \Delta_{1,\widehat{b}}\Delta_{2,\widehat{p}} \right] - \mathbb{E}_\theta \left[ \Pi_\theta \left[ \Delta_{1,\widehat{b}}|\bar{v}_k \right] (X)^\top \Pi_\theta [\Delta_{2,\widehat{p}}|\bar{v}_k] (X) \right] \\
&= \mathbb{E}_\theta \left[ \Pi_\theta^\perp \left[ \Delta_{1,\widehat{b}}|\bar{v}_k \right] (X) \right]^\top \mathbb{E}_\theta \left[ \Pi_\theta^\perp [\Delta_{2,\widehat{p}}|\bar{v}_k] (X) \right].
\end{aligned}$$

Finally, we prove equation (3.10):

$$\begin{aligned}
\text{Bias}_\theta(\widehat{\psi}_{2,k}(\Sigma_k^{-1})) &= \mathbb{E}_\theta \left[ \widehat{\psi}_{2,k}(\Sigma_k^{-1}) - \psi(\theta) \right] \\
&= \mathbb{E}_\theta \left[ \widehat{\psi}_1 - \psi(\theta) - \widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1}) \right] \\
&= \text{Bias}_\theta(\widehat{\psi}_1) - \text{Bias}_{\theta,k}(\widehat{\psi}_1) = \text{TB}_k(\theta)
\end{aligned}$$

where the third line follows from the definition of  $\text{Bias}_\theta(\widehat{\psi}_1)$  and equation (3.8). ■

**Comment 3.1.** By  $\text{Bias}_{\theta,k}(\widehat{\psi}_1) = \mathbb{E}_\theta \left[ \widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1}) \right] = \beta_{\widehat{b},k}^\top \Sigma_k^{-1} \beta_{\widehat{p},k}$ , we can view  $\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})$  as an unbiased estimator of the bilinear functional  $\beta_{\widehat{b},k}^\top \Sigma_k^{-1} \beta_{\widehat{p},k}$ , where  $\beta_{\widehat{b},k}$  and  $\beta_{\widehat{p},k}$  are the population coefficients of the following linear regressions:

$$\begin{aligned}
\lambda(X)^{1/2}(\widehat{b}(X) - b(X)) &\sim \lambda(X)^{1/2}\bar{z}_k(X), \\
\lambda(X)^{1/2}(\widehat{p}(X) - p(X)) &\sim \lambda(X)^{1/2}\bar{z}_k(X).
\end{aligned}$$

Under the above setup, in Appendix A.1, we relate  $\widehat{\psi}_{2,k}(\Sigma_k^{-1}) - \psi(\theta) = \widehat{\psi}_1 - \widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1}) - \psi(\theta)$  to the second order influence function of the truncated parameter  $\widetilde{\psi}_k(\theta)$  of a DR functional, following Robins et al. (2008, Definition 2.1). In particular, we have the following.

**Proposition 3.1.**

- $\widehat{\psi}_{2,k}(\Sigma_k^{-1}) - \psi(\theta)$  is the second order influence function of  $\widetilde{\psi}_k(\theta)$  under the law  $P_{\widehat{\theta}}$ ;
- $\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})$  is the second order influence function of  $\text{Bias}_{\theta,k}(\widehat{\psi}_1)$  under the law  $P_{\widehat{\theta}}$ .

**Comment 3.2.** For the definition of second order or even higher order influence functions, we refer the readers to [Robins et al. \(2008\)](#) or [Mukherjee et al. \(2017\)](#). Following ([Robins et al., 2008](#), Theorem 2.3), the second order influence function of  $\widetilde{\psi}_k(\theta)$  is “the influence function of the influence function” of  $\widetilde{\psi}_k(\theta)$ . Granted, one can derive the second order influence functions of  $\widetilde{\psi}_k(\theta)$  and  $\text{Bias}_{\theta,k}(\widehat{\psi}_1)$  for DR functionals from first principles (i.e. the definition of influence functions and scores, see [Appendix A.1](#)), but it is intuitively more appealing to consider higher order influence functions as a “bias corrector”.

We introduce the following additional notation for various norms which will be useful for stating the statistical properties of  $\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})$  and  $\widehat{\psi}_{2,k}(\Sigma_k^{-1})$ :

$$\begin{aligned}\mathbb{L}_{\theta,2,\widehat{b}} &:= \left\{ \mathbb{E}_{\theta}[\Delta_{1,\widehat{b}}^2] \right\}^{1/2}, \quad \mathbb{L}_{\theta,2,\widehat{b},k} := \left\{ \mathbb{E}_{\theta}[\Pi_{\theta}[\Delta_{1,\widehat{b}}|\bar{v}_k](X)^2] \right\}^{1/2}, \\ \mathbb{L}_{\theta,2,\widehat{p}} &:= \left\{ \mathbb{E}_{\theta}[\Delta_{1,\widehat{p}}^2] \right\}^{1/2}, \quad \mathbb{L}_{\theta,2,\widehat{p},k} := \left\{ \mathbb{E}_{\theta}[\Pi_{\theta}[\Delta_{1,\widehat{p}}|\bar{v}_k](X)^2] \right\}^{1/2}, \\ \mathbb{L}_{\theta,2,\widehat{\Sigma},k} &:= \|\widehat{\Sigma}_k - \Sigma_k\|.\end{aligned}$$

Henceforth, we further impose the following weak conditions (Condition [W](#)).

**Condition W.**

- (1) All the eigenvalues of  $\Sigma_k$  are bounded away from 0 and  $\infty$
- (2)  $\|\bar{z}_k(x)^{\top} \bar{z}_k(x)\|_{\infty} \leq Bk$  for some constant  $B > 0$ ;
- (3) The observed data  $O = (W, X)$ , the true nuisance functions  $b(X)$  and  $p(X)$ , the estimated nuisance functions  $\widehat{b}(X)$  and  $\widehat{p}(X)$ , and  $\lambda(X)$  are all bounded with  $P_{\theta}$ -probability 1.

**Comment 3.3.** Condition [W](#)(2) holds for Cohen-Vial-Daubechies wavelets, local polynomial partition and B-spline series ([Belloni et al., 2015](#)). But it does not hold in general for polynomial or Fourier series ([Belloni et al., 2015](#)).

The statistical properties of the oracle estimator  $\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})$  and the oracle biased-corrected estimator  $\widehat{\psi}_{2,k}(\Sigma_k^{-1})$  are given by the following theorem.

**Theorem 3.1.** Under the conditions of [Theorem 2.2](#) and Condition [W](#), with  $k, n \rightarrow \infty$ , and  $k = o(n^2)$ , conditional on the training sample, we have

- (1)  $\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})$  is unbiased for  $\text{Bias}_{\theta,k}(\widehat{\psi}_1)$  with variance of order

$$\frac{1}{n} \left( \frac{k}{n} + \mathbb{L}_{\theta,2,\widehat{b},k}^2 + \mathbb{L}_{\theta,2,\widehat{p},k}^2 \right).$$

- (2)  $\frac{\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1}) - \text{Bias}_{\theta,k}(\widehat{\psi}_1)}{\text{s.e.}_{\theta}[\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})]} \xrightarrow{d} N(0, 1)$ , where  $\xrightarrow{d}$  stands for convergence in distribution. Further,  $\text{s.e.}_{\theta}[\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})] := \text{var}_{\theta}^{1/2}[\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})]$  can be estimated by the bootstrap estimator  $\widehat{\text{s.e.}}[\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})] := \widehat{\text{var}}^{1/2}[\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})]$  defined in Appendix A.3 satisfying  $\frac{\widehat{\text{s.e.}}[\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})]}{\text{s.e.}_{\theta}[\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})]} = 1 + o_{\mathbb{P}_{\theta}}(1)$ .
- (3)  $\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1}) \pm z_{\alpha^{\dagger}/2} \widehat{\text{s.e.}}[\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})]$  is a  $(1 - \alpha^{\dagger})$  asymptotic two-sided Wald CI for  $\text{Bias}_{\theta,k}(\widehat{\psi}_1)$  with length of order

$$\frac{1}{\sqrt{n}} \left( \sqrt{\frac{k}{n}} + \mathbb{L}_{\theta,2,\widehat{b},k} + \mathbb{L}_{\theta,2,\widehat{p},k} \right).$$

*Proof.* The variance order of  $\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})$  is proved in Appendix A.2. When  $k = o(n^2)$  and  $k \rightarrow \infty$  as  $n \rightarrow \infty$ , the conditional asymptotic normality of  $\frac{\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1}) - \text{Bias}_{\theta,k}(\widehat{\psi}_1)}{\text{s.e.}_{\theta}[\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})]}$  follows directly from Hoeffding decomposition, with the conditional asymptotic normality of the degenerate second-order U-statistic part implied by Bhattacharya and Ghosh (1992, Corollary 1.2). Appendix A.3 proves that the bootstrap standard error estimators  $\widehat{\text{s.e.}}[\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})]$  satisfy  $\frac{\widehat{\text{s.e.}}[\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})]}{\text{s.e.}_{\theta}[\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})]} = 1 + o_{\mathbb{P}_{\theta}}(1)$ . ■

**Comment 3.4.** When  $k \gtrsim n^2$ , the Gaussian limit of  $\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})$  does not hold (rather one expect a Poisson limit). Further if  $k \gg n^2$ ,  $\text{var}_{\theta}[\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})] = O(\frac{k}{n^2})$  is of order greater than 1, and therefore  $\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})$  cannot consistently estimate  $\text{Bias}_{\theta,k}(\widehat{\psi}_1)$  regardless of its magnitude.

**3.2. Which null hypothesis should be tested?** Recall that in the introduction, we pose the following question:

*Given the training sample used to estimate the nuisance functions  $b$  and  $p$ , can we use empirical information to test the null hypothesis  $\text{NH}_0 : \text{Bias}_{\theta}(\widehat{\psi}_1) = o(n^{-1/2})$ , the violation of which implies under coverage of  $\text{CI}_{\alpha}(\widehat{\psi}_1)$ ?*

In fact, we shall have to settle for statements that are “in dialogue” with current practices and literature. In the current literature, Theorem 2.3 provides the theoretical guarantees on the coverage validity of a nominal  $(1 - \alpha)$  two-sided Wald confidence interval  $\text{CI}_{\alpha}(\widehat{\psi}_1)$  based on a DRML estimator  $\widehat{\psi}_1$ , under the “hypothesis” that  $\text{NH}_0 : \text{Bias}_{\theta}(\widehat{\psi}_1) = o(n^{-1/2})$  i.e. the bias of the DRML estimator  $\widehat{\psi}_1$  for  $\psi(\theta)$  is dominated by the standard error of  $\widehat{\psi}_1$  of order  $n^{-1/2}$ . This “hypothesis” is often proved by imposing untestable complexity-reducing assumptions on the nuisance parameters  $b$  and  $p$ , such as smoothness or sparsity. Moreover, Theorem 2.3 also requires that the training sample size  $n_{\text{tr}}$  is sufficiently large so that the ratio  $\text{Bias}_{\theta}(\widehat{\psi}_1)/\text{s.e.}_{\theta}(\widehat{\psi}_1)$  is close to zero, a consequence of  $\text{Bias}_{\theta}(\widehat{\psi}_1) = o(n^{-1/2})$ .

However a clause in  $\text{NH}_0$  involves an asymptotic statement which cannot be quantified or operationalized by a data analyst with a sample of finite size. We propose the following strategy of turning  $\text{NH}_0$  into a statement that is easy to deal with: whenever a null hypothesis is defined in terms of an asymptotic rate of convergence such as  $o(n^{-1/2})$  in the training sample data, we will (1) ask the authors to specify a positive number  $\delta = \delta(N)$  possibly depending on the actual sample size  $N$  and (2) then operationalize the asymptotic null hypothesis  $\text{Bias}_{\theta}(\widehat{\psi}_1) = o(n^{-1/2})$  as the null

hypothesis  $H_0(\delta)$ . That is, we have the following pair

$$(3.11) \quad \text{NH}_0 : \text{Bias}_\theta(\hat{\psi}_1) = o(n^{-1/2});$$

$$(3.12) \quad H_0(\delta) : \frac{|\text{Bias}_\theta(\hat{\psi}_1)|}{\text{s.e.}_\theta(\hat{\psi}_1)} < \delta.$$

by which we mean that if  $H_0(\delta)$  is rejected (accepted), we, by convention, will declare  $\text{NH}_0$  rejected (accepted). The authors' choice of  $\delta$  depends on the degree of under coverage they are willing to tolerate. For example, as we mentioned in Section 1, if 80% is the minimum actual coverage they would tolerate for a 90% two-sided Wald confidence interval  $\text{CI}_{\alpha=0.05}(\hat{\psi}_1)$  (2.4), then under normality they would select  $\delta = 0.75$ .

Similarly, we have the *surrogate* pair

$$(3.13) \quad \text{NH}_{0,k} : \text{Bias}_{\theta,k}(\hat{\psi}_1) = o(n^{-1/2});$$

$$(3.14) \quad H_{0,k}(\delta) : \frac{|\text{Bias}_{\theta,k}(\hat{\psi}_1)|}{\text{s.e.}_\theta(\hat{\psi}_1)} < \delta.$$

Note that  $H_{0,k}(\delta)$  specifies  $|\text{Bias}_{\theta,k}(\hat{\psi}_1)|$ , *not*  $|\text{Bias}_\theta(\hat{\psi}_1)|$ , to be less than  $\delta$  fraction of  $\text{s.e.}_\theta(\hat{\psi}_1)$ . Given the statistical properties of  $\hat{\mathbb{M}}_{22,k}(\Sigma_k^{-1})$  summarized in Theorem 3.1,  $\text{Bias}_{\theta,k}(\hat{\psi}_1)$  can be unbiasedly estimated by  $\hat{\mathbb{M}}_{22,k}(\Sigma_k^{-1})$ , whose standard error shrinks to zero at rate at most  $1/\sqrt{n}$ . This observation motivates the following oracle test of  $H_{0,k}(\delta)$  (3.14):

$$(3.15) \quad \hat{\chi}_{2,k}(\Sigma_k^{-1}; \varsigma_k, \delta) := \mathbb{1} \left\{ \frac{|\hat{\mathbb{M}}_{22,k}(\Sigma_k^{-1})|}{\widehat{\text{s.e.}}[\hat{\psi}_1]} - \varsigma_k \frac{\widehat{\text{s.e.}}[\hat{\mathbb{M}}_{22,k}(\Sigma_k^{-1})]}{\widehat{\text{s.e.}}[\hat{\psi}_1]} > \delta \right\},$$

for user-specified  $\varsigma_k$  and  $\delta > 0$ . In fact, later we will show in Section 3.3 that  $\hat{\chi}_{2,k}(\Sigma_k^{-1}; \varsigma_k, \delta)$  is a consistent test of  $H_{0,k}(\delta)$  (3.14). For now, we take this conclusion as given and discuss the implication of the test result on the original null hypothesis of interest  $\text{NH}_0$  (3.11) in Sections 3.2.1 and 3.2.2.

**3.2.1. Faithfulness assumption.** Suppose now a data analyst agrees that in reporting  $\text{CI}_\alpha(\hat{\psi}_1)$  (2.4) as a  $(1 - \alpha)$  two sided confidence interval for  $\psi(\theta)$ , their implicit or explicit null hypothesis is  $\text{NH}_0$  (3.11), that is, the bias  $\text{Bias}_\theta(\hat{\psi}_1)$  of their estimator is  $o(n^{-1/2})$ . Further suppose the test  $\hat{\chi}_k^{(2)}(\Sigma_k^{-1}; z_{\alpha^\dagger/2}, \delta)$  rejects  $H_{0,k}(\delta)$  (3.14), equivalently  $\text{NH}_{0,k}$  (3.13). However, rejecting  $\text{NH}_{0,k}$  (3.13) does *not* logically imply rejecting  $\text{NH}_0$  (3.11) in general. One “solution” is to adopt an additional “faithfulness” assumption under which rejection of  $\text{NH}_{0,k}$  does logically imply rejection of  $\text{NH}_0$ .

**Condition Faithfulness.** Given a fixed  $k$ ,  $\frac{\text{Bias}_\theta(\hat{\psi}_1)}{\text{Bias}_{\theta,k}(\hat{\psi}_1)} = 1 + \frac{\text{TB}_k(\theta)}{\text{Bias}_{\theta,k}(\hat{\psi}_1)}$  is not  $o(1)$ .

One might find this assumption rather natural because it holds unless  $\text{TB}_k(\theta)$  and  $\text{Bias}_{\theta,k}(\hat{\psi}_1)$  are of the same order and their leading constants sum to zero, which seems highly unlikely to be the case. In finite sample, we can also operationalize the above asymptotic faithfulness condition by choosing some  $\delta' > 0$  and imposing:

**Condition Faithfulness( $\delta'$ ).** For a given  $k$ ,  $\left| \frac{\text{Bias}_\theta(\widehat{\psi}_1)}{\text{Bias}_{\theta,k}(\widehat{\psi}_1)} \right| = \left| 1 + \frac{\text{TB}_k(\theta)}{\text{Bias}_{\theta,k}(\widehat{\psi}_1)} \right| \geq \delta'$ .

Under Condition **Faithfulness( $\delta'$ )**, rejection of  $H_{0,k}(\delta)$  implies rejection of  $H_0(\delta\delta')$ . If we choose  $\delta' = 0.15$ , Condition **Faithfulness( $\delta'$ )** holds unless  $-1.15 \leq \frac{\text{TB}_k(\theta)}{\text{Bias}_{\theta,k}(\widehat{\psi}_1)} \leq -0.85$ . When we reject  $H_{0,k}(\delta)$  for some large  $\delta$ , say  $\delta = 10$ , we will reject  $H_0(\delta\delta' = 1.5)$ , suggesting that the true asymptotic coverage of a 90% two-sided Wald confidence interval should be lower than 55.6%. To some extent, imposing Condition **Faithfulness** or Condition **Faithfulness( $\delta'$ )** may seem inconsistent with the goal of falsifying the validity of reported Wald confidence intervals without unverifiable assumptions.

**3.2.2. Can we avoid assuming faithfulness?** In this section, we discuss an alternative “solution” that avoids imposing Condition **Faithfulness** or Condition **Faithfulness( $\delta'$ )**. This “solution” was described in Section 1 as the ground of this paper.

We shall assume that the goal in using a machine learning algorithm to predict the nuisance functions  $b(x)$  and  $p(x)$  is to construct predictors  $\widehat{b}(x)$  and  $\widehat{p}(x)$  that (nearly) minimize the conditional mean weighted square error losses  $E_\theta[\lambda(X)\{\widehat{b}(X) - b(X)\}^2]$  and  $E_\theta[\lambda(X)\{\widehat{p}(X) - p(X)\}^2]$  over the set of functions computable by the algorithm. In fact, researchers who use the algorithm described below in Comment 3.5 are explicitly acknowledging this as their goal.

**Comment 3.5** (Training data cross validation squared error loss minimization framework in machine learning). *Training data cross validation squared error loss minimization (Steinwart and Christmann, 2008; Vapnik, 1998, 2013) is one of the most common model fitting methods in machine learning.*

For DR functionals, Theorem 2.2(5) says that the nuisance functions  $b$  and  $p$  happen to be the solutions to the following minimization problems:

$$\begin{cases} b = \arg \min_{h \in L_2(P_{\theta,X})} E_\theta \left[ S_{bp} \frac{h^2}{2} + m_2(O, h) \right], \\ p = \arg \min_{h \in L_2(P_{\theta,X})} E_\theta \left[ S_{bp} \frac{h^2}{2} + m_1(O, h) \right]. \end{cases}$$

As a result, it is natural to minimize the following loss functions to estimated  $b$  and  $p$  separately:

$$(3.16) \quad \widehat{b} = \arg \min_{h \in \mathcal{F}} E_\theta \left[ S_{bp} \frac{h^2}{2} + m_2(O, h) \right],$$

$$(3.17) \quad \widehat{p} = \arg \min_{h \in \mathcal{F}} E_\theta \left[ S_{bp} \frac{h^2}{2} + m_1(O, h) \right].$$

where  $\mathcal{F}$  is some function class chosen by analysts.

At a first sight, it might not be obvious how the recommended losses relate to the squared error loss or weighted squared error loss. However, a simple application of equation (2.2) in Theorem 2.1 immediately reveals that the loss minimization criterion



in equations (3.16) and (3.17) are equivalent to the weighted squared error losses of estimating  $b$  and  $p$ , weighted by  $\lambda(\cdot)$ :

$$\begin{aligned} \frac{1}{2}\mathbb{E}_\theta[\lambda(X)(b(X) - h(X))^2] &= \frac{1}{2}\mathbb{E}_\theta[S_{bp}(b(X) - h(X))^2] \\ &= \mathbb{E}_\theta[S_{bp}\frac{h(X)^2}{2} - S_{bp}b(X)h(X)] + \frac{1}{2}\mathbb{E}_\theta[S_{bp}b(X)^2] \\ &= \mathbb{E}_\theta[S_{bp}\frac{h(X)^2}{2} + m_2(O, h)] + \frac{1}{2}\mathbb{E}_\theta[S_{bp}b(X)^2] \end{aligned}$$

where in the third line we use equation (2.2) and the extra term  $\frac{1}{2}\mathbb{E}_\theta[S_{bp}b(X)^2]$  is unrelated to minimizing  $\frac{1}{2}\mathbb{E}_\theta[\lambda(X)(b(X) - h(X))^2]$ . By symmetry,

$$\frac{1}{2}\mathbb{E}_\theta[\lambda(X)(p(X) - h(X))^2] = \mathbb{E}_\theta[S_{bp}\frac{h(X)^2}{2} + m_1(O, h)] + \frac{1}{2}\mathbb{E}_\theta[S_{bp}p(X)^2].$$

We now define a new quantity that will be important to develop the procedure that avoids the faithfulness Condition **Faithfulness** or Condition **Faithfulness**( $\delta'$ ).

**Definition 3.2.** The Cauchy-Schwarz (CS) bias of  $\hat{\psi}_1$  as an estimator of  $\psi(\theta)$  is defined as

$$(3.18) \quad \text{CSBias}_\theta(\hat{\psi}_1) := \{\mathbb{E}_\theta[\lambda(X)\{b(X) - \hat{b}(X)\}^2]\mathbb{E}_\theta[\lambda(X)\{p(X) - \hat{p}(X)\}^2]\}^{1/2} \equiv \mathbb{L}_{\theta,2,\hat{b}}\mathbb{L}_{\theta,2,\hat{p}}.$$

When a data analyst reports a nominal  $(1 - \alpha)$  Wald confidence interval based on  $\hat{\psi}_1$ , (s)he naturally appeals to the following Cauchy-Schwarz (CS) null hypothesis  $\text{NH}_{0,\text{CS}}$ , together with its operationalized version:

$$(3.19) \quad \text{NH}_{0,\text{CS}} : \text{CSBias}_\theta(\hat{\psi}_1) = o(n^{-1/2})$$

$$(3.20) \quad \text{H}_{0,\text{CS}}(\delta) : \frac{\text{CSBias}_\theta(\hat{\psi}_1)}{\text{s.e.}_\theta(\hat{\psi}_1)} < \delta,$$

as the *justification* of a validity claim that the Wald confidence interval covers  $\psi(\theta)$  within the tolerance level (set by  $\delta$ ) below the nominal coverage. We have the following logical orderings between the null hypotheses defined above:

**Lemma 3.2.**

- (1)  $\text{NH}_{0,\text{CS}} \Rightarrow \text{NH}_0$ , and similarly  $\text{H}_{0,\text{CS}}(\delta) \Rightarrow \text{H}_0(\delta)$ ;
- (2)  $\text{NH}_{0,\text{CS}} \Rightarrow \text{NH}_{0,k}$  for all  $k$ , and similarly  $\text{H}_{0,\text{CS}}(\delta) \Rightarrow \text{H}_{0,k}(\delta)$  for all  $k$ .

*Proof.* The first part simply follows from CS inequality. The second part follows from the derivation below:

$$\begin{aligned} (3.21) \quad |\text{Bias}_{\theta,k}(\hat{\psi}_1)| &= \left| \mathbb{E}_\theta \left[ \Pi_\theta[\lambda^{1/2}(b - \hat{b})|\bar{z}_k](X) \Pi_\theta[\lambda^{1/2}(p - \hat{p})|\bar{z}_k](X) \right] \right| \\ &\leq \mathbb{L}_{\theta,2,\hat{b},k} \mathbb{L}_{\theta,2,\hat{p},k} \leq \mathbb{L}_{\theta,2,\hat{b}} \mathbb{L}_{\theta,2,\hat{p}} \end{aligned}$$

where the first inequality follows from CS inequality and the second inequality is a consequence of the fact that a projection contracts  $L_2(\mathbb{P}_\theta)$  norms. ■

However the converse statements of Lemma 3.2 are not always true: for example,  $NH_0$  may be true (and thus, by Theorem 2.3 the above the Wald confidence interval centered at  $\hat{\psi}_1$  is valid) even when the CS null hypothesis is false. Suppose we empirically falsify the *justification*  $NH_{0,CS}(H_{0,CS}(\delta))$  for the null hypothesis of actual interest  $NH_0(H_0(\delta))$ . Then, although logically  $NH_0$  may be true, there seems, to us, neither a substantive nor a philosophical reason to assume  $NH_0$  is true in the absence of  $NH_{0,CS}$ . Thus we will take the following philosophical stance:

**Condition CS.** *If the CS null hypothesis  $NH_{0,CS}$  and  $H_{0,CS}(\delta)$  being true is used as the justification for the validity of the Wald confidence interval  $CI_\alpha(\hat{\psi}_1)$ , but in fact are false, one should refuse to support claims whose validity rests on the truth of  $NH_0$  or  $H_0(\delta)$ ; in particular, the claims that the Wald confidence intervals centered at  $\hat{\psi}_1$  have true coverage greater than or equal to their nominal.*

Clearly Condition CS will allow meaningful inferences regarding  $\psi(\theta)$  only if it is possible to empirically reject the CS null hypothesis  $NH_{0,CS}$ . Indeed, it follows from Lemma 3.2(2) that the rejection of the surrogate  $H_{0,k}(\delta)$  implies rejection of  $H_{0,CS}(\delta)$ . Therefore we are only left to show that  $\hat{\chi}_{2,k}(\Sigma_k^{-1}; \varsigma_k, \delta)$  is a “good” test of the surrogate null hypothesis  $H_{0,k}(\delta)$ .

**Comment 3.6.** *Note that  $CSBias_\theta(\hat{\psi}_1)$  is not the only upper bound of  $|Bias_\theta(\hat{\psi}_1)|$  by CS inequality. For example, the following upper bound also holds:*

$$\begin{aligned} Bias_\theta(\hat{\psi}_1) &= \left| E_\theta[\lambda(X)(b(X) - \hat{b}(X))(p(X) - \hat{p}(X))] \right| \\ (3.22) \quad &\leq \left\{ E_\theta[\lambda(X)^2(b(X) - \hat{b}(X))^2] E_\theta[(p(X) - \hat{p}(X))^2] \right\}^{1/2}. \end{aligned}$$

Under the conditions in Theorem 2.2 and Condition W, however, we show that there exists two positive constants  $c < C$  such that

$$cCSBias_\theta(\hat{\psi}_1) \leq RHS \text{ of eq. (3.22)} \leq CCSBias_\theta(\hat{\psi}_1).$$

To see this,

$$\begin{aligned} RHS \text{ of eq. (3.22)} &= \left\{ E_\theta[\lambda(X)\lambda(X)(b(X) - \hat{b}(X))^2] E_\theta[\lambda(X)^{-1}\lambda(X)(p(X) - \hat{p}(X))^2] \right\}^{1/2} \\ &\leq \sqrt{\frac{\text{ess sup}_x \lambda(x)}{\text{ess inf}_x \lambda(x)}} \left\{ E_\theta[\lambda(X)(b(X) - \hat{b}(X))^2] E_\theta[\lambda(X)(p(X) - \hat{p}(X))^2] \right\}^{1/2} \end{aligned}$$

and

$$\begin{aligned} RHS \text{ of eq. (3.22)} &= \left\{ E_\theta[\lambda(X)\lambda(X)(b(X) - \hat{b}(X))^2] E_\theta[\lambda(X)^{-1}\lambda(X)(p(X) - \hat{p}(X))^2] \right\}^{1/2} \\ &\geq \sqrt{\frac{\text{ess inf}_x \lambda(x)}{\text{ess sup}_x \lambda(x)}} \left\{ E_\theta[\lambda(X)(b(X) - \hat{b}(X))^2] E_\theta[\lambda(X)(p(X) - \hat{p}(X))^2] \right\}^{1/2} \end{aligned}$$

where the above two inequalities both follow from Condition **W**(3) on  $\lambda(\cdot)$ . Hence if we reject  $\text{NH}_{0,\text{CS}}$ , there is no reason to believe that RHS of equation (3.22) =  $o(n^{-1/2})$ .

**3.3. Statistical properties of  $\widehat{\chi}_{2,k}(\Sigma_k^{-1}; \varsigma_k, \delta)$ .** As mentioned above, based on the statistical properties of  $\widehat{\psi}_1$  and  $\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})$  summarized in Theorems 2.3 and 3.1, for a DR functional  $\psi(\theta)$ , we study the statistical property of the oracle two-sided test  $\widehat{\chi}_{2,k}(\Sigma_k^{-1}; \varsigma_k, \delta)$  as a test of the surrogate null hypothesis  $\text{H}_{0,k}(\delta)$ .

The following theorem characterizes the asymptotic level and power of  $\widehat{\chi}_{2,k}(\Sigma_k^{-1}; \varsigma_k, \delta)$  for  $\text{H}_{0,k}(\delta)$  (3.14).

**Theorem 3.2.** For a DR functional  $\psi(\theta)$ , under the conditions in Theorem 2.2, Theorem 3.1 and Condition **W**, when  $k \rightarrow \infty$  but  $k = o(n)$ , for any given  $\varsigma_k, \delta > 0$ , suppose that  $\frac{|\text{Bias}_{\theta,k}(\widehat{\psi}_1)|}{\text{s.e.}_{\theta}[\widehat{\psi}_1]} = \gamma$  for some (sequence)  $\gamma = \gamma(n)$  (where  $\gamma(n)$  can diverge with  $n$ ), then the rejection probability of  $\widehat{\chi}_{2,k}(\Sigma_k^{-1}; \varsigma_k, \delta)$  converges to

$$(3.23) \quad 2 - \Phi \left( \varsigma_k - \lim_{n \rightarrow \infty} (\gamma - \delta) \frac{\text{s.e.}_{\theta}[\widehat{\psi}_1]}{\text{s.e.}_{\theta}[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} \right) - \Phi \left( \varsigma_k + \lim_{n \rightarrow \infty} (\gamma + \delta) \frac{\text{s.e.}_{\theta}[\widehat{\psi}_1]}{\text{s.e.}_{\theta}[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} \right)$$

as  $n \rightarrow \infty$ . In particular,

- (1) under  $\text{H}_{0,k}(\delta) : \gamma \leq \delta$ ,  $\widehat{\chi}_{2,k}(\Sigma_k^{-1}; \varsigma_k, \delta)$  rejects the null with probability less than or equal to  $2(1 - \Phi(\varsigma_k))$ , as  $n \rightarrow \infty$ ;
- (2) under the following alternative to  $\text{H}_{0,k}(\delta)$ :  $\gamma = \delta + c$ , for any diverging sequence  $c = c(n) \rightarrow \infty$ ,  $\widehat{\chi}_{2,k}(\Sigma_k^{-1}; \varsigma_k, \delta)$  rejects the null with probability converging to 1, as  $n \rightarrow \infty$ .
- (2') If  $\mathbb{L}_{\theta,2,\widehat{b},k}$  and  $\mathbb{L}_{\theta,2,\widehat{p},k}$  converge to 0, under the following alternative to  $\text{H}_{0,k}(\delta)$ :  $\gamma = \delta + c$ , for any fixed  $c > 0$  or any diverging sequence  $c = c(n) \rightarrow \infty$ ,  $\widehat{\chi}_{2,k}(\Sigma_k^{-1}; \varsigma_k, \delta)$  has rejection probability converging to 1, as  $n \rightarrow \infty$ .

**Comment 3.7.** In Appendix A.4, we prove equation (3.23). We now prove that equation (3.23) implies Theorem 3.2(1)-(2) and (2').

- Regarding (1), under  $\text{H}_{0,k}(\delta) : \gamma \leq \delta$ ,

$$-(\gamma - \delta) \frac{\text{s.e.}_{\theta}[\widehat{\psi}_1]}{\text{s.e.}_{\theta}[\widehat{\mathbb{I}\mathbb{F}}_{22,k}]} \geq 0 \text{ and } (\gamma + \delta) \frac{\text{s.e.}_{\theta}[\widehat{\psi}_1]}{\text{s.e.}_{\theta}[\widehat{\mathbb{I}\mathbb{F}}_{22,k}]} \geq 0,$$

which implies that the rejection probability is less than or equal to  $2 - 2\Phi(\varsigma_k)$ . Choose  $\varsigma_k = z_{\alpha^\dagger/2}$ ,  $2(1 - \Phi(\varsigma_k)) = 2\alpha^\dagger/2 = \alpha^\dagger$  and conclude that the test is a valid level  $\alpha^\dagger$  test of the null.

- Under the alternative to  $\text{H}_{0,k}(\delta)$  with  $\gamma = \delta + c$  for some  $c > 0$ , it follows from Theorem 3.1 and equation (3.23) that the rejection probability of  $\widehat{\chi}_{2,k}(\Sigma_k^{-1}; \varsigma_k, \delta)$ , as  $n \rightarrow \infty$ , is no smaller than

$$2 - \Phi \left( \varsigma_k - c\Theta(b, p, \widehat{b}, \widehat{p}, f_X, \bar{z}_k) \left\{ \frac{k}{n} + \mathbb{L}_{\theta,2,\widehat{b},k} + \mathbb{L}_{\theta,2,\widehat{p},k} \right\}^{-1} \right) - \Phi(\infty)$$

where  $\Theta(b, p, \hat{b}, \hat{p}, f_X, \bar{Z}_k)$  is some positive constant depending on the true regression functions  $b$  and  $p$ , the estimated functions  $\hat{b}, \hat{p}$  from the training sample, the density  $f_X$  of  $X$  and the chosen basis functions  $\bar{z}_k$ . To have power approaching 1 to reject  $H_{0,k}(\delta)$ , we need one of the following:

- If one of  $\mathbb{L}_{\theta,2,\hat{b},k}$  and  $\mathbb{L}_{\theta,2,\hat{p},k}$  is  $O(1)$ , we need  $c \rightarrow \infty$  to guarantee that the rejection probability of  $\hat{\chi}_{2,k}(\Sigma_k^{-1}; \varsigma_k, \delta)$  converges to  $1 - \Phi(-\infty) = 1$ . Hence we have Theorem 3.2(2).
- If  $c$  is fixed, we need both  $\mathbb{L}_{\theta,2,\hat{b},k}$  and  $\mathbb{L}_{\theta,2,\hat{p},k}$  to be  $o(1)$  to guarantee that the rejection probability of  $\hat{\chi}_{2,k}(\Sigma_k^{-1}; \varsigma_k, \delta)$  converges to  $1 - \Phi(-\infty) = 1$ . Hence we have Theorem 3.2(2').

Theorem 3.2 implies that  $\hat{\chi}_{2,k}(\Sigma_k^{-1}; z_{\alpha^\dagger/2}, \delta)$  is an asymptotically valid level  $\alpha^\dagger$  two-sided test of the surrogate null  $H_{0,k}(\delta)$ , and hence by Lemma 3.2(2) it is also an asymptotically  $\alpha^\dagger$  level test of  $H_{0,CS}(\delta)$ . Thus when  $\hat{\chi}_{2,k}(\Sigma_k^{-1}; z_{\alpha^\dagger/2}, \delta)$  rejects  $H_{0,k}(\delta)$ , it also rejects  $H_{0,CS}(\delta)$  and by Condition CS, we conclude that we have no justification for the coverage validity of  $CI_\alpha(\hat{\psi}_1)$ . It should be noted that  $\hat{\chi}_{2,k}(\Sigma_k^{-1}; z_{\alpha^\dagger/2}, \delta)$  can be a powerless test for  $H_{0,CS}(\delta)$  under certain laws  $P_\theta$ .

If one is willing to assume the faithfulness Condition Faithfulness( $\delta'$ ), then  $\hat{\chi}_{2,k}(\Sigma_k^{-1}; z_{\alpha^\dagger/2}, \delta)$  is an asymptotically valid level  $\alpha^\dagger$  test of  $H_0(\delta\delta')$ . As above, it can be a powerless test of  $H_0(\delta\delta')$ .

#### 4. ASSUMPTION FREE TEST WHEN $\Sigma_k^{-1}$ IS UNKNOWN

Heretofore we have assumed that  $\Sigma_k^{-1}$  is known, or almost equivalently, that we have access to a semisupervised data of size much larger than  $N$  but containing information only on the components  $(S_{bp}, X)$  of the observed data vector  $O = (W, X)$ . In the example of the counterfactual mean  $E_\theta[Y(a=1)]$  under strong ignorability, the semisupervised data needs to include  $(A, X)$ , which might be difficult to collect in practice. In this section, we will focus on the impact of estimating  $\Sigma_k^{-1}$  on the test  $\hat{\chi}_{2,k}(\Sigma_k^{-1}; \varsigma_k, \delta)$  for testing the null hypothesis  $H_{0,k}(\delta)$ .

One approach to proceed is to construct an estimator of the density of  $X$  (Liu et al., 2016; Robins et al., 2008, 2009a, 2017). But when dimension of the covariates  $X$  is large, accurate density estimation is problematic. More recently, in the regime  $k = o(n)$ , Mukherjee et al. (2017) proposed to replace  $\Sigma_k^{-1}$  by  $\hat{\Sigma}_k^{-1}$  where  $\hat{\Sigma}_k = n^{-1} \sum_{i \in \text{tr}} S_{bp,i} \bar{z}_k(X_i) \bar{z}_k(X_i)^\top$  is the sample Gram matrix estimator from the training sample. They show that, unlike the oracle  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})$ ,  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\hat{\Sigma}_k^{-1})$  is a biased estimator of  $\text{Bias}_{\theta,k}(\hat{\psi}_1)$  with bias  $\text{EB}_{\theta,2,k}(\hat{\Sigma}_k^{-1}) \equiv E_\theta[\hat{\mathbb{I}}\mathbb{F}_{22,k}(\hat{\Sigma}_k^{-1}) - \hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})]$ <sup>4</sup> of order

$$O\left(\mathbb{L}_{\theta,2,\hat{b},k} \mathbb{L}_{\theta,2,\hat{p},k} \mathbb{L}_{\theta,2,\hat{\Sigma},k}\right) \lesssim O\left(\mathbb{L}_{\theta,2,\hat{b},k} \mathbb{L}_{\theta,2,\hat{p},k} \sqrt{\frac{k \log(k)}{n}}\right)$$

under Condition W. [Note that  $\text{EB}_{\theta,2,k}(\hat{\Sigma}_k^{-1})$  is also the bias of  $\hat{\psi}_{2,k}(\hat{\Sigma}_k^{-1})$  as an estimator of the  $\tilde{\psi}_k(\theta)$ .] It follows that  $\text{EB}_{\theta,2,k}(\hat{\Sigma}_k^{-1}) \rightarrow 0$  as  $n \rightarrow \infty$  if (i)  $k \log(k) = o(n)$  and (ii)  $\mathbb{L}_{\theta,2,\hat{b},k}$  and  $\mathbb{L}_{\theta,2,\hat{p},k}$  are bounded. However, as we will show

<sup>4</sup>We denote the bias due to estimating  $\Sigma_k^{-1}$  by  $\text{EB}_{\theta,2,k}(\hat{\Sigma}_k^{-1})$  because this bias is termed as Estimation Bias in Robins et al. (2008) and Mukherjee et al. (2017)

in this section, we need to add an extra bias correction term to ensure that the test  $\hat{\chi}_{2,k}(\Sigma_k^{-1}; \varsigma_k, \delta)$  is still valid under  $H_{0,k}(\delta)$  (3.14) when  $\Sigma_k^{-1}$  is estimated by  $\hat{\Sigma}_k^{-1}$ . We first define

$$(4.1) \quad \hat{\mathbb{I}\mathbb{F}}_{33,k}(\Sigma_k^{-1}) := \frac{1}{n(n-1)(n-2)} \sum_{1 \leq i_1 \neq i_2 \neq i_3 \leq n} \hat{\mathbb{I}\mathbb{F}}_{33,k,\bar{i}_3}(\Sigma_k^{-1})$$

where

$$(4.2) \quad \hat{\mathbb{I}\mathbb{F}}_{33,k,\bar{i}_3}(\Sigma_k^{-1}) := \left[ \mathcal{E}_{\hat{b},m_2}(\bar{z}_k)(O) \right]_{i_1}^\top \Sigma_k^{-1} \left\{ \left[ S_{b,p} \bar{z}_k(X) \bar{z}_k(X)^\top \right]_{i_3} - \Sigma_k \right\} \Sigma_k^{-1} \left[ \mathcal{E}_{\hat{p},m_1}(\bar{z}_k)(O) \right]_{i_2}.$$

Then  $\hat{\mathbb{I}\mathbb{F}}_{33,k}(\hat{\Sigma}_k^{-1})$  simply replaces  $\Sigma_k^{-1}$  by  $\hat{\Sigma}_k^{-1}$  in  $\hat{\mathbb{I}\mathbb{F}}_{33,k}(\Sigma_k^{-1})$ . Analogous to Proposition 3.1, we relate  $\hat{\psi}_{3,k}(\hat{\Sigma}_k^{-1}) - \psi(\theta) \equiv \hat{\psi}_1 - \hat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1}) - \psi(\theta)$  to the *third order* influence function of the truncated parameter  $\tilde{\psi}_k(\theta)$  of a DR functional, following Robins et al. (2008, Definition 2.1). In particular, we prove the following in Appendix A.1.

**Proposition 4.1.**

- $\hat{\psi}_{3,k}(\hat{\Sigma}_k^{-1}) - \psi(\theta)$  is the third order influence function of  $\tilde{\psi}_k(\theta)$  under the law  $P_{\hat{\theta}}$ ;
- $\hat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})$  is the third order influence function of  $\text{Bias}_{\theta,k}(\hat{\psi}_1)$  under the law  $P_{\hat{\theta}}$ .

Then we have the following results on the variances and the estimation biases of  $\hat{\mathbb{I}\mathbb{F}}_{22,k}(\hat{\Sigma}_k^{-1})$  and  $\hat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})$  as an estimator of  $\text{Bias}_{\theta,k}(\hat{\psi}_1) \equiv E_\theta[\hat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]$ : the estimation bias bounds are proved in Mukherjee et al. (2017) and the variance bounds are proved in Appendix A.2:

**Proposition 4.2.** Under the conditions of Theorem 3.2, when  $k = o(n)$ , the followings hold on the event that  $\hat{\Sigma}_k$  is invertible:

- (1) The estimation bias and variance of  $\hat{\mathbb{I}\mathbb{F}}_{22,k}(\hat{\Sigma}_k^{-1})$  are of the following order:

$$\begin{aligned} \left| \text{EB}_{\theta,2,k}(\hat{\Sigma}_k^{-1}) \right| &\lesssim \mathbb{L}_{\theta,2,\hat{b},k} \mathbb{L}_{\theta,2,\hat{p},k} \mathbb{L}_{\theta,2,\hat{\Sigma},k} \lesssim \mathbb{L}_{\theta,2,\hat{b},k} \mathbb{L}_{\theta,2,\hat{p},k} \sqrt{\frac{k \log(k)}{n}}, \\ \text{var}_\theta \left[ \hat{\mathbb{I}\mathbb{F}}_{22,k}(\hat{\Sigma}_k^{-1}) \right] &\lesssim \frac{1}{n} \left\{ \frac{k}{n} + \mathbb{L}_{\theta,2,\hat{b},k}^2 + \mathbb{L}_{\theta,2,\hat{p},k}^2 \right\}. \end{aligned}$$

- (2) The estimation bias and variance of  $\hat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1}) := \hat{\mathbb{I}\mathbb{F}}_{22,k}(\hat{\Sigma}_k^{-1}) + \hat{\mathbb{I}\mathbb{F}}_{33,k}(\hat{\Sigma}_k^{-1})$  are of the following order:

$$\begin{aligned} \left| \text{EB}_{\theta,3,k}(\hat{\Sigma}_k^{-1}) \right| &\lesssim \mathbb{L}_{\theta,2,\hat{b},k} \mathbb{L}_{\theta,2,\hat{p},k} \mathbb{L}_{\theta,2,\hat{\Sigma},k}^2 \lesssim \mathbb{L}_{\theta,2,\hat{b},k} \mathbb{L}_{\theta,2,\hat{p},k} \frac{k \log(k)}{n}, \\ \text{var}_\theta \left[ \hat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1}) \right] &\lesssim \frac{1}{n} \left\{ \frac{k}{n} + \mathbb{L}_{\theta,2,\hat{b},k}^2 + \mathbb{L}_{\theta,2,\hat{p},k}^2 \right\} \end{aligned}$$

where

$$\text{EB}_{\theta,3,k}(\hat{\Sigma}_k^{-1}) := E_\theta[\hat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})] - \text{Bias}_{\theta,k}(\hat{\psi}_1).$$

(3) As  $k \rightarrow \infty, n \rightarrow \infty$ , but  $k = o(n^2)$ ,

$$\frac{\widehat{\mathbb{IF}}_{22,k}(\widehat{\Sigma}_k^{-1}) - \mathbb{E}_\theta \left[ \widehat{\mathbb{IF}}_{22,k}(\widehat{\Sigma}_k^{-1}) \right]}{\text{s.e.}_\theta[\widehat{\mathbb{IF}}_{22,k}(\widehat{\Sigma}_k^{-1})]}, \frac{\widehat{\mathbb{IF}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1}) - \mathbb{E}_\theta \left[ \widehat{\mathbb{IF}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1}) \right]}{\text{s.e.}_\theta[\widehat{\mathbb{IF}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]} \xrightarrow{d} N(0, 1).$$

We now define the following candidate test statistics when  $\Sigma_k^{-1}$  is replaced by  $\widehat{\Sigma}_k^{-1}$ :

$$(4.3) \quad \widehat{\chi}_{2,k}(\widehat{\Sigma}_k^{-1}; \varsigma_k, \delta) := \mathbb{1} \left\{ \frac{|\widehat{\mathbb{IF}}_{22,k}(\widehat{\Sigma}_k^{-1})|}{\widehat{\text{s.e.}}[\widehat{\psi}_1]} - \varsigma_k \frac{\widehat{\text{s.e.}}[\widehat{\mathbb{IF}}_{22,k}(\widehat{\Sigma}_k^{-1})]}{\widehat{\text{s.e.}}[\widehat{\psi}_1]} > \delta \right\},$$

$$(4.4) \quad \widehat{\chi}_{3,k}(\widehat{\Sigma}_k^{-1}; \varsigma_k, \delta) := \mathbb{1} \left\{ \frac{|\widehat{\mathbb{IF}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})|}{\widehat{\text{s.e.}}[\widehat{\psi}_1]} - \varsigma_k \frac{\widehat{\text{s.e.}}[\widehat{\mathbb{IF}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]}{\widehat{\text{s.e.}}[\widehat{\psi}_1]} > \delta \right\},$$

As a consequence of the asymptotic statistical properties of  $\widehat{\mathbb{IF}}_{22,k}(\widehat{\Sigma}_k^{-1})$  and  $\widehat{\mathbb{IF}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})$  in Proposition 4.2, theoretically one can guarantee that  $\widehat{\chi}_{3,k}(\widehat{\Sigma}_k^{-1}; \varsigma_k, \delta)$ , but not  $\widehat{\chi}_{2,k}(\widehat{\Sigma}_k^{-1}; \varsigma_k, \delta)$ , is asymptotically equivalent to the oracle test  $\widehat{\chi}_{2,k}(\Sigma_k^{-1}; \varsigma_k, \delta)$ . The proof can be found in Appendix A.5.

**Proposition 4.3.** Under the conditions in Proposition 4.2,  $k(\log(k))^2 = o(n)$  and together with the additional restriction

$$(4.5) \quad \text{Bias}_{\theta,k}(\widehat{\psi}_1) \neq o \left( \mathbb{L}_{\theta,2,\widehat{b},k} \mathbb{L}_{\theta,2,\widehat{p},k} \right),$$

for any given  $\varsigma_k, \delta > 0$ , suppose that  $\frac{|\text{Bias}_{\theta,k}(\widehat{\psi}_1)|}{\text{s.e.}_\theta[\widehat{\psi}_1]} = \gamma$  for some (sequence)  $\gamma = \gamma(n)$  (where  $\gamma(n)$  can diverge with  $n$ ), then the rejection probability of  $\widehat{\chi}_{3,k}(\widehat{\Sigma}_k^{-1}; \varsigma_k, \delta)$  converges to

$$(4.6) \quad 2 - \Phi \left( \varsigma_k - \lim_{n \rightarrow \infty} (\gamma - \delta) \frac{\text{s.e.}_\theta[\widehat{\psi}_1]}{\text{s.e.}_\theta[\widehat{\mathbb{IF}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]} \right) - \Phi \left( \varsigma_k + \lim_{n \rightarrow \infty} (\gamma + \delta) \frac{\text{s.e.}_\theta[\widehat{\psi}_1]}{\text{s.e.}_\theta[\widehat{\mathbb{IF}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]} \right)$$

as  $n \rightarrow \infty$ . In particular,

- (1) under  $\mathbf{H}_{0,k}(\delta) : \gamma \leq \delta$ ,  $\widehat{\chi}_{3,k}(\widehat{\Sigma}_k^{-1}; \varsigma_k, \delta)$  rejects the null with probability less than or equal to  $2(1 - \Phi(\varsigma_k))$ , as  $n \rightarrow \infty$ ;
- (2) under the following alternative to  $\mathbf{H}_{0,k}(\delta)$ :  $\gamma = \delta + c$ , for any diverging sequence  $c = c(n) \rightarrow \infty$ ,  $\widehat{\chi}_{3,k}(\widehat{\Sigma}_k^{-1}; \varsigma_k, \delta)$  rejects the null with probability converging to 1, as  $n \rightarrow \infty$ .
- (2') If  $\mathbb{L}_{\theta,2,\widehat{b},k}$  and  $\mathbb{L}_{\theta,2,\widehat{p},k}$  converge to 0, under the following alternative to  $\mathbf{H}_{0,k}(\delta)$ :  $\gamma = \delta + c$ , for any fixed  $c > 0$  or any diverging sequence  $c = c(n) \rightarrow \infty$ ,  $\widehat{\chi}_{3,k}(\widehat{\Sigma}_k^{-1}; \varsigma_k, \delta)$  has rejection probability converging to 1, as  $n \rightarrow \infty$ .

**Comment 4.1** (Comments on Proposition 4.3).

- (1) As a consequence, under the conditions of Proposition 4.3,  $\widehat{\chi}_{3,k}(\widehat{\Sigma}_k^{-1}; z_{\alpha^\dagger/2}, \delta)$  is an asymptotically valid level  $\alpha^\dagger$  two-sided test of the surrogate null  $\mathbf{H}_{0,k}(\delta)$ , and hence by Lemma 3.2(2) it is also an asymptotically  $\alpha^\dagger$  level test of  $\mathbf{H}_{0,\text{CS}}(\delta)$ . Thus when  $\widehat{\chi}_{3,k}(\widehat{\Sigma}_k^{-1}; z_{\alpha^\dagger/2}, \delta)$  rejects  $\mathbf{H}_{0,k}(\delta)$ , it also rejects  $\mathbf{H}_{0,\text{CS}}(\delta)$  and by Condition CS, we conclude that we have



no justification for the coverage validity of  $\text{CI}_\alpha(\hat{\psi}_1)$ . It should be noted that  $\hat{\chi}_{3,k}(\hat{\Sigma}_k^{-1}; z_{\alpha^\dagger/2}, \delta)$  can be a powerless test for  $\text{H}_{0,\text{CS}}(\delta)$  under certain laws  $\text{P}_\theta$ .

If one is willing to assume the faithfulness Condition **Faithfulness**( $\delta'$ ), then  $\hat{\chi}_{3,k}(\hat{\Sigma}_k^{-1}; z_{\alpha^\dagger/2}, \delta)$  is an asymptotically valid level  $\alpha^\dagger$  test of  $\text{H}_0(\delta\delta')$ . As above, it can be a powerless test of  $\text{H}_0(\delta\delta')$ .

- (2) The additional restriction  $\text{Bias}_{\theta,k}(\hat{\psi}_1) \neq o(\mathbb{L}_{\theta,2,\hat{b},k} \mathbb{L}_{\theta,2,\hat{p},k})$  will often hold when we estimate both nuisance functions  $b$  and  $p$  from a single training sample. However when we further split the training sample to estimate the nuisance functions from separate independent subsamples, it is possible to construct estimators that violate this restriction. For example, see [Kennedy \(2020\)](#); [Newey and Robins \(2018\)](#). In Section 5, we will propose a procedure that successively relaxes this additional restriction but at the expense of higher computational cost.
- (3) We cannot show that  $\hat{\chi}_{2,k}(\hat{\Sigma}_k^{-1}; \varsigma_k, \delta)$  is asymptotically equivalent to the oracle test  $\hat{\chi}_{2,k}(\Sigma_k^{-1}; \varsigma_k, \delta)$ . It would be true if the bias  $\text{EB}_{\theta,2,k}(\hat{\Sigma}_k^{-1})$  is dominated by  $\text{s.e.}_\theta[\hat{\mathbb{I}\mathbb{F}}_{22,k}(\hat{\Sigma}_k^{-1})] \asymp \frac{1}{\sqrt{n}} \left( \sqrt{\frac{k}{n}} + \mathbb{L}_{\theta,2,\hat{b},k} + \mathbb{L}_{\theta,2,\hat{p},k} \right)$ . Consider the following scenario: Under  $\text{H}_{0,k}(\delta)$  (3.14), the further restriction  $\text{Bias}_{\theta,k}(\hat{\psi}_1) \neq o(\mathbb{L}_{\theta,2,\hat{b},k} \mathbb{L}_{\theta,2,\hat{p},k})$  and  $\mathbb{L}_{\theta,2,\hat{b},k} \asymp \mathbb{L}_{\theta,2,\hat{p},k} \leq n^{-1/4}$ , however, we are only able to obtain the following upper bound:

$$\begin{aligned} \text{EB}_{\theta,2,k}(\hat{\Sigma}_k^{-1}) &\lesssim \mathbb{L}_{\theta,2,\hat{b},k} \mathbb{L}_{\theta,2,\hat{p},k} \sqrt{\frac{k \log(k)}{n}} \\ &\ll \sqrt{\frac{1}{n}} \sqrt{\frac{k \log(k)}{n}} = \frac{\sqrt{k \log(k)}}{n} \end{aligned}$$

with an extra  $\sqrt{\log(k)}$  factor compared to the desired rate  $\sqrt{k}/n$ . It is possible that our upper bound on  $\text{EB}_{\theta,2,k}(\hat{\Sigma}_k^{-1})$  is not tight, but this remains an open theoretical problem.

## 5. HIGHER ORDER TESTS TO RELAX THE ADDITIONAL RESTRICTION (4.5)

In this section, we propose a procedure that successively relax the additional restriction  $\text{Bias}_{\theta,k}(\hat{\psi}_1) \neq o(\mathbb{L}_{\theta,2,\hat{b},k} \mathbb{L}_{\theta,2,\hat{p},k})$  imposed in Proposition 4.3 but at the expense of greater computational cost. This procedure depends on the following quantities that generalize  $\hat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\Sigma_k^{-1})$  to higher order U-statistics: for  $m \geq 2$ ,

$$\hat{\mathbb{I}\mathbb{F}}_{22 \rightarrow mm,k}(\Sigma_k^{-1}) := \sum_{j=2}^m \hat{\mathbb{I}\mathbb{F}}_{jj,k}(\Sigma_k^{-1})$$

and

$$(5.1) \quad \hat{\mathbb{I}\mathbb{F}}_{jj,k}(\Sigma_k^{-1}) := \frac{(n-j)!}{n!} \sum_{1 \leq i_1 \neq \dots \neq i_j \leq n} \hat{\mathbb{I}\mathbb{F}}_{jj,k,\bar{i}_j}(\Sigma_k^{-1})$$

where

$$(5.2) \quad \hat{\mathbb{I}\mathbb{F}}_{jj,k,\bar{i}_j}(\Sigma_k^{-1}) := (-1)^j \left[ \mathcal{E}_{b,m_2}(\bar{z}_k)(O) \right]_{i_1}^\top \left\{ \prod_{s=3}^j \Sigma_k^{-1} \left( [S_{bp} \bar{z}_k(X) \bar{z}_k(X)^\top]_{i_s} - \Sigma_k \right) \right\} \Sigma_k^{-1} [\mathcal{E}_{\hat{p},m_1}(\bar{z}_k)(O)]_{i_2}.$$

$\widehat{\mathbb{IF}}_{22 \rightarrow mm, k}(\widehat{\Sigma}_k^{-1})$  and  $\widehat{\psi}_{m, k}(\widehat{\Sigma}_k^{-1})$  are defined in the same way, except that  $\Sigma_k^{-1}$  is replaced by  $\widehat{\Sigma}_k^{-1}$ . Similar to Propositions 3.1 and 4.1, in Appendix A.1, we prove the following

**Proposition 5.1.**

- $\widehat{\psi}_{m, k}(\Sigma_k^{-1}) - \psi(\theta)$  is the  $m$ -th order influence function of  $\widetilde{\psi}_k(\theta)$  under the law  $\mathbf{P}_{\widehat{\theta}}$ , where  $\widehat{\psi}_{m, k}(\Sigma_k^{-1}) := \widehat{\psi}_1 - \widehat{\mathbb{IF}}_{22 \rightarrow mm, k}(\Sigma_k^{-1})$ ;
- $\widehat{\mathbb{IF}}_{22 \rightarrow mm, k}(\Sigma_k^{-1})$  is the  $m$ -th order influence function of  $\text{Bias}_{\theta, k}(\widehat{\psi}_1)$  under the law  $\mathbf{P}_{\widehat{\theta}}$ .

**Comment 5.1.** Note that  $\widehat{\mathbb{IF}}_{22 \rightarrow mm, k}(\Sigma_k^{-1})$  is an  $m$ -th order U-statistics. The greater  $m$  is, the higher computational cost is incurred for computing  $\widehat{\mathbb{IF}}_{\ell \ell \rightarrow mm, k}(\Sigma_k^{-1})$ .

Similar to Proposition 4.2, we also obtain the following bounds on the estimation bias and variance of  $\widehat{\mathbb{IF}}_{22 \rightarrow mm, k}(\Sigma_k^{-1})$  as an estimator of  $\text{Bias}_{\theta, k}(\widehat{\psi}_1) \equiv \mathbf{E}_{\theta}[\widehat{\mathbb{IF}}_{22, k}(\Sigma_k^{-1})]$ :

**Proposition 5.2.** Under the conditions of Theorem 3.2, when  $k = o(n)$ , the followings hold on the event that  $\widehat{\Sigma}_k$  is invertible,

$$\begin{aligned} \left| \mathbf{EB}_{\theta, mm, k}(\widehat{\Sigma}_k^{-1}) \right| &\lesssim \mathbb{L}_{\theta, 2, \widehat{b}, k} \mathbb{L}_{\theta, 2, \widehat{p}, k} \mathbb{L}_{\theta, 2, \widehat{\Sigma}, k}^{m-1} \lesssim \mathbb{L}_{\theta, 2, \widehat{b}, k} \mathbb{L}_{\theta, 2, \widehat{p}, k} \left( \frac{k \log(k)}{n} \right)^{\frac{m-1}{2}}, \\ \text{var}_{\theta} \left[ \widehat{\mathbb{IF}}_{22 \rightarrow mm, k}(\widehat{\Sigma}_k^{-1}) \right] &\lesssim \frac{1}{n} \left\{ \frac{k}{n} + \mathbb{L}_{\theta, 2, \widehat{b}, k}^2 + \mathbb{L}_{\theta, 2, \widehat{p}, k}^2 \right\} \end{aligned}$$

where

$$\mathbf{EB}_{\theta, m, k}(\widehat{\Sigma}_k^{-1}) := \mathbf{E}_{\theta}[\widehat{\mathbb{IF}}_{22 \rightarrow mm, k}(\widehat{\Sigma}_k^{-1})] - \text{Bias}_{\theta, k}(\widehat{\psi}_1).$$

In addition, for any  $m \geq 3$ ,

$$\frac{\widehat{\mathbb{IF}}_{22 \rightarrow mm, k}(\widehat{\Sigma}_k^{-1}) - \mathbf{E}_{\theta}[\widehat{\mathbb{IF}}_{22 \rightarrow mm, k}(\widehat{\Sigma}_k^{-1})]}{\text{s.e.}_{\theta}[\widehat{\mathbb{IF}}_{22 \rightarrow mm, k}(\widehat{\Sigma}_k^{-1})]} \xrightarrow{d} N(0, 1).$$

We first define the following test that has the same property as  $\widehat{\chi}_{3, k}(\widehat{\Sigma}_k^{-1}; \varsigma_k, \delta)$  under a weakened version of the additional restriction  $\text{Bias}_{\theta, k}(\widehat{\psi}_1) \neq o(\mathbb{L}_{\theta, 2, \widehat{b}, k} \mathbb{L}_{\theta, 2, \widehat{p}, k})$  imposed in Proposition 4.3: for any fixed  $m \geq 3$ , define

$$(5.3) \quad \widehat{\chi}_{m, k}(\widehat{\Sigma}_k^{-1}; \varsigma_k, \delta) := \mathbb{1} \left\{ \frac{|\widehat{\mathbb{IF}}_{22 \rightarrow mm, k}(\widehat{\Sigma}_k^{-1})|}{\widehat{\text{s.e.}}[\widehat{\psi}_1]} - \varsigma_k \frac{\widehat{\text{s.e.}}[\widehat{\mathbb{IF}}_{22 \rightarrow mm, k}(\widehat{\Sigma}_k^{-1})]}{\widehat{\text{s.e.}}[\widehat{\psi}_1]} > \delta \right\}.$$

Proposition 5.2 immediately implies the following:

**Proposition 5.3.** Every statements in Proposition 4.3 about the test  $\widehat{\chi}_{3, k}(\widehat{\Sigma}_k^{-1}; \varsigma_k, \delta)$  hold for  $\widehat{\chi}_{m, k}(\widehat{\Sigma}_k^{-1}; \varsigma_k, \delta)$ , under the same conditions in Proposition 4.3, except that the additional restriction  $\text{Bias}_{\theta, k}(\widehat{\psi}_1) \neq o(\mathbb{L}_{\theta, 2, \widehat{b}, k} \mathbb{L}_{\theta, 2, \widehat{p}, k})$  is relaxed to  $\text{Bias}_{\theta, k}(\widehat{\psi}_1) \neq o\left(\mathbb{L}_{\theta, 2, \widehat{b}, k} \mathbb{L}_{\theta, 2, \widehat{p}, k} \left(\frac{k \log(k)}{n}\right)^{\frac{m-3}{2}}\right)$ .

**Comment 5.2.** Note that the greater  $m$  is, the less restrictive  $\text{Bias}_{\theta, k}(\widehat{\psi}_1) \neq o\left(\mathbb{L}_{\theta, 2, \widehat{b}, k} \mathbb{L}_{\theta, 2, \widehat{p}, k} \left(\frac{k \log(k)}{n}\right)^{\frac{m-3}{2}}\right)$  becomes.

As remarked in Comment 5.1, although increasing  $m$  relaxes the restriction (4.5), it incurs higher computational cost. In practice, one might want  $m$  to be as small as possible. In view of such computational concern in practice, we propose the following “early-stopping” test that stops increasing  $m$  if  $\hat{\chi}_{m,k}(\hat{\Sigma}_k^{-1}; \varsigma_k, \delta)$  fails to reject  $H_{0,k}(\delta)$  at a smaller  $m$ : for a fixed integer  $M > 0$ , define

$$(5.4) \quad \hat{\chi}_{M,k}^{es}(\hat{\Sigma}_k^{-1}; \varsigma_k, \delta) := \mathbb{1} \left\{ \frac{|\hat{\mathbb{I}}\mathbb{F}_{22 \rightarrow mm,k}(\hat{\Sigma}_k^{-1})|}{\widehat{\text{s.e.}}[\hat{\psi}_1]} - \varsigma_k \frac{\widehat{\text{s.e.}}[\hat{\mathbb{I}}\mathbb{F}_{22 \rightarrow mm,k}(\hat{\Sigma}_k^{-1})]}{\widehat{\text{s.e.}}[\hat{\psi}_1]} > \delta : \forall m = 2, \dots, M \right\}.$$

If  $\hat{\chi}_{M,k}^{es}(\hat{\Sigma}_k^{-1}; \varsigma_k, \delta)$  fails to reject  $H_{0,k}(\delta)$  at some  $m \leq M$ , we stop the test at  $m$  and claim that we fail to reject  $H_{0,k}(\delta)$ . This “early stopping” procedure, with  $\varsigma_k$  set to  $z_{\alpha^\dagger/2}$  is a conservative  $\alpha^\dagger$ -level test under much weaker conditions:

**Proposition 5.4.** *Under the conditions of Proposition 4.3, but with (4.5) replaced by*

$$\text{Bias}_{\theta,k}(\hat{\psi}_1) \neq o \left( \mathbb{L}_{\theta,2,\hat{b},k} \mathbb{L}_{\theta,2,\hat{p},k} \left( \frac{k \log(k)}{n} \right)^{\frac{M-3}{2}} \right),$$

*under  $H_{0,k}(\delta)$ :  $\hat{\chi}_{M,k}^{es}(\hat{\Sigma}_k^{-1}; \varsigma_k, \delta)$  rejects the null with probability less than or equal to  $2(1 - \Phi(\varsigma_k))$ , as  $n \rightarrow \infty$ .*

The power of the above “early-stopping” procedure is more challenging to characterize and it may require more assumptions on the estimation biases  $\text{EB}_{\theta,m,k}(\hat{\Sigma}_k^{-1})$  for each  $m$ . We leave this problem to the future.

## 6. MONTE CARLO EXPERIMENTS

In the Monte Carlo (MC) experiments, we focus on parameter the counterfactual mean of  $Y$  when  $\{0, 1\}$ -valued  $A$  is set to 1:  $\psi^\dagger(\theta) \equiv \text{E}_\theta[Y(a = 1)] \equiv \text{E}_\theta[b(X)]$  under ignorability. Here  $p(X) = 1/\text{E}_\theta[A|W]$  is the inverse propensity score and  $b(X) = \text{E}_\theta[Y|A = 1, X]$  is the conditional mean of the outcome in the treatment group. For this parameter,  $\Sigma_k = \text{E}_\theta[A \bar{z}_k(X) \bar{z}_k(X)^\top]$  and  $\hat{\Sigma}_k = n^{-1} \sum_{i \in \text{tr}} A_i \bar{z}_k(X_i) \bar{z}_k(X_i)^\top$ . We choose  $\psi^\dagger(\theta) \equiv 0$ .

We consider two simulation setups. In simulation setup I, we draw  $N = 100,000$  i.i.d.  $X_j$  for  $j = 1, \dots, 4$  (so  $d = 4$ ) and the correlations between every two dimensions but the same marginal density  $f$  supported on  $[0, 1]$  with  $f \in \text{H\"older}(s_f = 0.1)$ , according to the algorithm described in Appendix B.1. The specific form of  $f$  is provided in Appendix B. We then simulate  $Y$  and  $A$  according to the following data generating mechanism:

$$Y \sim b(X) + N(0, 1) = \sum_{j=1}^4 \tau_{b,j} h_b(X_j; 0.25) + N(0, 1)$$

and

$$A \sim \text{Bernoulli} \left( 1/p(X) \equiv \text{expit} \left\{ \sum_{j=1}^4 \tau_{p,j} h_p(X_j; 0.25) \right\} \right)$$

where  $h_b(\cdot; 0.25)$  and  $h_p(\cdot; 0.25)$  have the forms as defined in Appendix B and hence both belong to  $\text{H\"older}(0.25+c)$  for some very small  $c > 0$ . The numerical values for  $(\tau_{b,j}, \tau_{p,j})_{j=1}^4$  are provided in Table 5. We fix half of the  $N = 100,000$

samples as the training sample so  $n_{\text{tr}} = n = 50,000$  and only consider the randomness from the estimation sample in the simulation. In simulation setup II, we consider the same data generating mechanism as in setup I except that we choose  $b(X) = \sum_{j=1}^4 \tau_{b,j} h_b(X_j; 0.6)$  and  $1/p(X) \equiv \text{expit} \left\{ \sum_{j=1}^4 \tau_{p,j} h_p(X_j; 0.6) \right\}$  where  $h_b(\cdot; 0.6)$  and  $h_p(\cdot; 0.6)$  have the forms as defined in Appendix B and hence both belong to Hölder( $0.6 + c$ ) for some very small  $c > 0$ . Thus in simulation setup II, we expect that the DRML estimator should work well and our test should not reject with high probability.

We choose D12 (or equivalently db6) Daubechies wavelets (Daubechies, 1992) at resolutions  $\ell \in (6, 7, 8)$  to form the basis functions

$$\bar{z}_k(X) = (\bar{z}_{k'}(X_1)^\top, \bar{z}_{k'}(X_2)^\top, \bar{z}_{k'}(X_3)^\top, \bar{z}_{k'}(X_4)^\top)^\top,$$

with the corresponding  $k' \in \{2^6 = 64, 2^7 = 128, 2^8 = 256\}$  and  $k \in \{64 \cdot 4 = 256, 128 \cdot 4 = 512, 256 \cdot 4 = 1024\}$ . To compute the oracle statistics and tests, we evaluate  $\Sigma_k$  through Monte Carlo integration by simulating  $L = 10^7$  independent  $(A, X)$  from the true data generating law. This same strategy will be applied when a semisupervised dataset with  $L = 10^7$  observations of  $(A, X)$  but no information on  $Y$  is available (see Section 3). To investigate the finite sample performance of the statistical procedures developed in this article, all the summary statistics of the Monte Carlo experiments are calculated based on 100 replicates. We estimate the nuisance functions  $1/p(x)$  and  $b(x)$  using generalized additive models (GAMs) (Hastie and Tibshirani, 1986, 1987). In particular, the smoothing parameters were selected by generalized cross validation, the default setup in gam function from R package mgcv (Wood et al., 2016). We choose db6 father wavelets to construct HOIF estimators when analyzing the simulated data.

**6.1. Finite sample performance of tests for  $H_{0,k}(\delta)$  (3.14).** In this section, we consider testing the null hypothesis  $H_{0,k}(\delta)$  (3.14). Henceforth we investigate the finite sample performance of the oracle statistics and tests  $\hat{\Pi}\hat{\mathbb{F}}_{22,k}(\Sigma_k^{-1})$ ,  $\hat{\psi}_{2,k}(\Sigma_k^{-1})$ , and  $\hat{\chi}_{2,k}(\Sigma_k^{-1})$ , together with the statistics and tests relying on  $\hat{\Sigma}_k^{-1}$ :  $\hat{\Pi}\hat{\mathbb{F}}_{22,k}(\hat{\Sigma}_k^{-1})$ ,  $\hat{\Pi}\hat{\mathbb{F}}_{33,k}(\hat{\Sigma}_k^{-1})$ ,  $\hat{\Pi}\hat{\mathbb{F}}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})$ ,  $\hat{\psi}_{2,k}(\hat{\Sigma}_k^{-1})$ ,  $\hat{\psi}_{3,k}(\hat{\Sigma}_k^{-1})$ ,  $\hat{\chi}_{2,k}(\hat{\Sigma}_k^{-1})$ , and  $\hat{\chi}_{3,k}(\hat{\Sigma}_k^{-1})$ .

First, we check the asymptotic normalities of  $\frac{\hat{\Pi}\hat{\mathbb{F}}_{22,k}(\Sigma_k^{-1})}{\widehat{\text{s.e.}}[\hat{\Pi}\hat{\mathbb{F}}_{22,k}(\Sigma_k^{-1})]}$ ,  $\frac{\hat{\Pi}\hat{\mathbb{F}}_{22,k}(\hat{\Sigma}_k^{-1})}{\widehat{\text{s.e.}}[\hat{\Pi}\hat{\mathbb{F}}_{22,k}(\hat{\Sigma}_k^{-1})]}$ ,  $\frac{\hat{\Pi}\hat{\mathbb{F}}_{33,k}(\hat{\Sigma}_k^{-1})}{\widehat{\text{s.e.}}[\hat{\Pi}\hat{\mathbb{F}}_{33,k}(\hat{\Sigma}_k^{-1})]}$ , and  $\frac{\hat{\Pi}\hat{\mathbb{F}}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})}{\widehat{\text{s.e.}}[\hat{\Pi}\hat{\mathbb{F}}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})]}$  through normal qq-plots displayed in Figures 1 (simulation setup I) and 2 (simulation setup II). We observe that the distributions of most of these statistics are close to normal at different  $k$ 's ( $k = 256$ : left panels;  $k = 512$ : middle panels;  $k = 1024$ : right panels). In the third row of Figures 1 and 2), we observe some deviation from normality in the tails of  $\frac{\hat{\Pi}\hat{\mathbb{F}}_{33,k}(\hat{\Sigma}_k^{-1})}{\widehat{\text{s.e.}}[\hat{\Pi}\hat{\mathbb{F}}_{33,k}(\hat{\Sigma}_k^{-1})]}$ , but it does not have significant influence on the asymptotic distribution of  $\frac{\hat{\Pi}\hat{\mathbb{F}}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})}{\widehat{\text{s.e.}}[\hat{\Pi}\hat{\mathbb{F}}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})]}$ .

In the simulation, we use nonparametric bootstrap to estimate the standard errors of  $\hat{\Pi}\hat{\mathbb{F}}_{22,k}(\Sigma_k^{-1})$ ,  $\hat{\Pi}\hat{\mathbb{F}}_{22,k}(\hat{\Sigma}_k^{-1})$ ,  $\hat{\Pi}\hat{\mathbb{F}}_{33,k}(\hat{\Sigma}_k^{-1})$  and  $\hat{\Pi}\hat{\mathbb{F}}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})$ , as described in Appendix A.3. Hence we study if the estimated standard errors of  $\hat{\Pi}\hat{\mathbb{F}}_{22,k}(\Sigma_k^{-1})$ ,  $\hat{\Pi}\hat{\mathbb{F}}_{22,k}(\hat{\Sigma}_k^{-1})$ ,  $\hat{\Pi}\hat{\mathbb{F}}_{33,k}(\hat{\Sigma}_k^{-1})$ , and  $\hat{\Pi}\hat{\mathbb{F}}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})$  by nonparametric bootstrap are close to their true standard errors (calibrated by the MC variance from 100 replicates in the simulation). We use  $B = 100$  bootstrap samples to

compute the bootstrapped standard errors as the estimated standard errors for all four statistics. In Tables 1 (simulation setup I) and 3 (simulation setup II), we display the MC standard deviations (the upper numerical values in each cell), accompanied with the MC averages of the estimated standard errors (the lower numerical values outside the parenthesis in each cell) and *MC standard deviations* of the estimated standard errors (the lower numerical values inside the parenthesis in each cell) of all three statistics at  $k = 256$  (left panel),  $k = 512$  (middle panel) and  $k = 1024$  (right panel). From Table 1 and Table 3, we observe that the estimated standard errors only slightly differ from the MC standard deviations.

Then we investigate (1) the finite sample performance of  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})$ ,  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\hat{\Sigma}_k^{-1})$ , and  $\hat{\mathbb{I}}\mathbb{F}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})$  and evaluate how close they are to  $\text{Bias}_\theta(\hat{\psi}_1)$ , which is evaluated based on the MC bias of 100 replicates, and (2) the rejection rate of the tests  $\hat{\chi}_{2,k}(\Sigma_k^{-1}; z_{0.10/2}, \delta)$ ,  $\hat{\chi}_{2,k}(\hat{\Sigma}_k^{-1}; z_{0.10/2}, \delta)$  and  $\hat{\chi}_{3,k}(\Sigma_k^{-1}; z_{0.10/2}, \delta)$  for the null hypothesis  $H_{0,k}(\delta)$  (3.14). The numerical results are shown in Tables 2 (simulation setup I) and 4 (simulation setup II).

- In the upper panel of Table 2 (simulation setup I), the first row and the third column, we display the MC bias of the DRML estimator  $\hat{\psi}_1$  and the MC average of  $\widehat{\text{s.e.}}[\hat{\psi}_1]$ , which are  $-34.26 \times 10^{-3}$  and  $8.77 \times 10^{-3}$ . Since the ratio between the bias and standard error is around 4, we expect the associated 90% Wald confidence interval  $\hat{\psi}_1 \pm z_{0.05}\widehat{\text{s.e.}}(\hat{\psi}_1)$  does not have the nominal coverage. This is indeed the case by reading from the first row, the second column of the upper panel of Table 2, showing the MC coverage probability for the 90% Wald confidence interval of  $\hat{\psi}_1$  is 0%.

Similarly, in the upper panel of Table 4 (simulation setup II), the first row and the third column, we display the MC bias of the DRML estimator  $\hat{\psi}_1$  and the MC average of  $\widehat{\text{s.e.}}[\hat{\psi}_1]$ , which are  $-8.88 \times 10^{-3}$  and  $8.10 \times 10^{-3}$ . Since the ratio between the bias and standard error is around 1, we expect the associated 90% Wald confidence interval  $\hat{\psi}_1 \pm z_{0.05}\widehat{\text{s.e.}}(\hat{\psi}_1)$  has close to nominal coverage. This is indeed the case by reading from the first row, the second column of the upper panel of Table 4, showing the MC coverage probability for the 90% Wald confidence interval of  $\hat{\psi}_1$  is 83%.

- In the second column of the upper panel of Table 2, we display the MC averages of  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})$  and the MC averages of its estimated standard errors; in the middle column, we display those of  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\hat{\Sigma}_k^{-1})$ ; and in the lower panel, we display those of  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\hat{\Sigma}_k^{-1}) + \hat{\mathbb{I}}\mathbb{F}_{33,k}(\hat{\Sigma}_k^{-1})$  after third-order bias correction. In the upper panel, we observe that increasing  $k$  from 256 to 512 does improve the amount of bias recovered by  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})$  (from  $-18.76 \times 10^{-3}$  to  $-25.34 \times 10^{-3}$ ), but there is no obvious difference in the MC averages of  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})$  between  $k = 512$  and  $k = 1024$  ( $-25.34 \times 10^{-3}$  and  $-25.43 \times 10^{-3}$ , about 75% of the total bias). The MC average of the estimated standard error increases with  $k$ , from  $2.45 \times 10^{-3}$  at  $k = 256$  to  $3.43 \times 10^{-3}$  at  $k = 1024$ . This is consistent with the theoretical prediction that the variability of  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})$  should grow with  $k$ . In the middle, we observe very similar numerical results between  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})$  and  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\hat{\Sigma}_k^{-1})$  for all the statistics that we are interested in. Interestingly, we did not see much improvement of using  $\hat{\mathbb{I}}\mathbb{F}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})$  instead of  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\hat{\Sigma}_k^{-1})$ , when compared to  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})$ , at least in this example.

In the second column of the upper panel of Table 4, in the upper panel, we also observe that increasing  $k$  from 256 to 512 also slightly improve the amount of bias recovered by  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})$  (from  $-4.54 \times 10^{-3}$  to  $-4.94 \times 10^{-3}$ ). The MC average of the estimated standard error increases with  $k$ , from  $1.51 \times 10^{-3}$  at  $k = 256$  to  $2.43 \times 10^{-3}$  at  $k = 1024$ . In the middle panel, we observe very similar numerical results between  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})$  and  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\hat{\Sigma}_k^{-1})$  for all the statistics that we are interested in. Interestingly, we again did not see much improvement of using  $\hat{\mathbb{I}}\mathbb{F}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})$  instead of  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\hat{\Sigma}_k^{-1})$ , when compared to  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})$ , at least in this example.

- In the third columns of Table 2, we display the MC coverage probabilities of the 90% two-sided confidence intervals of  $\hat{\psi}_{2,k}(\Sigma_k^{-1})$  (upper panel) and  $\hat{\psi}_{3,k}(\hat{\Sigma}_k^{-1})$  (lower panel). Comparing the oracle bias-corrected estimator  $\hat{\psi}_{2,k}(\Sigma_k^{-1})$  vs.  $\hat{\psi}_1$ , the coverage probability improves from 0% to 33% (82%) at  $k = 256$  (at  $k = 1024$ ). Similar observations can be made for  $\hat{\psi}_{2,k}(\hat{\Sigma}_k^{-1})$  and  $\hat{\psi}_{3,k}(\hat{\Sigma}_k^{-1})$  as well when  $\Sigma_k$  is unknown.

In the third columns of Table 4, we display the MC coverage probabilities of the 90% two-sided confidence intervals of  $\hat{\psi}_{2,k}(\Sigma_k^{-1})$  (upper panel) and  $\hat{\psi}_{3,k}(\hat{\Sigma}_k^{-1})$  (lower panel). Comparing the oracle bias-corrected estimator  $\hat{\psi}_{2,k}(\Sigma_k^{-1})$  vs.  $\hat{\psi}_1$ , the coverage probability improves from 83% to 100% at  $k = 256$  and  $k = 1024$ . Similar observations can be made for  $\hat{\psi}_{2,k}(\hat{\Sigma}_k^{-1})$  and  $\hat{\psi}_{3,k}(\hat{\Sigma}_k^{-1})$  as well when  $\Sigma_k$  is unknown.

- In the fourth columns of Tables 2 and 4, we display the MC biases of  $\hat{\psi}_{2,k}(\Sigma_k^{-1})$  (upper panel) and  $\hat{\psi}_{3,k}(\hat{\Sigma}_k^{-1})$  (lower panel), together with their MC standard deviations (in the parentheses). As expected, the standard deviations of the estimators after bias correction are very similar to those of  $\hat{\psi}_1$ . This indeed confirms that we are able to correct bias without drastically inflating the variance.
- In the fifth column of Table 2, we display the MC rejection rates of the test statistic: the upper panel shows the MC rejection rates of the oracle test  $\hat{\chi}_{2,k}(\Sigma_k^{-1}; z_{0.10/2}, \delta = 3/4)$  and the lower panel shows the MC rejection rates of the test  $\hat{\chi}_{3,k}(\hat{\Sigma}_k^{-1}; z_{0.10/2}, \delta = 3/4)$ . These simulation results demonstrate that when a “good” basis (db6 father wavelets) is used, our proposed test does have power to reject the null hypothesis of actual interest  $H_0(\delta)$  (3.12) that the ratio between the bias of  $\hat{\psi}_1$  and its standard error is lower than  $\delta$ . All the rejection rates in Table 2 are 100% when  $\delta = 3/4$ , regardless of whether we are using the oracle test or not. When  $\delta = 2$ , we obtain nontrivial rejection rates. This is as expected because the ratio between  $\hat{\mathbb{I}}\mathbb{F}_{22,k}$  and  $\widehat{\text{s.e.}}(\hat{\psi}_1)$  is around 2 to 3. In the fifth column of Table 4, all the rejection rates are 0% when  $\delta = 3/4$ . This is also as expected because the ratio between  $\hat{\mathbb{I}}\mathbb{F}_{22,k}$  and  $\widehat{\text{s.e.}}(\hat{\psi}_1)$  is close to 1/2.

## 7. CONCLUDING REMARKS

In this paper, we developed almost assumption-free tests that can criticize if the Wald confidence interval  $\text{Cl}_\alpha(\hat{\psi}_1)$  associated with an DRML estimator  $\hat{\psi}_1$  of a DR functional has the nominal coverage. We also develop a test hierarchy that can perform similar tasks for higher-order estimators  $\hat{\psi}_{m,k}$  based on higher-order influence function estimators (Robins et al., 2008). In this paper, we fix the basis  $\bar{z}$  but one might perform multiple testing if a family of candidate



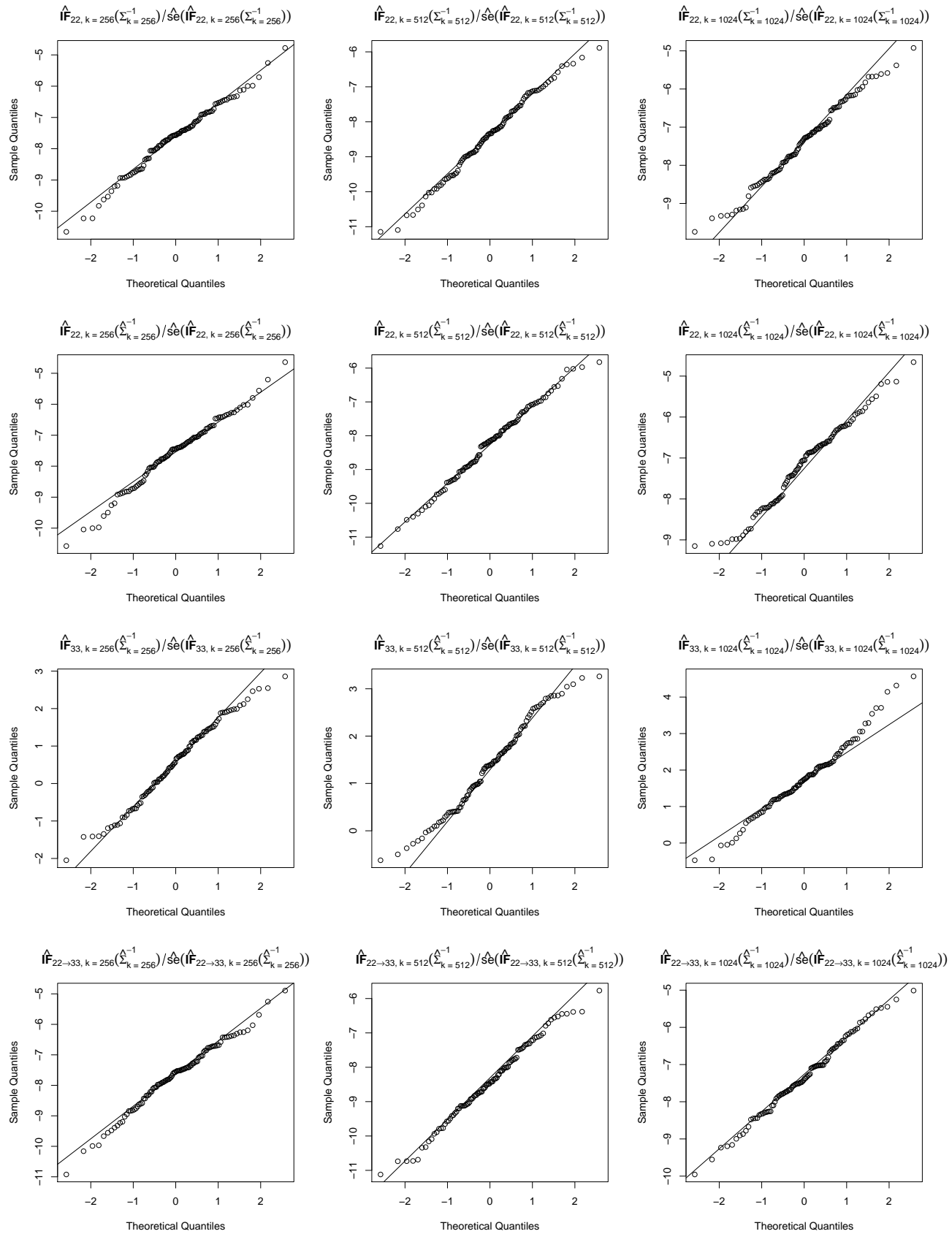


FIGURE 1. For data generating mechanism in simulation setup I: normal qq-plots for  $\hat{\mathbb{F}}_{22, k}(\hat{\Sigma}_k^{-1}) / \widehat{\text{se}}[\hat{\mathbb{F}}_{22, k}(\hat{\Sigma}_k^{-1})]$  (the first row),  $\hat{\mathbb{F}}_{22, k}(\hat{\Sigma}_k^{-1}) / \widehat{\text{se}}[\hat{\mathbb{F}}_{22, k}(\hat{\Sigma}_k^{-1})]$  (the second row),  $\hat{\mathbb{F}}_{33, k}(\hat{\Sigma}_k^{-1}) / \widehat{\text{se}}[\hat{\mathbb{F}}_{33, k}(\hat{\Sigma}_k^{-1})]$  (the third row) and  $\hat{\mathbb{F}}_{22 \rightarrow 33, k}(\hat{\Sigma}_k^{-1}) / \widehat{\text{se}}[\hat{\mathbb{F}}_{22 \rightarrow 33, k}(\hat{\Sigma}_k^{-1})]$  (the fourth row) for  $k = 256$  (left panels),  $k = 512$  (middle panels) and  $k = 1024$  (right panels).

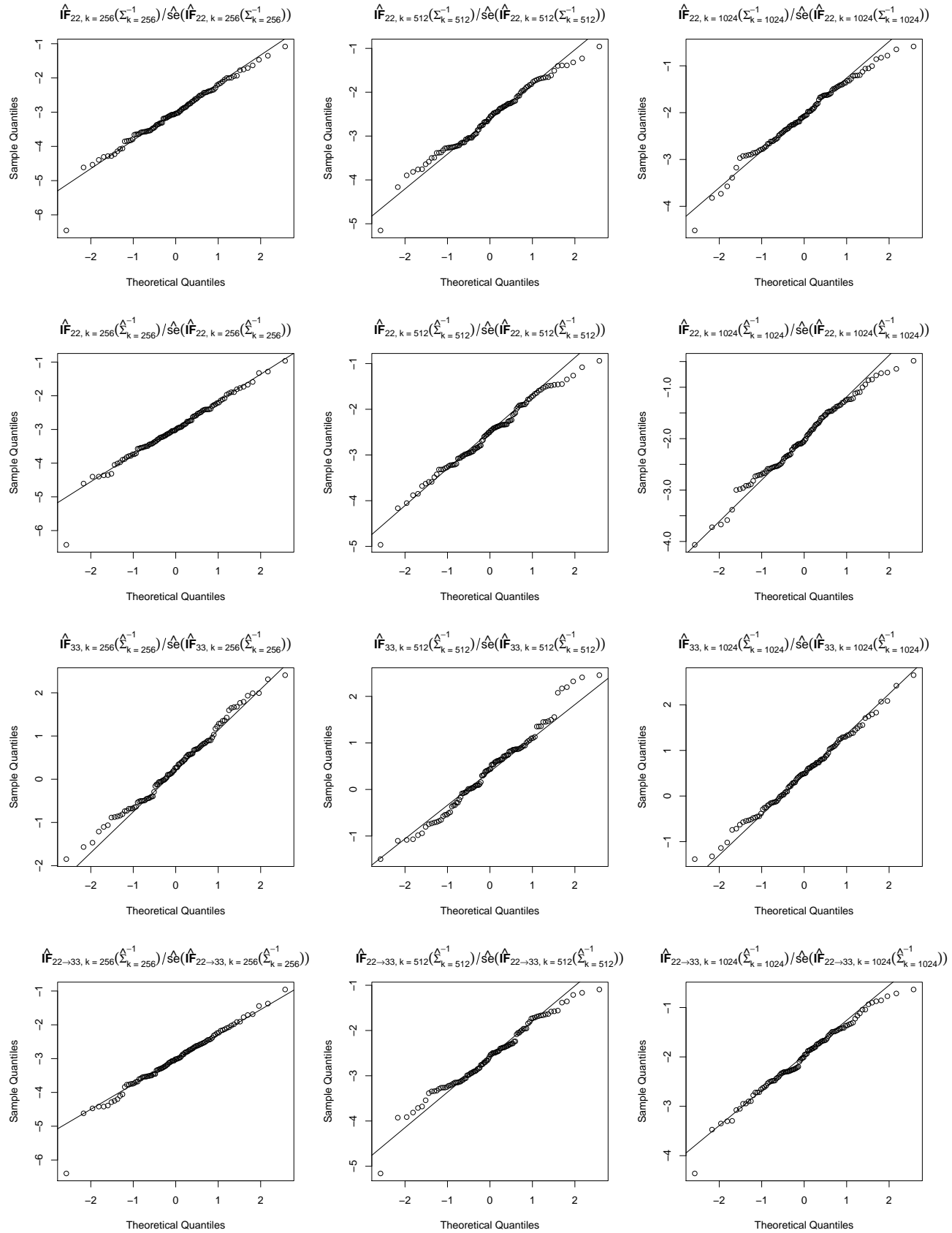


FIGURE 2. For data generating mechanism in simulation setup II: normal qq-plots for  $\hat{\mathbb{IF}}_{22,k}(\hat{\Sigma}_k^{-1})/\widehat{\text{se}}[\hat{\mathbb{IF}}_{22,k}(\hat{\Sigma}_k^{-1})]$  (the first row),  $\hat{\mathbb{IF}}_{22,k}(\hat{\Sigma}_k^{-1})/\widehat{\text{se}}[\hat{\mathbb{IF}}_{22,k}(\hat{\Sigma}_k^{-1})]$  (the second row),  $\hat{\mathbb{IF}}_{33,k}(\hat{\Sigma}_k^{-1})/\widehat{\text{se}}[\hat{\mathbb{IF}}_{33,k}(\hat{\Sigma}_k^{-1})]$  (the third row) and  $\hat{\mathbb{IF}}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})/\widehat{\text{se}}[\hat{\mathbb{IF}}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})]$  (the fourth row) for  $k = 256$  (left panels),  $k = 512$  (middle panels) and  $k = 1024$  (right panels).

$k$	256	512	1024
$\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})$	2.437 2.268 (0.218)	3.011 2.716 (0.268)	3.247 2.841 (0.322)
$\hat{\mathbb{I}}\mathbb{F}_{22,k}(\hat{\Sigma}_k^{-1})$	2.456 2.271 (0.218)	3.075 2.778 (0.278)	3.546 3.032 (0.378)
$\hat{\mathbb{I}}\mathbb{F}_{33,k}(\hat{\Sigma}_k^{-1})$	0.495 0.603 (0.0598)	0.814 1.126 (0.126)	1.518 1.998 (0.285)
$\hat{\mathbb{I}}\mathbb{F}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})$	2.282 2.251 (0.227)	2.745 2.738 (0.290)	2.743 3.009 (0.427)

TABLE 1. For data generating mechanism in simulation setup I: We reported the MC standard deviations  $\times 10^{-3}$  (upper values in each cell), the MC averages of the estimated standard errors  $\times 10^{-3}$  (lower values in each cell outside the parenthesis) and the MC standard deviations of the estimated standard errors  $\times 10^{-3}$  (lower values in each cell inside the parenthesis) through nonparametric bootstrap resampling for  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})$ ,  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\hat{\Sigma}_k^{-1})$ ,  $\hat{\mathbb{I}}\mathbb{F}_{33,k}(\hat{\Sigma}_k^{-1})$  and  $\hat{\mathbb{I}}\mathbb{F}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})$ .

basis functions is available to the analysts. It is also interesting to investigate the possibility and consequence of selecting basis functions  $\bar{z}(\cdot)$  from a huge dictionary  $\mathcal{V}$  of functions using data-driven methods, the goal of which is to maximize the power of the assumption-free test.

We conclude the paper with another interesting open problem. The procedure developed in this paper relies heavily on higher-order U-statistics (in particular for Sections 4 and 5) to correct for the estimation bias, which are difficult to compute when the order becomes large. There have been a series of papers (Kennedy, 2020; Newey and Robins, 2018) advocating the use of multiple sample splitting to avoid higher order U-statistics in the problem of estimating DR functionals. However, the analyses and sample splitting schemes depend on the specific functional of interest and it generally requires stronger untestable assumptions in order to completely replace higher order U-statistics. It will be interesting to develop a unified treatment of multiple sample splitting in the assumption-free setup for the entire class of DR functionals as an alternative route that computes faster than higher order U-statistics.

#### REFERENCES

- Miguel A Arcones and Evarist Gine. On the bootstrap of  $U$  and  $V$  statistics. *The Annals of Statistics*, 20(2):655–674, 1992.
- Rose Baker. An order-statistics-based method for constructing multivariate distributions with fixed marginals. *Journal of Multivariate Analysis*, 99(10):2312–2327, 2008.
- Heejung Bang and James M Robins. Doubly robust estimation in missing data and causal inference models. *Biometrics*, 61(4):962–973, 2005.

k	$\hat{\mathbb{F}}_{22,k}(\Sigma_k^{-1}) \times 10^{-3}$	MC Coverage		$\hat{\chi}_{2,k}(\Sigma_k^{-1}; z_{0.10/2}, \delta = 3/4)$
		$(\hat{\psi}_{2,k}(\Sigma_k^{-1}) \text{ 90\% Wald CI})$	$\text{Bias}(\hat{\psi}_{2,k}(\Sigma_k^{-1})) \times 10^{-3}$	
0	NA (NA)	0%	-34.26 (8.77)	NA
256	-18.76 (2.27)	31%	-15.50 (8.60)	100%
512	-25.34 (2.72)	83%	-8.92 (8.61)	100%
1024	-25.43 (2.84)	82%	-8.83 (8.79)	100%
k	$\hat{\mathbb{F}}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1}) \times 10^{-3}$	MC Coverage		$\hat{\chi}_{3,k}(\hat{\Sigma}_k^{-1}; z_{0.10/2}, \delta = 3/4)$
		$(\hat{\psi}_{3,k}(\hat{\Sigma}_k^{-1}) \text{ 90\% Wald CI})$	$\text{Bias}(\hat{\psi}_{3,k}(\hat{\Sigma}_k^{-1})) \times 10^{-3}$	
256	-18.54 (2.25)	30%	-15.72 (8.60)	100%
512	-24.65 (2.74)	81%	-9.61 (8.60)	100%
1024	-23.82 (3.01)	76%	-10.44 (8.76)	100%

TABLE 2. For data generating mechanism in simulation setup I: We reported the MC averages of point estimates and standard errors (the first column in each panel) of  $\hat{\mathbb{F}}_{22,k}(\Sigma_k^{-1})$  (upper panel) and  $\hat{\mathbb{F}}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})$  (lower panel), together with the coverage probabilities of two-sided 90% Wald confidence intervals (the second column in each panel) of  $\hat{\psi}_{2,k}(\Sigma_k^{-1})$  (upper panel) and  $\hat{\psi}_{3,k}(\hat{\Sigma}_k^{-1})$  (lower panel), the MC biases and MC standard deviations (the third column in each panel) of  $\hat{\psi}_{2,k}(\Sigma_k^{-1})$  (upper panel) and  $\hat{\psi}_{3,k}(\hat{\Sigma}_k^{-1})$  (lower panel) and the empirical rejection rates (the fourth column in each panel) of  $\hat{\chi}_{2,k}^{(1)}(\Sigma_k^{-1}; z_{0.10/2}, \delta = 3/4)$  (upper panel) and  $\hat{\chi}_{3,k}^{(1)}(\hat{\Sigma}_k^{-1}; z_{0.10/2}, \delta = 3/4)$  (lower panel). In the upper panel, for  $k = 0$ ,  $\hat{\psi}_{2,k=0} \equiv \hat{\psi}_1$ .

Alexandre Belloni, Victor Chernozhukov, Denis Chetverikov, and Kengo Kato. Some new asymptotic theory for least squares series: Pointwise and uniform results. *Journal of Econometrics*, 186(2):345–366, 2015.

Rabi N Bhattacharya and Jayanta K Ghosh. A class of  $U$ -statistics and asymptotic normality of the number of  $k$ -clusters. *Journal of Multivariate Analysis*, 43(2):300–330, 1992.

Peter J Bickel, Chris A J Klaassen, Ya’acov Ritov, and Jon A Wellner. *Efficient and Adaptive Estimation for Semiparametric Models*. Johns Hopkins series in the mathematical sciences. Springer New York, 1998. ISBN 9780387984735.

Jelena Bradic, Victor Chernozhukov, Whitney K Newey, and Yinchu Zhu. Minimax semiparametric learning with approximate sparsity. *arXiv preprint arXiv:1912.12213*, 2019.

Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins. Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1):C1–C68, 2018a.

Victor Chernozhukov, Whitney Newey, and James Robins. Double/de-biased machine learning using regularized Riesz representers. *arXiv preprint arXiv:1802.08667*, 2018b.

$k$	256	512	1024
$\hat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})$	1.166 1.209 (0.156)	1.367 1.408 (0.191)	1.673 1.584 (0.267)
$\hat{\mathbb{I}\mathbb{F}}_{22,k}(\hat{\Sigma}_k^{-1})$	1.187 1.219 (0.159)	1.433 1.439 (0.198)	1.828 1.697 (0.285)
$\hat{\mathbb{I}\mathbb{F}}_{33,k}(\hat{\Sigma}_k^{-1})$	0.251 0.331 (0.0451)	0.417 0.629 (0.0997)	0.854 1.251 (0.247)
$\hat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})$	1.109 1.186 (0.156)	1.256 1.346 (0.191)	1.488 1.560 (0.383)

TABLE 3. For data generating mechanism in simulation setup II: We reported the MC standard deviations  $\times 10^{-3}$  (upper values in each cell), the MC averages of the estimated standard errors  $\times 10^{-3}$  (lower values in each cell outside the parenthesis) and the MC standard deviations of the estimated standard errors  $\times 10^{-3}$  (lower values in each cell inside the parenthesis) through nonparametric bootstrap resampling for  $\hat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})$ ,  $\hat{\mathbb{I}\mathbb{F}}_{22,k}(\hat{\Sigma}_k^{-1})$ ,  $\hat{\mathbb{I}\mathbb{F}}_{33,k}(\hat{\Sigma}_k^{-1})$  and  $\hat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})$ .

Victor Chernozhukov, Whitney K Newey, and Rahul Singh. Learning L2 continuous regression functionals via regularized Riesz representers. *arXiv preprint arXiv:1809.05224*, 2018c.

Ingrid Daubechies. *Ten lectures on wavelets*, volume 61. SIAM, 1992.

Max H Farrell, Tengyuan Liang, and Sanjog Misra. Deep neural networks for estimation and inference. *arXiv preprint arXiv:1809.09953*, 2019.

Wolfgang Härdle, Gerard Kerkycharian, Dominique Picard, and Alexander Tsybakov. *Wavelets, approximation, and statistical applications*, volume 129. Springer Science & Business Media, 1998.

Trevor Hastie and Robert Tibshirani. Generalized additive models. *Statistical Science*, 1(3):297–310, 1986.

Trevor Hastie and Robert Tibshirani. Generalized additive models: some applications. *Journal of the American Statistical Association*, 82(398):371–386, 1987.

Wassily Hoeffding. A class of statistics with asymptotically normal distribution. *The Annals of Mathematical Statistics*, 19(3):293–325, 1948.

Marie Huskova and Paul Janssen. Consistency of the generalized bootstrap for degenerate  $U$ -statistics. *The Annals of Statistics*, 21(4):1811–1823, 1993.

Edward H Kennedy. Optimal doubly robust estimation of heterogeneous causal effects. *arXiv preprint arXiv:2004.14497*, 2020.

Lin Liu, Rajarsh Mukherjee, James M Robins, and Eric Tchetgen Tchetgen. Adaptive estimation of nonparametric functionals. *arXiv preprint arXiv:1608.01364*, 2016.

Stéphane Mallat. *A wavelet tour of signal processing*. Elsevier, 1999.

k	MC Coverage			
	$\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1}) \times 10^{-3}$	$(\hat{\psi}_{2,k}(\Sigma_k^{-1}) \text{ 90\% Wald CI})$	$\text{Bias}(\hat{\psi}_{2,k}(\Sigma_k^{-1})) \times 10^{-3}$	$\hat{\chi}_{2,k}(\Sigma_k^{-1}; z_{0.10/2}, \delta = 3/4)$
0	NA (NA)	83%	-8.88 (8.10)	NA
256	-4.54 (1.21)	99%	-4.34 (8.13)	0%
512	-4.89 (1.41)	99%	-3.99 (8.21)	0%
1024	-4.94 (1.58)	100%	-3.94 (8.35)	0%

---

k	MC Coverage			
	$\hat{\mathbb{I}}\mathbb{F}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1}) \times 10^{-3}$	$(\hat{\psi}_{3,k}(\hat{\Sigma}_k^{-1}) \text{ 90\% Wald CI})$	$\text{Bias}(\hat{\psi}_{3,k}(\hat{\Sigma}_k^{-1})) \times 10^{-3}$	$\hat{\chi}_{3,k}(\hat{\Sigma}_k^{-1}; z_{0.10/2}, \delta = 3/4)$
256	-4.49 (1.19)	99%	-4.39 (8.13)	0%
512	-4.76 (1.35)	99%	-4.12 (8.20)	0%
1024	-4.62 (1.56)	100%	-4.26 (8.32)	0%

TABLE 4. For data generating mechanism in simulation setup II: We reported the MC averages of point estimates and standard errors (the first column in each panel) of  $\hat{\mathbb{I}}\mathbb{F}_{22,k}(\Sigma_k^{-1})$  (upper panel) and  $\hat{\mathbb{I}}\mathbb{F}_{22 \rightarrow 33,k}(\hat{\Sigma}_k^{-1})$  (lower panel), together with the coverage probabilities of two-sided 90% Wald confidence intervals (the second column in each panel) of  $\hat{\psi}_{2,k}(\Sigma_k^{-1})$  (upper panel) and  $\hat{\psi}_{3,k}(\hat{\Sigma}_k^{-1})$  (lower panel), the MC biases and MC standard deviations (the third column in each panel) of  $\hat{\psi}_{2,k}(\Sigma_k^{-1})$  (upper panel) and  $\hat{\psi}_{3,k}(\hat{\Sigma}_k^{-1})$  (lower panel) and the empirical rejection rates (the fourth column in each panel) of  $\hat{\chi}_{2,k}^{(1)}(\Sigma_k^{-1}; z_{0.10/2}, \delta = 3/4)$  (upper panel) and  $\hat{\chi}_{3,k}^{(1)}(\hat{\Sigma}_k^{-1}; z_{0.10/2}, \delta = 3/4)$  (lower panel). In the upper panel, for  $k = 0$ ,  $\hat{\psi}_{2,k=0} \equiv \hat{\psi}_1$ .

James E Mosimann. On the compound multinomial distribution, the multivariate  $\beta$ -distribution, and correlations among proportions. *Biometrika*, 49(1-2):65–82, 1962.

Rajarshi Mukherjee, Whitney K Newey, and James M Robins. Semiparametric efficient empirical higher order influence function estimators. *arXiv preprint arXiv:1705.07577*, 2017.

Justin T Newcomer, Nagaraj K Neerchal, and Jorge G Morel. Computation of higher order moments from two multinomial overdispersion likelihood models. Technical report, Department of Mathematics and Statistics, University of Maryland, Baltimore, USA, 2008.

Whitney K Newey. Semiparametric efficiency bounds. *Journal of Applied Econometrics*, 5(2):99–135, 1990.

Whitney K Newey and James M Robins. Cross-fitting and fast remainder rates for semiparametric estimation. *arXiv preprint arXiv:1801.09138*, 2018.

Ya’acov Ritov, Peter J Bickel, Anthony C Gamst, and Bastiaan Jan Korneel Kleijn. The Bayesian analysis of complex, high-dimensional models: Can it be CODA? *Statistical Science*, 29(4):619–639, 2014.

James Robins, Lingling Li, Eric Tchetgen Tchetgen, and Aad van der Vaart. Higher order influence functions and minimax estimation of nonlinear functionals. In *Probability and Statistics: Essays in Honor of David A. Freedman*,



- pages 335–421. Institute of Mathematical Statistics, 2008.
- James Robins, Lingling Li, Eric Tchetgen Tchetgen, and Aad W van der Vaart. Quadratic semiparametric von Mises calculus. *Metrika*, 69(2-3):227–247, 2009a.
- James Robins, Eric Tchetgen Tchetgen, Lingling Li, and Aad van der Vaart. Semiparametric minimax rates. *Electronic Journal of Statistics*, 3:1305–1321, 2009b.
- James Robins, Lingling Li, Eric Tchetgen Tchetgen, and Aad van der Vaart. Technical report: Higher order influence functions and minimax estimation of nonlinear functionals. *arXiv preprint arXiv:1601.05820*, 2016.
- James M Robins and Ya’acov Ritov. Toward a curse of dimensionality appropriate (CODA) asymptotic theory for semi-parametric models. *Statistics in Medicine*, 16(3):285–319, 1997.
- James M Robins and Andrea Rotnitzky. Comments on “Inference for semiparametric models: some questions and an answer”. *Statistica Sinica*, 11(4):920–936, 2001.
- James M Robins, Lingling Li, Rajarshi Mukherjee, Eric Tchetgen Tchetgen, and Aad van der Vaart. Minimax estimation of a functional on a structured high-dimensional model. *The Annals of Statistics*, 45(5):1951–1987, 2017.
- Andrea Rotnitzky, Ezequiel Smucler, and James M Robins. Characterization of parameters with a mixed bias property. *arXiv preprint arXiv:1904.03725*, 2019.
- Daniel O Scharfstein, Andrea Rotnitzky, and James M Robins. Adjusting for nonignorable drop-out using semiparametric nonresponse models. *Journal of the American Statistical Association*, 94(448):1096–1120, 1999a.
- Daniel O Scharfstein, Andrea Rotnitzky, and James M Robins. Rejoinder. *Journal of the American Statistical Association*, 94(448):1135–1146, 1999b.
- Ezequiel Smucler, Andrea Rotnitzky, and James M Robins. A unifying approach for doubly-robust  $\ell_1$  regularized estimation of causal contrasts. *arXiv preprint arXiv:1904.03737*, 2019.
- Ingo Steinwart and Andreas Christmann. *Support vector machines*. Springer Science & Business Media, 2008.
- Aad van der Vaart. Part III: Semiparametric statistics. *Lectures on Probability Theory and Statistics*, pages 331–457, 2002.
- Aad van der Vaart. Higher order tangent spaces and influence functions. *Statistical Science*, 29(4):679–686, 2014.
- Aad W van der Vaart. *Asymptotic statistics*, volume 3. Cambridge University Press, 1998.
- Vladimir N Vapnik. *Statistical learning theory*. Adaptive and learning systems for signal processing, communications, and control. Wiley, 1998. ISBN 9780471030034.
- Vladimir N Vapnik. *The nature of statistical learning theory*. Springer Science & Business Media, 2013.
- Simon N Wood, Natalya Pya, and Benjamin Säfken. Smoothing parameter and model selection for general smooth models. *Journal of the American Statistical Association*, 111(516):1548–1563, 2016.

## APPENDIX A. DERIVATIONS AND PROOFS

**A.1. Higher order influence functions of  $\tilde{\psi}_k(\theta)$ .** In this section, we give an almost self-contained derivation of HOIFs of  $\tilde{\psi}_k(\theta)$ , except for [Robins et al. \(2008, Theorem 2.3\)](#), which provides the “calculus” of deriving  $m$ -th order influence functions from  $(m - 1)$ -th order influence functions. We give the formulae for HOIFs of  $\tilde{\psi}_k(\theta)$  in the theorem below.

**Theorem A.1** (HOIFs of  $\tilde{\psi}_k(\theta)$ <sup>5</sup>). *Under the conditions of Theorem 2.2, the HOIFs of order  $m$  of  $\tilde{\psi}_k(\theta)$  are  $\mathbb{IF}_{m,k}(\Sigma_k^{-1}) = \mathbb{IF}_1 - \mathbb{IF}_{22 \rightarrow mm,k}(\Sigma_k^{-1})$ , where  $\mathbb{IF}_{\ell \ell \rightarrow mm,k}(\Sigma_k^{-1}) := \sum_{j=\ell}^m \mathbb{IF}_{jj,k}(\Sigma_k^{-1})$ ,*

$$(A.1) \quad \mathbb{IF}_{22,k}(\Sigma_k^{-1}) := \frac{1}{n(n-1)} \sum_{1 \leq i_1 \neq i_2 \leq n} \text{IF}_{22,k,\bar{i}_2}(\Sigma_k^{-1})$$

and for  $j > 2$ ,

$$(A.2) \quad \mathbb{IF}_{jj,k}(\Sigma_k^{-1}) := \frac{(n-j)!}{n!} \sum_{1 \leq i_1 \neq \dots \neq i_j \leq n} \text{IF}_{jj,k,\bar{i}_j}(\Sigma_k^{-1}).$$

Here

$$(A.3) \quad \text{IF}_{22,k,\bar{i}_2}(\Sigma_k^{-1}) \equiv \text{IF}_{22,\tilde{\psi}_k,\bar{i}_2}(\theta) = \left[ \mathcal{E}_{1,\tilde{b}_{k,\theta}}(\bar{z}_k)(O) \right]_{i_1}^\top \Sigma_k^{-1} \left[ \mathcal{E}_{2,\tilde{p}_{k,\theta}}(\bar{z}_k)(O) \right]_{i_2},$$

and

$$(A.4) \quad \begin{aligned} \text{IF}_{jj,k,\bar{i}_j}(\Sigma_k^{-1}) &\equiv \text{IF}_{jj,\tilde{\psi}_k,\bar{i}_j}(\theta) \\ &= (-1)^j \left[ \mathcal{E}_{1,\tilde{b}_{k,\theta}}(\bar{z}_k)(O) \right]_{i_1}^\top \left\{ \prod_{s=3}^j \Sigma_k^{-1} \left( [S_{bp}\bar{z}_k(X)\bar{z}_k(X)^\top]_{i_s} - \Sigma_k \right) \right\} \Sigma_k^{-1} \left[ \mathcal{E}_{2,\tilde{p}_{k,\theta}}(\bar{z}_k)(O) \right]_{i_2}. \end{aligned}$$

are the kernels of second order influence function  $\mathbb{IF}_{22,k}(\Sigma_k^{-1})$  and  $j$ -th order influence function  $\mathbb{IF}_{jj,k}(\Sigma_k^{-1})$  respectively, where  $\bar{i}_j := \{i_1, \dots, i_j\}$ .

*Proof.* Without loss of generality, we assume  $P_\theta(S_{bp} \geq 0) = 1$ . First, let’s recall the results from Theorem 2.1 (or [Rotnitzky et al. \(2019, Theorems 2 \(ii\)\)](#)), the influence function of  $\psi(\theta)$  has the following form

$$\text{IF}_{1,\psi(\theta)} = S_{bp}b(X)p(X) + m_1(O, b) + m_2(O, p) + S_0 - \psi(\theta),$$

where the linear functionals  $m_1(O, h)$  and  $m_2(O, h)$  have the Riesz representers  $\mathcal{R}_1(X)$  and  $\mathcal{R}_2(X)$  respectively. By definition of Riesz representer, they must obey the following identities:

$$\begin{cases} E_\theta[m_1(O, h)] = E_\theta[h(X)\mathcal{R}_1(X)] \\ E_\theta[m_2(O, h)] = E_\theta[h(X)\mathcal{R}_2(X)] \end{cases}$$

for any  $h \in L_2(g)$ .

<sup>5</sup> The definitions of  $\mathbb{IF}_{jj,k}(\Sigma_k^{-1})$  for  $j \geq 2$  differ from those in [Robins et al. \(2008\)](#) in the sign for notational convenience.

To derive the HOIFs for  $\tilde{\psi}_k(\theta)$ , we need its influence function. Recall from Definition 3.1 that  $\tilde{\psi}_k(\theta)$  is of the following form:

$$\begin{aligned}
& \tilde{\psi}_k(\theta) \\
& \equiv \mathbb{E}_\theta \left[ \mathcal{H} \left( \tilde{b}_{k,\theta}, \tilde{p}_{k,\theta} \right) \right] \\
& = \mathbb{E}_\theta \left[ S_{bp} \tilde{b}_{k,\theta}(X) \tilde{p}_{k,\theta}(X) + m_1(O, \tilde{b}_{k,\theta}) + m_2(O, \tilde{p}_{k,\theta}) + S_0 \right] \\
& = \mathbb{E}_\theta \left[ \begin{aligned} & S_{bp} \left( \hat{b}(X) + \tilde{\zeta}_{b,k}^\top(\theta) \bar{z}_k(X) \right) \left( \hat{p}(X) + \tilde{\zeta}_{p,k}^\top(\theta) \bar{z}_k(X)^\top \right) \\ & + m_1(O, \hat{b} + \tilde{\zeta}_{b,k}^\top(\theta) \bar{z}_k) + m_2(O, \hat{p} + \tilde{\zeta}_{p,k}^\top(\theta) \bar{z}_k) + S_0 \end{aligned} \right] \\
& = \mathbb{E}_\theta \left[ \begin{aligned} & S_{bp} \left( \hat{b}(X) - \Sigma_k^{-1} \mathbb{E}_\theta \left[ S_{bp} \hat{b}(X) \bar{z}_k(X)^\top + m_2(O, \bar{z}_k(X))^\top \right] \bar{z}_k(X) \right) \\ & \times \left( \hat{p}(X) - \Sigma_k^{-1} \mathbb{E}_\theta \left[ S_{bp} \hat{p}(X) \bar{z}_k(X)^\top + m_1(O, \bar{z}_k(X))^\top \right] \bar{z}_k(X) \right) \\ & + m_1(O, \hat{b} - \Sigma_k^{-1} \mathbb{E}_\theta \left[ S_{bp} \hat{b}(X) \bar{z}_k(X)^\top + m_2(O, \bar{z}_k(X))^\top \right] \bar{z}_k) \\ & + m_2(O, \hat{p} - \Sigma_k^{-1} \mathbb{E}_\theta \left[ S_{bp} \hat{p}(X) \bar{z}_k(X)^\top + m_1(O, \bar{z}_k(X))^\top \right] \bar{z}_k) + S_0 \end{aligned} \right].
\end{aligned}$$

Then

$$\begin{aligned}
\text{IF}_{1, \tilde{\psi}_k(\theta), i}(\theta) &= \mathcal{H}(\tilde{b}_k(X, \theta), \tilde{p}_k(X, \theta)) - \tilde{\psi}_k(\theta) \\
&+ \mathbb{E}_\theta \left[ \partial \mathcal{H} \left( b_k^*(X, \tilde{\zeta}_{b,k}(\theta)), p_k^*(X, \tilde{\zeta}_{b,k}(\theta)) \right) / \partial \bar{\zeta}_{b,k}^\top \right] \text{IF}_{1, \tilde{\zeta}_{b,k}}(\theta) \\
&+ \mathbb{E}_\theta \left[ \partial \mathcal{H} \left( b_k^*(X, \tilde{\zeta}_{b,k}(\theta)), p_k^*(X, \tilde{\zeta}_{p,k}(\theta)) \right) / \partial \bar{\zeta}_{p,k}^\top \right] \text{IF}_{1, \tilde{\zeta}_{p,k}}(\theta) \\
&= \mathcal{H}(\tilde{b}_{k,\theta}, \tilde{p}_{k,\theta}) - \tilde{\psi}_k(\theta) \\
&= S_{bp,i} \tilde{b}_{k,\theta}(X_i) \tilde{p}_{k,\theta}(X_i) + m_1(O_i, \tilde{b}_{k,\theta}) + m_2(O_i, \tilde{p}_{k,\theta}) + S_{0,i} - \tilde{\psi}_k(\theta).
\end{aligned}$$

where the second equality follows from the definitions of  $\tilde{\zeta}_{b,k}(\theta)$  and  $\tilde{\zeta}_{p,k}(\theta)$  in equation (3.1). Then the influence function of  $-\text{IF}_{1, \tilde{\psi}_k(\theta), i}(\theta)$ , denoted as  $\text{IF}_{1, -\text{if}_{1, \tilde{\psi}_k(\theta), i_1}, i_2}(\theta)$ , can be calculated as follows:

$$\text{IF}_{1, -\text{if}_{1, \tilde{\psi}_k(\theta), i_1}, i_2}(\theta) = - \frac{\partial \mathcal{H}_{i_1} \left( \tilde{b}_{k,\theta}(X_{i_1}), \tilde{p}_{k,\theta}(X_{i_1}) \right)}{\partial \bar{\zeta}_{b,k}^\top} \text{IF}_{1, \tilde{\zeta}_{b,k}, i_2}(\theta) - \frac{\partial \mathcal{H}_{i_1} \left( \tilde{b}_{k,\theta}(X_{i_1}), \tilde{p}_{k,\theta}(X_{i_1}) \right)}{\partial \bar{\zeta}_{p,k}^\top} \text{IF}_{1, \tilde{\zeta}_{p,k}, i_2}(\theta)$$

where

$$\begin{cases} \frac{\partial \mathcal{H}_{i_1}(\tilde{b}_{k,\theta}(X_{i_1}), \tilde{p}_{k,\theta}(X_{i_1}))}{\partial \bar{\zeta}_{b,k}} = S_{bp,i_1} \tilde{p}_{k,\theta}(X_{i_1}) \bar{z}_k(X_{i_1}) + m_1(O_{i_1}, \bar{z}_k), \\ \frac{\partial \mathcal{H}_{i_1}(\tilde{b}_{k,\theta}(X_{i_1}), \tilde{p}_{k,\theta}(X_{i_1}))}{\partial \bar{\zeta}_{p,k}} = S_{bp,i_1} \tilde{b}_{k,\theta}(X_{i_1}) \bar{z}_k(X_{i_1}) + m_2(O_{i_1}, \bar{z}_k) \end{cases}$$

and by Lemma A.1

$$\begin{cases} \text{IF}_{1,\tilde{\zeta}_{b,k},i_2}(\theta) = -\Sigma_k^{-1} \left( S_{bp,i_2} \tilde{b}_{k,\theta}(X_{i_2}) \bar{z}_k(X_{i_2}) + m_2(O_{i_2}, \bar{z}_k) \right), \\ \text{IF}_{1,\tilde{\zeta}_{p,k},i_2}(\theta) = -\Sigma_k^{-1} \left( S_{bp,i_2} \tilde{p}_{k,\theta}(X_{i_2}) \bar{z}_k(X_{i_2}) + m_1(O_{i_2}, \bar{z}_k) \right). \end{cases}$$

Then

$$\begin{aligned} & \text{IF}_{1,-\text{if}_{1,\tilde{\psi}_k(\theta),i_1},i_2}(\theta) \\ &= [S_{bp} \tilde{p}_{k,\theta}(X) \bar{z}_k(X) + m_1(O, \bar{z}_k)]_{i_1}^\top \Sigma_k^{-1} \left[ S_{bp} \tilde{b}_{k,\theta}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) \right]_{i_2} \\ &+ \left[ S_{bp} \tilde{b}_{k,\theta}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) \right]_{i_1}^\top \Sigma_k^{-1} [S_{bp} \tilde{p}_{k,\theta}(X) \bar{z}_k(X) + m_1(O, \bar{z}_k)]_{i_2}. \end{aligned}$$

Then by Robins et al. (2008, part 5.c of Theorem 2.3),

$$\mathbb{IF}_{22,k}(\Sigma_k^{-1}) = \frac{1}{2} \Pi_\theta \left[ \mathbb{V}_2 \left[ \text{IF}_{1,-\text{if}_{1,\tilde{\psi}_k(\theta),i_1},i_2}(\theta) \right] |\mathcal{U}_1^{\perp,2}(\theta) \right].$$

Here  $\mathcal{U}_m(\theta)$  denotes the Hilbert space of all  $m$ -th order U-statistics with mean zero and finite variance with inner product defined by covariances w.r.t. the  $n$ -fold product measure  $P_\theta(\cdot)^n$  and  $\mathcal{U}_m^{\perp,m+1}(\theta)$  denotes the orthocomplement of  $\mathcal{U}_m(\theta)$  in  $\mathcal{U}_{m+1}(\theta)$ . The notation  $\mathbb{V}_m[U_{i_1,\dots,i_m}]$  maps an  $m$ -th order U-statistic kernel  $U_{i_1,\dots,i_m}$  to an  $m$ -th order U-statistic:

$$\mathbb{V}_m[U_{i_1,\dots,i_m}] = \frac{(n-m)!}{n!} \sum_{1 \leq i_1 \neq \dots \neq i_m \leq n} U_{i_1,\dots,i_m}.$$

For the 2nd order U-statistic kernel  $\text{IF}_{1,\text{if}_{1,\tilde{\psi}_k(\theta),i_1},i_2}(\theta)$ , we thus have

$$\begin{aligned} & \mathbb{V}_2 \left[ \text{IF}_{1,-\text{if}_{1,\tilde{\psi}_k(\theta),i_1},i_2}(\theta) \right] \\ &= \mathbb{V}_2 \left[ \begin{aligned} & [S_{bp} \tilde{p}_{k,\theta}(X) \bar{z}_k(X) + m_1(O, \bar{z}_k)]_{i_1}^\top \Sigma_k^{-1} \left[ S_{bp} \tilde{b}_{k,\theta}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) \right]_{i_2} \\ & + \left[ S_{bp} \tilde{b}_{k,\theta}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) \right]_{i_1}^\top \Sigma_k^{-1} [S_{bp} \tilde{p}_{k,\theta}(X) \bar{z}_k(X) + m_1(O, \bar{z}_k)]_{i_2} \end{aligned} \right] \\ &= \mathbb{V}_2 \left[ 2 [S_{bp} \tilde{p}_{k,\theta}(X) \bar{z}_k(X) + m_1(O, \bar{z}_k)]_{i_1}^\top \Sigma_k^{-1} \left[ S_{bp} \tilde{b}_{k,\theta}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) \right]_{i_2} \right] \end{aligned}$$

where the last equality is due to the symmetrization of  $\mathbb{V}_2[\cdot]$ . In addition, we have

$$\Pi_\theta \left[ \mathbb{V}_2 \left[ \text{IF}_{1,-\text{if}_{1,\tilde{\psi}_k(\theta),i_1},i_2}(\theta) \right] |\mathcal{U}_1(\theta) \right] = 0$$

since

$$\begin{cases} \mathbb{E}_\theta [S_{bp} \tilde{p}_{k,\theta}(X) \bar{z}_k(X) + m_1(O, \bar{z}_k)] = 0, \\ \mathbb{E}_\theta [S_{bp} \tilde{b}_{k,\theta}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k)] = 0. \end{cases}$$

Thus

$$\begin{aligned} \mathbb{IF}_{22,k}(\Sigma_k^{-1}) &= \frac{1}{2} \mathbb{V}_2 \left[ \mathbf{IF}_{1,-\text{if}_{1,\tilde{\psi}_k(\theta),i_1},i_2}(\theta) \right] \\ &= \frac{1}{n(n-1)} \sum_{1 \leq i_1 \neq i_2 \leq n} [S_{bp} \tilde{p}_{k,\theta}(X) \bar{z}_k(X) + m_1(O, \bar{z}_k)]_{i_1}^\top \Sigma_k^{-1} \left[ S_{bp} \tilde{b}_{k,\theta}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) \right]_{i_2}. \end{aligned}$$

Hence formally we have

$$\mathbf{IF}_{22,k,\bar{i}_2}(\Sigma_k^{-1}) = [S_{bp} \tilde{p}_{k,\theta}(X) \bar{z}_k(X) + m_1(O, \bar{z}_k)]_{i_1}^\top \Sigma_k^{-1} \left[ S_{bp} \tilde{b}_{k,\theta}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) \right]_{i_2}$$

and

$$\widehat{\mathbf{IF}}_{22,k,\bar{i}_2}(\Sigma_k^{-1}) = [S_{bp} \hat{p}(X) \bar{z}_k(X) + m_1(O, \bar{z}_k)]_{i_1}^\top \Sigma_k^{-1} \left[ S_{bp} \hat{b}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) \right]_{i_2}.$$

We now complete the proof by induction. We assume that the form for  $\widehat{\mathbf{IF}}_{jj,k,\bar{i}_j}(\Sigma_k^{-1})$  is true and thus by the induction hypothesis

$$\begin{aligned} \mathbf{IF}_{jj,k,\bar{i}_j}(\Sigma_k^{-1}) &= (-1)^j [S_{bp} \tilde{p}_{k,\theta}(X) \bar{z}_k(X) + m_1(O, \bar{z}_k)]_{i_1}^\top \prod_{s=3}^j \Sigma_k^{-1} \left( [S_{bp} \bar{z}_k(X) \bar{z}_k(X)^\top]_{i_s} - \Sigma_k \right) \\ &\quad \times \Sigma_k^{-1} \left[ S_{bp} \tilde{b}_{k,\theta}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) \right]_{i_2} \end{aligned}$$

and

$$\mathbb{IF}_{jj,k}(\Sigma_k^{-1}) = \mathbb{V}_j \left[ \mathbf{IF}_{jj,k,\bar{i}_j}(\Sigma_k^{-1}) \right].$$

Then

$$\begin{aligned} &\mathbb{V}_{j+1} \left[ \mathbf{IF}_{(j+1),(j+1),k,\bar{i}_{j+1}}(\Sigma_k^{-1}) \right] \\ &= \frac{1}{j+1} \mathbb{V}_{j+1} \left[ \Pi_\theta \left[ \mathbf{IF}_{1,\text{if}_{jj,k,\bar{i}_j},i_{j+1}}(\theta) | \mathcal{U}_j^{\perp,j+1}(\theta) \right] \right]. \end{aligned}$$

Then the derivatives w.r.t.  $\theta$ 's in  $\tilde{\zeta}_{b,k}(\theta), \tilde{\zeta}_{p,k}(\theta)$  and  $(j-1)$  terms of  $\Sigma_k^{-1}$  will be contributing terms to  $\mathbb{V}_{j+1} \left[ \mathbf{IF}_{(j+1),(j+1),k,\bar{i}_{j+1}}(\Sigma_k^{-1}) \right]$ . But differentiating w.r.t.  $(j-2)$  terms of  $\Sigma_k$  will not contribute terms to  $\mathbb{V}_{j+1} \left[ \mathbf{IF}_{(j+1),(j+1),k,\bar{i}_{j+1}}(\Sigma_k^{-1}) \right]$  as the contribution of these  $(j-2)$  terms to  $\mathbf{IF}_{1,\text{if}_{jj,k,\bar{i}_j},i_{j+1}}(\theta)$  is a function of  $j$  units' data and is thus an element of  $\mathcal{U}_j(\theta)$ . Then as in the proof of Lemma A.1,

$$\mathbf{IF}_{1,\Sigma_k^{-1}}(\theta) = -\Sigma_k^{-1} (S_{bp} \bar{z}_k(X) \bar{z}_k(X)^\top - \Sigma_k) \Sigma_k^{-1}$$

and upon permuting the indices, the contribution of each of these  $(j - 1)$  terms to  $\text{IF}_{1, \text{if}_{jj,k,\bar{i}_j}, i_{j+1}}(\theta)$  is

$$\begin{aligned}
 & - (-1)^j [S_{bp} \tilde{p}_{k,\theta}(X) \bar{z}_k(X) + m_1(O, \bar{z}_k)]_{i_1}^\top \\
 & \times \sum_{s=3}^{j+1} \Sigma_k^{-1} \left\{ (S_{bp} \bar{z}_k(X) \bar{z}_k(X)^\top)_{i_s} - \Sigma_k \right\} \\
 & \times \Sigma_k^{-1} [S_{bp} \tilde{b}_{k,\theta}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k)]_{i_2}
 \end{aligned}
 \tag{A.5}$$

which is degenerate and thus orthogonal to  $\mathcal{U}_j(\theta)$ . Notice that equation (A.5) is exactly the kernel of  $\text{IF}_{(j+1),(j+1),k}(\Sigma_k^{-1})$ .

Then differentiating w.r.t.  $\theta$ 's in  $\tilde{p}_{k,\theta}$  and  $\tilde{b}_{k,\theta}$ , we further obtain terms of the following form contributing to  $\text{IF}_{1, \text{if}_{jj,k,\bar{i}_j}, i_{j+1}}(\theta)$ :

$$\begin{aligned}
 & (-1)^j \text{IF}_{1, \tilde{\zeta}_{p,k}, i_{j+1}}(\theta)^\top [S_{bp} \bar{z}_k(X) \bar{z}_k(X)^\top]_{i_1} \left[ \sum_{s=3}^j \Sigma_k^{-1} \left\{ (S_{bp} \bar{z}_k(X) \bar{z}_k(X)^\top)_{i_s} - \Sigma_k \right\} \right] \\
 & \times \Sigma_k^{-1} [S_{bp} \tilde{b}_{k,\theta}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k)]_{i_2} \\
 & + (-1)^j [S_{bp} \tilde{p}_{k,\theta}(X) \bar{z}_k(X) + m_1(O, \bar{z}_k)]_{i_1}^\top \left[ \sum_{s=3}^j \Sigma_k^{-1} \left\{ (S_{bp} \bar{z}_k(X) \bar{z}_k(X)^\top)_{i_s} - \Sigma_k \right\} \right] \\
 & \times \Sigma_k^{-1} [S_{bp} \bar{z}_k(X) \bar{z}_k(X)^\top]_{i_2} \text{IF}_{1, \tilde{\zeta}_{b,k}, i_{j+1}}(\theta).
 \end{aligned}
 \tag{A.6}$$

Again applying Lemma A.1 and projecting onto  $\mathcal{U}_{j,\theta}^{\perp, j+1}$ , we obtain that

$$\Pi_\theta \left[ \text{equation (A.6)} | \mathcal{U}_{j,\theta}^{\perp, j+1} \right] \equiv 2 \times \text{equation (A.5)}.$$

Recall that there are  $(j - 1)$  terms of equation (A.5) in total contributing to  $\text{IF}_{1, \text{if}_{jj,k,\bar{i}_j}, i_{j+1}}(\theta)$ , thus eventually there are  $(j + 1)$  terms of equation (A.5) in total contributing to  $\mathbb{V}_{j+1} \left[ \Pi_\theta \left[ \text{IF}_{1, \text{if}_{jj,k,\bar{i}_j}, i_{j+1}}(\theta) | \mathcal{U}_j^{\perp, j+1}(\theta) \right] \right]$ . Further recall that

$$\begin{aligned}
 & \mathbb{V}_{j+1} [\text{IF}_{(j+1),(j+1),k,\bar{i}_{j+1}}(\Sigma_k^{-1})] \\
 & = \frac{1}{j+1} \mathbb{V}_{j+1} \left[ \Pi_\theta \left[ \text{IF}_{1, \text{if}_{jj,k,\bar{i}_j}, i_{j+1}}(\theta) | \mathcal{U}_j^{\perp, j+1}(\theta) \right] \right].
 \end{aligned}$$

Therefore we have

$$\begin{aligned}
 \text{IF}_{(j+1),(j+1),k}(\Sigma_k^{-1}) & \equiv \mathbb{V}_{j+1} [\text{IF}_{(j+1),(j+1),k,\bar{i}_{j+1}}(\Sigma_k^{-1})] \\
 & = \frac{1}{j+1} \mathbb{V}_{j+1} [(j+1) \times \text{equation (A.5)}] \\
 & = \mathbb{V}_{j+1} [\text{equation (A.5)}].
 \end{aligned}$$

We thus prove the results for general  $j$  by induction. ■

**Lemma A.1.** As defined in equation (3.1), the influence functions of  $\tilde{\zeta}_{b,k}(\theta)$  and  $\tilde{\zeta}_{p,k}(\theta)$  are

$$(A.7) \quad \begin{cases} \text{IF}_{1,\tilde{\zeta}_{b,k}}(\theta) = -\Sigma_k^{-1} \left( S_{bp} \tilde{b}_{k,\theta}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) \right), \\ \text{IF}_{1,\tilde{\zeta}_{p,k}}(\theta) = -\Sigma_k^{-1} \left( S_{bp} \tilde{p}_{k,\theta}(X) \bar{z}_k(X) + m_1(O, \bar{z}_k) \right). \end{cases}$$

*Proof.* We only need to show  $\widehat{\text{IF}}_{1,\tilde{\zeta}_{b,k}}(\theta)$  and  $\widehat{\text{IF}}_{1,\tilde{\zeta}_{p,k}}(\theta)$  will follow by symmetry. Recall that

$$\tilde{\zeta}_{b,k}(\theta) = -\left\{ \mathbb{E}_\theta \left[ S_{bp} \bar{z}_k(X) \bar{z}_k(X)^\top \right] \right\}^{-1} \mathbb{E}_\theta \left[ S_{bp} \widehat{b}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) \right].$$

Then

$$\begin{aligned} & \text{IF}_{1,\tilde{\zeta}_{b,k}}(\theta) \\ &= -\left\{ \mathbb{E}_\theta \left[ S_{bp} \bar{z}_k(X) \bar{z}_k(X)^\top \right] \right\}^{-1} \text{IF}_{1,\mathbb{E}_\theta \left[ S_{bp} \widehat{b}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) \right]}(\theta) \\ &\quad - \text{IF}_{1,\left\{ \mathbb{E}_\theta \left[ S_{bp} \bar{z}_k(X) \bar{z}_k(X)^\top \right] \right\}^{-1}(\theta) \mathbb{E}_\theta \left[ S_{bp} \widehat{b}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) \right]} \\ &= -\Sigma_k^{-1} \left( S_{bp} \widehat{b}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) - \mathbb{E}_\theta \left[ S_{bp} \widehat{b}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) \right] \right) \\ &\quad + \Sigma_k^{-1} \left( S_{bp} \bar{z}_k(X) \bar{z}_k(X)^\top - \Sigma_k \right) \Sigma_k^{-1} \mathbb{E}_\theta \left[ S_{bp} \widehat{b}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) \right] \\ &= -\Sigma_k^{-1} \left( S_{bp} \widehat{b}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) - S_{bp} \bar{z}_k(X)^\top \Sigma_k^{-1} \mathbb{E}_\theta \left[ S_{bp} \widehat{b}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) \right] \bar{z}_k(X) \right) \\ &= -\Sigma_k^{-1} \left( S_{bp} \widehat{b}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) - S_{bp} \bar{z}_k(X)^\top \tilde{\zeta}_{b,k}(\theta) \bar{z}_k(X) \right) \\ &= -\Sigma_k^{-1} \left\{ S_{bp} \left( \widehat{b}(X) - \bar{z}_k(X)^\top \tilde{\zeta}_{b,k}(\theta) \right) \bar{z}_k(X) + m_2(O, \bar{z}_k) \right\} \\ &= -\Sigma_k^{-1} \left( S_{bp} \tilde{b}_{k,\theta}(X) \bar{z}_k(X) + m_2(O, \bar{z}_k) \right) \end{aligned}$$

where the fourth equality follows from the definition of  $\tilde{\zeta}_{b,k}(\theta)$  and the sixth equality follows from the definition of  $\tilde{b}_{k,\theta}(X)$ .  $\blacksquare$

Under conditions in Theorem 2.2,  $\tilde{b}_{k,\hat{\theta}} = \widehat{b}$  and  $\tilde{p}_{k,\hat{\theta}} = \widehat{p}$  or equivalently  $\tilde{\zeta}_{b,k}(\widehat{\theta}) = \tilde{\zeta}_{p,k}(\widehat{\theta}) = 0$ . This follows from results in Theorem 2.1 with  $\theta$  evaluated at  $\widehat{\theta}$ :

$$\mathbb{E}_{\widehat{\theta}} \left[ S_{bp} \widehat{b}(X) h(X) + m_2(O, h) \right] = 0 \text{ for all } h \in \mathcal{B}, \mathbb{E}_{\widehat{\theta}} \left[ S_{bp} \widehat{p}(X) h(X) + m_1(O, h) \right] = 0 \text{ for all } h \in \mathcal{P}.$$

Thus under  $\widehat{\theta}$ , the HOIFs of  $\widehat{\psi}_k(\widehat{\theta})$  are:  $\widehat{\mathbb{IF}}_{m,k}(\Sigma_k^{-1}) = \widehat{\mathbb{IF}}_1 - \widehat{\mathbb{IF}}_{22 \rightarrow mm,k}(\Sigma_k^{-1})$ , where  $\widehat{\mathbb{IF}}_{\ell\ell \rightarrow mm,k}(\Sigma_k^{-1}) := \sum_{j=\ell}^m \widehat{\mathbb{IF}}_{jj,k}(\Sigma_k^{-1})$ ,

$$(A.8) \quad \widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1}) := \frac{1}{n(n-1)} \sum_{1 \leq i_1 \neq i_2 \leq n} \widehat{\text{IF}}_{22,k,\bar{i}_2}(\Sigma_k^{-1})$$



and for  $j > 2$ ,

$$(A.9) \quad \hat{\mathbb{I}\mathbb{F}}_{jj,k}(\Sigma_k^{-1}) := \frac{(n-j)!}{n!} \sum_{1 \leq i_1 \neq \dots \neq i_j \leq n} \hat{\mathbb{I}\mathbb{F}}_{jj,k,\bar{i}_j}(\Sigma_k^{-1}).$$

Here

$$(A.10) \quad \hat{\mathbb{I}\mathbb{F}}_{22,k,\bar{i}_2}(\Sigma_k^{-1}) \equiv \mathbb{I}\mathbb{F}_{22,\tilde{\psi}_k,\bar{i}_2}(\hat{\theta}) = \left[ \mathcal{E}_{\hat{b},m_2}(\bar{\mathbf{z}}_k)(O) \right]_{i_1}^\top \Sigma_k^{-1} [\mathcal{E}_{\hat{p},m_1}(\bar{\mathbf{z}}_k)(O)]_{i_2},$$

and

$$(A.11) \quad \begin{aligned} \hat{\mathbb{I}\mathbb{F}}_{jj,k,\bar{i}_j}(\Sigma_k^{-1}) &\equiv \mathbb{I}\mathbb{F}_{jj,\tilde{\psi}_k,\bar{i}_j}(\hat{\theta}) \\ &= (-1)^j \left[ \mathcal{E}_{\hat{b},m_2}(\bar{\mathbf{z}}_k)(O) \right]_{i_1}^\top \left\{ \prod_{s=3}^j \Sigma_k^{-1} \left( [S_{\hat{b}\hat{p}} \bar{\mathbf{z}}_k(X) \bar{\mathbf{z}}_k(X)^\top]_{i_s} - \Sigma_k \right) \right\} \Sigma_k^{-1} [\mathcal{E}_{\hat{p},m_1}(\bar{\mathbf{z}}_k)(O)]_{i_2}. \end{aligned}$$

are the kernels of second order influence function  $\hat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})$  and  $j$ -th order influence function  $\hat{\mathbb{I}\mathbb{F}}_{jj,k}(\Sigma_k^{-1})$  respectively, where  $\mathcal{E}_{\hat{b},m_2}$  and  $\mathcal{E}_{\hat{p},m_1}$  are defined in equations (3.6) and (3.7).

**A.2. General formula for the variance order of  $\hat{\mathbb{I}\mathbb{F}}_{mm,k}(\hat{\Sigma}_k^{-1})$ .** In this section we derive the following formula for the order of  $\text{var}_\theta \left[ \hat{\mathbb{I}\mathbb{F}}_{mm,k}(\hat{\Sigma}_k^{-1}) \right]$  for general  $m$ . To do this, we need the following additional notation:  $\mathbb{L}_{\theta,4,\hat{b}\hat{p},k} := \left\{ \mathbb{E}_\theta [\Pi_\theta[\Delta_{1,\hat{b}}|\bar{v}_k](X)^2 \Pi_\theta[\Delta_{1,\hat{p}}|\bar{v}_k](X)^2] \right\}^{1/4}$ . Then we have

**Lemma A.2.** When  $m \geq 3$ ,

$$(A.12) \quad \begin{aligned} \text{var}_\theta \left[ \hat{\mathbb{I}\mathbb{F}}_{mm,k}(\hat{\Sigma}_k^{-1}) \right] &\lesssim \frac{1}{n} \left\{ \begin{aligned} &\left( \frac{k}{n} \right)^{m-1} + \left( \frac{k}{n} \right)^{m-2} \sum_{f \in \{\hat{b}, \hat{p}, \hat{\Sigma}\}} \mathbb{L}_{\theta,2,f,k}^2 \\ &+ \sum_{j=3}^{m-1} \left( \frac{k}{n} \right)^{m-j} \left( \mathbb{L}_{\theta,2,\hat{\Sigma},k}^2 \right)^{2(\lfloor \frac{j}{2} \rfloor - 1)} \left\{ \sum_{f_1 \neq f_2 \in \{\hat{b}, \hat{p}, \hat{\Sigma}\}} \mathbb{L}_{\theta,2,f_1,k}^2 \mathbb{L}_{\theta,2,f_2,k}^2 + \mathbb{L}_{\theta,2,\hat{\Sigma},k}^4 \right\} \\ &+ \left( \mathbb{L}_{\theta,2,\hat{\Sigma},k}^2 \right)^{2(\lfloor \frac{m}{2} \rfloor - 1)} \left\{ \sum_{f_1 \neq f_2 \in \{\hat{b}, \hat{p}, \hat{\Sigma}\}} \mathbb{L}_{\theta,2,f_1,k}^2 \mathbb{L}_{\theta,2,f_2,k}^2 + \mathbb{L}_{\theta,4,\hat{b}\hat{p},k}^4 \right\} \end{aligned} \right\} \\ &\lesssim \frac{1}{n} \left\{ \begin{aligned} &\left( \frac{k}{n} \right)^{m-1} + \left( \frac{k}{n} \right)^{m-2} \sum_{f \in \{\hat{b}, \hat{p}, \hat{\Sigma}\}} \mathbb{L}_{\theta,2,f,k}^2 \\ &+ \sum_{j=3}^{m-1} \left( \frac{k}{n} \right)^{m-j} \left( \mathbb{L}_{\theta,2,\hat{\Sigma},k}^2 \right)^{2(\lfloor \frac{j}{2} \rfloor - 1)} \left\{ \sum_{f_1 \neq f_2 \in \{\hat{b}, \hat{p}, \hat{\Sigma}\}} \mathbb{L}_{\theta,2,f_1,k}^2 \mathbb{L}_{\theta,2,f_2,k}^2 + \mathbb{L}_{\theta,2,\hat{\Sigma},k}^4 \right\} \\ &+ \left( \mathbb{L}_{\theta,2,\hat{\Sigma},k}^2 \right)^{2(\lfloor \frac{m}{2} \rfloor - 1)} \left\{ \sum_{f_1 \neq f_2 \in \{\hat{b}, \hat{p}, \hat{\Sigma}\}} \mathbb{L}_{\theta,2,f_1,k}^2 \mathbb{L}_{\theta,2,f_2,k}^2 + \mathbb{L}_{\theta,2,\hat{b},k}^2 + \mathbb{L}_{\theta,2,\hat{p},k}^2 \right\} \end{aligned} \right\}. \end{aligned}$$

When  $m = 2$ ,

$$(A.13) \quad \text{var}_\theta \left[ \widehat{\mathbb{IF}}_{22,k}(\widehat{\Sigma}_k^{-1}) \right] \lesssim \frac{1}{n} \left\{ \frac{k}{n} + \sum_{f \in \{\widehat{b}, \widehat{p}\}} \mathbb{L}_{\theta,2,f,k}^2 \right\}.$$

Under Condition **W**,  $\mathbb{L}_{\theta,4,\widehat{b},k}^4 \lesssim \mathbb{L}_{\theta,2,\widehat{b},k}^2 + \mathbb{L}_{\theta,2,\widehat{p},k}^2$ .

*Proof.* Equation **A.13** will follow from the celebrated Hoeffding's decomposition (Hoeffding, 1948) of the U-statistics  $\widehat{\mathbb{IF}}_{mm,k}(\widehat{\Sigma}_k^{-1}) - \mathbb{E}_\theta \left[ \widehat{\mathbb{IF}}_{mm,k}(\widehat{\Sigma}_k^{-1}) \right]$  and then we simply bound each term in the decomposition. The detailed calculations can be found in the forthcoming updated version of Mukherjee et al. (2017). ■

**A.3. Bootstrapping higher order influence functions.** In this section, we will use higher order moments of multinomial distributions frequently. Two good references are Mosimann (1962); Newcomer et al. (2008). For notational convenience, we suppress the dependence on  $\Sigma_k^{-1}$  or  $\widehat{\Sigma}_k^{-1}$  in  $\widehat{\mathbb{IF}}_{mm,k}(\Sigma_k^{-1})$  and  $\widehat{\mathbb{IF}}_{mm,k}(\widehat{\Sigma}_k^{-1})$  as the proposed bootstrap resampling procedure will not depend on  $\Sigma_k^{-1}$  or the training sample. We propose the following bootstrap estimator of  $\text{var}_\theta[\widehat{\mathbb{IF}}_{22,k}]$ : choose some large integer  $M$  as the number of bootstrap resamples, and define

$$(A.14) \quad \widehat{\text{var}}[\widehat{\mathbb{IF}}_{22,k}] := \underbrace{\frac{1}{M-1} \sum_{m=1}^M \left( \widehat{\mathbb{IF}}_{22,k}^{(m)} - \frac{1}{M} \sum_{m=1}^M \widehat{\mathbb{IF}}_{22,k}^{(m)} \right)^2}_{T_1} - \underbrace{\frac{2}{M-1} \sum_{m=1}^M \left( \widehat{\mathbb{IF}}_{22,k}^{(m),c} - \frac{1}{M} \sum_{m=1}^M \widehat{\mathbb{IF}}_{22,k}^{(m),c} \right)^2}_{T_2}$$

where

$$\begin{aligned} \widehat{\mathbb{IF}}_{22,k}^{(m)} &= \frac{1}{n(n-1)} \sum_{1 \leq i_1 \neq i_2 \leq n} W_{i_1}^{(m)} W_{i_2}^{(m)} \widehat{\mathbb{F}}_{22,k,\bar{i}_2} \\ \widehat{\mathbb{IF}}_{22,k}^{(m),c} &= \frac{1}{n(n-1)} \sum_{1 \leq i_1 \neq i_2 \leq n} (W_{i_1}^{(m)} - 1)(W_{i_2}^{(m)} - 1) \widehat{\mathbb{F}}_{22,k,\bar{i}_2} \end{aligned}$$

and  $\mathbf{W}^{(m)} = (W_1^{(m)}, \dots, W_n^{(m)}) \stackrel{i.i.d.}{\sim} \text{Multinom}(n, \underbrace{n^{-1}, \dots, n^{-1}}_{(n-1)'s})$  for  $m = 1, \dots, M$  are  $M$  multinomial weights independent of the observed data  $\{O_i, i = 1, \dots, N\}$ .

**Comment A.1.** We now explain why  $\widehat{\text{var}}[\widehat{\mathbb{IF}}_{22,k}]$  is proposed as an estimator of  $\text{var}_\theta[\widehat{\mathbb{IF}}_{22,k}]$ . Conditional on the observed data (including both the training and estimation samples), the expectation over the random multinomial weights  $\mathbf{W}$  of the term  $T_1$  is:

$$\begin{aligned} & \mathbb{E}_{\mathbf{W}}[T_1 | \{O_i\}_{i=1}^N] \\ &= \mathbb{E}_{\mathbf{W}} \left[ \left( \frac{1}{n(n-1)} \sum_{1 \leq i_1 \neq i_2 \leq n} W_{i_1}^{(m)} W_{i_2}^{(m)} \widehat{\mathbb{F}}_{22,k,i_1,i_2} \right)^2 | \{O_i\}_{i=1}^N \right] \end{aligned}$$

$$\begin{aligned}
& - \left( \mathbf{E}_{\mathbf{W}} \left[ \frac{1}{n(n-1)} \sum_{1 \leq i_1 \neq i_2 \leq n} W_{i_1}^{(m)} W_{i_2}^{(m)} \widehat{\mathbf{F}}_{22,k,i_1,i_2} | \{O_i\}_{i=1}^N \right] \right)^2 \\
& = \frac{1}{n^2(n-1)^2} \sum_{1 \leq i_1 \neq i_2 \leq n} \{ \widehat{\mathbf{F}}_{22,k,i_1,i_2} \widehat{\mathbf{F}}_{22,k,i_1,i_2} + \widehat{\mathbf{F}}_{22,k,i_1,i_2} \widehat{\mathbf{F}}_{22,k,i_2,i_1} \} \{ \mathbf{E}_{\mathbf{W}}[W_{i_1}^2 W_{i_2}^2] - (\mathbf{E}_{\mathbf{W}}[W_{i_1} W_{i_2}])^2 \} \\
& + \frac{1}{n^2(n-1)^2} \sum_{1 \leq i_1 \neq i_2 \neq i_3 \leq n} \left\{ \begin{aligned} & \widehat{\mathbf{F}}_{22,k,i_1,i_2} \widehat{\mathbf{F}}_{22,k,i_1,i_3} \{ \mathbf{E}_{\mathbf{W}}[W_{i_1}^2 W_{i_2} W_{i_3}] - \mathbf{E}_{\mathbf{W}}[W_{i_1} W_{i_2}] \mathbf{E}_{\mathbf{W}}[W_{i_1} W_{i_3}] \} \\ & + \widehat{\mathbf{F}}_{22,k,i_1,i_2} \widehat{\mathbf{F}}_{22,k,i_3,i_1} \{ \mathbf{E}_{\mathbf{W}}[W_{i_1}^2 W_{i_2} W_{i_3}] - \mathbf{E}_{\mathbf{W}}[W_{i_1} W_{i_2}] \mathbf{E}_{\mathbf{W}}[W_{i_3} W_{i_1}] \} \\ & + \widehat{\mathbf{F}}_{22,k,i_1,i_2} \widehat{\mathbf{F}}_{22,k,i_2,i_3} \{ \mathbf{E}_{\mathbf{W}}[W_{i_1} W_{i_2}^2 W_{i_3}] - \mathbf{E}_{\mathbf{W}}[W_{i_1} W_{i_2}] \mathbf{E}_{\mathbf{W}}[W_{i_2} W_{i_3}] \} \\ & + \widehat{\mathbf{F}}_{22,k,i_1,i_2} \widehat{\mathbf{F}}_{22,k,i_3,i_2} \{ \mathbf{E}_{\mathbf{W}}[W_{i_1} W_{i_2}^2 W_{i_3}] - \mathbf{E}_{\mathbf{W}}[W_{i_1} W_{i_2}] \mathbf{E}_{\mathbf{W}}[W_{i_3} W_{i_2}] \} \end{aligned} \right\} \\
& + \frac{1}{n^2(n-1)^2} \sum_{1 \leq i_1 \neq i_2 \neq i_3 \neq i_4 \leq n} \widehat{\mathbf{F}}_{22,k,i_1,i_2} \widehat{\mathbf{F}}_{22,k,i_3,i_4} \{ \mathbf{E}_{\mathbf{W}}[W_{i_1} W_{i_2} W_{i_3} W_{i_4}] - \mathbf{E}_{\mathbf{W}}[W_{i_1} W_{i_2}] \mathbf{E}_{\mathbf{W}}[W_{i_3} W_{i_4}] \}.
\end{aligned}$$

Plugging in the following higher order moments of Multinom( $n, \underbrace{n^{-1}, \dots, n^{-1}}_{(n-1)'s}$ ) ([Newcomer et al., 2008](#)):

$$\begin{aligned}
\mathbf{E}_{\mathbf{W}}[W_1 W_2] &= \frac{n-1}{n}, \mathbf{E}_{\mathbf{W}}[W_1 W_2 W_3] = \frac{(n-1)(n-2)}{n^2}, \mathbf{E}_{\mathbf{W}}[W_1 W_2 W_3 W_4] = \frac{(n-1)(n-2)(n-3)}{n^3}, \\
\mathbf{E}_{\mathbf{W}}[W_1^2 W_2] &= \frac{2(n-1)^2}{n^2}, \mathbf{E}_{\mathbf{W}}[W_1^2] = \frac{2n-1}{n}, \\
\mathbf{E}_{\mathbf{W}}[W_1^2 W_2 W_3] &= \frac{(n-1)(n-2)(n-3) + n(n-1)(n-2)}{n^3} \\
&= \frac{(n-1)(2n^2 - 7n + 6)}{n^3}, \\
\mathbf{E}_{\mathbf{W}}[W_1^2 W_2^2] &= \frac{(n-1)(n-2)(n-3) + 2n(n-1)(n-2) + n^2(n-1)}{n^3} \\
&= \frac{(n-1)(4n^2 - 9n + 6)}{n^3},
\end{aligned}$$

we have

$$\begin{aligned}
& \mathbf{E}_{\mathbf{W}}[T_1 | \{O_i\}_{i=1}^N] \\
\text{(A.15)} \quad & = \frac{1}{n^2(n-1)^2} \sum_{1 \leq i_1 \neq i_2 \leq n} \{ \widehat{\mathbf{F}}_{22,k,i_1,i_2} \widehat{\mathbf{F}}_{22,k,i_1,i_2} + \widehat{\mathbf{F}}_{22,k,i_1,i_2} \widehat{\mathbf{F}}_{22,k,i_2,i_1} \} \left( 3 - \frac{11}{n} + \frac{14}{n^2} - \frac{6}{n^3} \right) \\
& + \frac{1}{n^2(n-1)^2} \sum_{1 \leq i_1 \neq i_2 \neq i_3 \leq n} \left\{ \begin{aligned} & \widehat{\mathbf{F}}_{22,k,i_1,i_2} \widehat{\mathbf{F}}_{22,k,i_1,i_3} + \widehat{\mathbf{F}}_{22,k,i_1,i_2} \widehat{\mathbf{F}}_{22,k,i_3,i_1} \\ & + \widehat{\mathbf{F}}_{22,k,i_1,i_2} \widehat{\mathbf{F}}_{22,k,i_2,i_3} + \widehat{\mathbf{F}}_{22,k,i_1,i_2} \widehat{\mathbf{F}}_{22,k,i_3,i_2} \end{aligned} \right\} \left( 1 - \frac{7}{n} + \frac{12}{n^2} - \frac{6}{n^3} \right) \\
& + \frac{1}{n^2(n-1)^2} \sum_{1 \leq i_1 \neq i_2 \neq i_3 \neq i_4 \leq n} \widehat{\mathbf{F}}_{22,k,i_1,i_2} \widehat{\mathbf{F}}_{22,k,i_3,i_4} \left( -\frac{4}{n} + \frac{10}{n^2} - \frac{6}{n^3} \right)
\end{aligned}$$

The term  $T_2$  is the bootstrap variance estimator for second order degenerate U-statistics and it has been proposed previously in [Arcones and Gine \(1992\)](#); [Huskova and Janssen \(1993\)](#). Similarly, conditional on the observed data (including both the training and estimation samples), the expectation over the random multinomial weights  $\mathbf{W}$  of the term  $T_2$  is:

$$\begin{aligned}
& \mathbf{E}_{\mathbf{W}}[T_2 | \{O_i\}_{i=1}^N] \\
&= \frac{2}{n^2(n-1)^2} \sum_{1 \leq i_1 \neq i_2 \leq n} \{ \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_1,i_2} + \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_2,i_1} \} \{ \mathbf{E}_{\mathbf{W}}[(W_{i_1}-1)^2(W_{i_2}-1)^2] - (\mathbf{E}_{\mathbf{W}}[(W_{i_1}-1)(W_{i_2}-1)])^2 \} \\
&+ \frac{2}{n^2(n-1)^2} \sum_{1 \leq i_1 \neq i_2 \neq i_3 \leq n} \left\{ \begin{aligned} & \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_1,i_3} \{ \mathbf{E}_{\mathbf{W}}[(W_{i_1}-1)^2(W_{i_2}-1)(W_{i_3}-1)] - \mathbf{E}_{\mathbf{W}}[(W_{i_1}-1)(W_{i_2}-1)]\mathbf{E}_{\mathbf{W}}[(W_{i_1}-1)(W_{i_3}-1)] \} \\ & + \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_3,i_1} \{ \mathbf{E}_{\mathbf{W}}[(W_{i_1}-1)^2(W_{i_2}-1)(W_{i_3}-1)] - \mathbf{E}_{\mathbf{W}}[(W_{i_1}-1)(W_{i_2}-1)]\mathbf{E}_{\mathbf{W}}[(W_{i_3}-1)(W_{i_1}-1)] \} \\ & + \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_2,i_3} \{ \mathbf{E}_{\mathbf{W}}[(W_{i_1}-1)(W_{i_2}-1)^2(W_{i_3}-1)] - \mathbf{E}_{\mathbf{W}}[(W_{i_1}-1)(W_{i_2}-1)]\mathbf{E}_{\mathbf{W}}[(W_{i_2}-1)(W_{i_3}-1)] \} \\ & + \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_3,i_2} \{ \mathbf{E}_{\mathbf{W}}[(W_{i_1}-1)(W_{i_2}-1)^2(W_{i_3}-1)] - \mathbf{E}_{\mathbf{W}}[(W_{i_1}-1)(W_{i_2}-1)]\mathbf{E}_{\mathbf{W}}[(W_{i_3}-1)(W_{i_2}-1)] \} \end{aligned} \right\} \\
&+ \frac{2}{n^2(n-1)^2} \sum_{1 \leq i_1 \neq i_2 \neq i_3 \neq i_4 \leq n} \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_3,i_4} \{ \mathbf{E}_{\mathbf{W}}[(W_{i_1}-1)(W_{i_2}-1)(W_{i_3}-1)(W_{i_4}-1)] - \mathbf{E}_{\mathbf{W}}[(W_{i_1}-1)(W_{i_2}-1)]\mathbf{E}_{\mathbf{W}}[(W_{i_3}-1)(W_{i_4}-1)] \} \\
&= \frac{2}{n^2(n-1)^2} \sum_{1 \leq i_1 \neq i_2 \leq n} \{ \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_1,i_2} + \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_2,i_1} \} \{ \mathbf{E}_{\mathbf{W}}[(W_1^2 - 2W_1 + 1)(W_2^2 - 2W_2 + 1)] - (\mathbf{E}_{\mathbf{W}}[(W_1 - 1)(W_2 - 1)])^2 \} \\
&+ \frac{2}{n^2(n-1)^2} \sum_{1 \leq i_1 \neq i_2 \neq i_3 \leq n} \left\{ \begin{aligned} & \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_1,i_3} + \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_3,i_1} \\ & + \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_2,i_3} + \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_3,i_2} \end{aligned} \right\} \{ \mathbf{E}_{\mathbf{W}}[(W_1 - 1)^2(W_2 - 1)(W_3 - 1)] - (\mathbf{E}_{\mathbf{W}}[(W_1 - 1)(W_2 - 1)])^2 \} \\
&+ \frac{2}{n^2(n-1)^2} \sum_{1 \leq i_1 \neq i_2 \neq i_3 \neq i_4 \leq n} \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_3,i_4} \{ \mathbf{E}_{\mathbf{W}}[(W_1 - 1)(W_2 - 1)(W_3 - 1)(W_4 - 1)] - (\mathbf{E}_{\mathbf{W}}[(W_1 - 1)(W_2 - 1)])^2 \} \\
&= \frac{2}{n^2(n-1)^2} \sum_{1 \leq i_1 \neq i_2 \leq n} \{ \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_1,i_2} + \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_2,i_1} \} \left( 1 - \frac{3}{n} + \frac{7}{n^2} - \frac{6}{n^3} - \frac{1}{n^2} \right) \\
&+ \frac{2}{n^2(n-1)^2} \sum_{1 \leq i_1 \neq i_2 \neq i_3 \leq n} \left\{ \begin{aligned} & \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_1,i_3} + \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_3,i_1} \\ & + \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_2,i_3} + \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_3,i_2} \end{aligned} \right\} \left( -\frac{1}{n} + \frac{5}{n^2} - \frac{6}{n^3} - \frac{1}{n^2} \right) \\
&+ \frac{2}{n^2(n-1)^2} \sum_{1 \leq i_1 \neq i_2 \neq i_3 \neq i_4 \leq n} \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_3,i_4} \left( \frac{3}{n^2} - \frac{6}{n^3} - \frac{1}{n^2} \right).
\end{aligned}$$

Hence  $T_2$  serves as a correction term of the inflated factor 3 appeared in equation (A.15).

Combining the above computations, we have

$$\begin{aligned}
& \mathbf{E}_{\mathbf{W}}[\widehat{\text{var}}[\widehat{\mathbb{IF}}_{22,k}] | \{O_i\}_{i=1}^N] = \mathbf{E}_{\mathbf{W}}[T_1 | \{O_i\}_{i=1}^N] - \mathbf{E}_{\mathbf{W}}[T_2 | \{O_i\}_{i=1}^N] \\
&= \frac{1}{n^2(n-1)^2} \sum_{1 \leq i_1 \neq i_2 \leq n} \{ \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_1,i_2} + \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_2,i_1} \} \left( 1 - \frac{5}{n} + \frac{2}{n^2} + \frac{6}{n^3} \right)
\end{aligned}$$

$$\begin{aligned}
& + \frac{1}{n^2(n-1)^2} \sum_{1 \leq i_1 \neq i_2 \neq i_3 \leq n} \left\{ \begin{aligned} & \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_1,i_3} + \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_3,i_1} \\ & + \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_2,i_3} + \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_3,i_2} \end{aligned} \right\} \left( 1 - \frac{5}{n} + \frac{6}{n^3} \right) \\
& - \frac{4}{n} \frac{1}{n^2(n-1)^2} \sum_{1 \leq i_1 \neq i_2 \neq i_3 \neq i_4 \leq n} \widehat{\mathbb{F}}_{22,k,i_1,i_2} \widehat{\mathbb{F}}_{22,k,i_3,i_4} \left( 1 - \frac{1.5}{n} - \frac{1.5}{n^2} \right).
\end{aligned}$$

Finally, it is not difficult to see that

$$\mathbb{E}_\theta[\mathbb{E}_W[\widehat{\text{var}}[\widehat{\mathbb{IF}}_{22,k}] | \{O_i\}_{i=1}^N]] = \text{var}_\theta[\widehat{\mathbb{IF}}_{22,k}] (1 + O(n^{-1})).$$

**Proposition A.1.** Under Condition **W**, as  $n \rightarrow \infty$  and  $M \rightarrow \infty$ ,

$$\frac{\widetilde{\text{var}}[\widehat{\mathbb{IF}}_{22,k}]}{\text{var}_\theta[\widehat{\mathbb{IF}}_{22,k}]} = 1 + o_{\mathbb{P}_\theta}(1).$$

*Proof.* Let the  $\widetilde{\text{var}}[\widehat{\mathbb{IF}}_{22,k}]$  be the limit (in  $\mathbb{P}_W$ -probability) of  $\widehat{\text{var}}[\widehat{\mathbb{IF}}_{22,k}]$  as  $M \rightarrow \infty$  given all the observed data. Using the moments of multinomial distributions (Newcomer et al., 2008), together with the explanation in Comment **A.1**, it is easy to show that  $\mathbb{E}_\theta \left[ \frac{\widetilde{\text{var}}[\widehat{\mathbb{IF}}_{22,k}]}{\text{var}_\theta[\widehat{\mathbb{IF}}_{22,k}]} \right] = 1 + o(1)$  and  $\text{var}_\theta \left[ \frac{\widetilde{\text{var}}[\widehat{\mathbb{IF}}_{22,k}]}{\text{var}_\theta[\widehat{\mathbb{IF}}_{22,k}]} \right] = o(1)$ . Combining the above two claims, we have  $\frac{\widetilde{\text{var}}[\widehat{\mathbb{IF}}_{22,k}]}{\text{var}_\theta[\widehat{\mathbb{IF}}_{22,k}]} = 1 + o_{\mathbb{P}_\theta}(1)$ .  $\blacksquare$

Following computations similar in spirit to but more tedious than the computations above, we propose the following estimator of  $\text{var}_\theta[\widehat{\mathbb{IF}}_{33,k}]$ :

$$\begin{aligned}
\widehat{\text{var}}[\widehat{\mathbb{IF}}_{33,k}] &:= \underbrace{\frac{1}{M-1} \sum_{m=1}^M \left( \widehat{\mathbb{IF}}_{33,k}^{(m)} - \frac{1}{M} \sum_{m=1}^M \widehat{\mathbb{IF}}_{33,k}^{(m)} \right)^2}_{S_1} \\
&\quad - \underbrace{\frac{1}{M-1} \sum_{m=1}^M \left( \widehat{\mathbb{IF}}_{33,k}^{(m),c(b,p)} - \frac{1}{M} \sum_{m=1}^M \widehat{\mathbb{IF}}_{33,k}^{(m),c(b,p)} \right)^2}_{S_2} \\
&\quad - \underbrace{\frac{1}{M-1} \sum_{m=1}^M \left( \widehat{\mathbb{IF}}_{33,k}^{(m),c(b,\Sigma)} - \frac{1}{M} \sum_{m=1}^M \widehat{\mathbb{IF}}_{33,k}^{(m),c(b,\Sigma)} \right)^2}_{S_3} \\
&\quad - \underbrace{\frac{1}{M-1} \sum_{m=1}^M \left( \widehat{\mathbb{IF}}_{33,k}^{(m),c(p,\Sigma)} - \frac{1}{M} \sum_{m=1}^M \widehat{\mathbb{IF}}_{33,k}^{(m),c(p,\Sigma)} \right)^2}_{S_4}
\end{aligned}
\tag{A.16}$$

where

$$\begin{aligned}\widehat{\mathbb{IF}}_{33,k}^{(m)} &= \frac{1}{n(n-1)(n-2)} \sum_{1 \leq i_1 \neq i_2 \neq i_3 \leq n} W_{i_1}^{(m)} W_{i_2}^{(m)} W_{i_3}^{(m)} \widehat{\mathbb{IF}}_{33,k,\bar{i}_3} \\ \widehat{\mathbb{IF}}_{33,k}^{(m),c(b,p)} &= \frac{1}{n(n-1)(n-2)} \sum_{1 \leq i_1 \neq i_2 \neq i_3 \leq n} (W_{i_1}^{(m)} - 1) W_{i_2}^{(m)} (W_{i_3}^{(m)} - 1) \widehat{\mathbb{IF}}_{33,k,\bar{i}_3} \\ \widehat{\mathbb{IF}}_{33,k}^{(m),c(b,\Sigma)} &= \frac{1}{n(n-1)(n-2)} \sum_{1 \leq i_1 \neq i_2 \neq i_3 \leq n} (W_{i_1}^{(m)} - 1) (W_{i_2}^{(m)} - 1) W_{i_3}^{(m)} \widehat{\mathbb{IF}}_{33,k,\bar{i}_3} \\ \widehat{\mathbb{IF}}_{33,k}^{(m),c(p,\Sigma)} &= \frac{1}{n(n-1)(n-2)} \sum_{1 \leq i_1 \neq i_2 \neq i_3 \leq n} W_{i_1}^{(m)} (W_{i_2}^{(m)} - 1) (W_{i_3}^{(m)} - 1) \widehat{\mathbb{IF}}_{33,k,\bar{i}_3}.\end{aligned}$$

We then have:

**Proposition A.2.** Under Condition **W**, as  $n \rightarrow \infty$  and  $M \rightarrow \infty$ ,

$$\frac{\widehat{\text{var}}[\widehat{\mathbb{IF}}_{33,k}]}{\text{var}_{\theta}[\widehat{\mathbb{IF}}_{33,k}]} = 1 + o_{\mathbf{P}_{\theta}}(1).$$

Finally, we consider estimating  $\text{var}_{\theta}[\widehat{\mathbb{IF}}_{22 \rightarrow 33,k}]$  through nonparametric bootstrap. It is not difficult to see that the following estimator has the desired property  $\frac{\widehat{\text{var}}[\widehat{\mathbb{IF}}_{22 \rightarrow 33,k}]}{\text{var}_{\theta}[\widehat{\mathbb{IF}}_{22 \rightarrow 33,k}]} = 1 + o_{\mathbf{P}_{\theta}}(1)$  as  $M \rightarrow \infty$ :

$$\widehat{\text{var}}[\widehat{\mathbb{IF}}_{22 \rightarrow 33,k}] = \widehat{\text{var}}[\widehat{\mathbb{IF}}_{22,k}] + \widehat{\text{var}}[\widehat{\mathbb{IF}}_{33,k}] + 2\widehat{\text{cov}}[\widehat{\mathbb{IF}}_{22,k}, \widehat{\mathbb{IF}}_{33,k}]$$

where

$$\begin{aligned}\widehat{\text{cov}}[\widehat{\mathbb{IF}}_{22,k}, \widehat{\mathbb{IF}}_{33,k}] &= \frac{1}{M-1} \sum_{m=1}^M \left( \widehat{\mathbb{IF}}_{22,k}^{(m)} - \frac{1}{M} \sum_{m'=1}^M \widehat{\mathbb{IF}}_{22,k}^{(m')} \right) \left( \widehat{\mathbb{IF}}_{33,k}^{(m)} - \frac{1}{M} \sum_{m'=1}^M \widehat{\mathbb{IF}}_{33,k}^{(m')} \right) \\ &\quad - \frac{2}{M-1} \sum_{m=1}^{M-1} \left( \widehat{\mathbb{IF}}_{22,k}^{(m),c} - \frac{1}{M} \sum_{m'=1}^M \widehat{\mathbb{IF}}_{22,k}^{(m'),c} \right) \left( \widehat{\mathbb{IF}}_{33,k}^{(m),c(b,p)} - \frac{1}{M} \sum_{m'=1}^M \widehat{\mathbb{IF}}_{33,k}^{(m'),c(b,p)} \right).\end{aligned}$$

It is not difficult to generalize the above arguments to construct bootstrapped estimators of  $\text{var}_{\theta}[\widehat{\mathbb{IF}}_{mm,k}]$  for general  $m \geq 2$ . We plan to report the general construction elsewhere.

#### A.4. Proof of Theorem 3.2.

*Proof.* The rejection probability of  $\widehat{\chi}_{2,k}(\Sigma_k^{-1}; \varsigma_k, \delta)$  is computed as follows:

$$\lim_{n \rightarrow \infty} \mathbf{P}_{\theta} \left( \frac{|\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})|}{\widehat{\text{s.e.}}[\widehat{\psi}_1]} - \varsigma_k \frac{\widehat{\text{s.e.}}[\widehat{\mathbb{IF}}_{22,k}(\Sigma_k^{-1})]}{\widehat{\text{s.e.}}[\widehat{\psi}_1]} > \delta \right)$$

$$\begin{aligned}
&= \lim_{n \rightarrow \infty} \left\{ \begin{aligned} &\mathbf{P}_\theta \left( \frac{\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})}{\widehat{\mathbf{s.e.}}[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} > \varsigma_k + \delta \frac{\widehat{\mathbf{s.e.}}[\widehat{\psi}_1]}{\widehat{\mathbf{s.e.}}[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} \right) \\ &+ \mathbf{P}_\theta \left( \frac{\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})}{\widehat{\mathbf{s.e.}}[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} < -\varsigma_k - \delta \frac{\widehat{\mathbf{s.e.}}[\widehat{\psi}_1]}{\widehat{\mathbf{s.e.}}[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} \right) \end{aligned} \right\} \\
&= \lim_{n \rightarrow \infty} \mathbf{P}_\theta \left( \frac{\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1}) - \text{Bias}_{\theta,k}(\widehat{\psi}_1)}{\widehat{\mathbf{s.e.}}[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} > \varsigma_k - \frac{\text{Bias}_{\theta,k}(\widehat{\psi}_1)}{\widehat{\mathbf{s.e.}}[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} + \delta \frac{\widehat{\mathbf{s.e.}}[\widehat{\psi}_1]}{\widehat{\mathbf{s.e.}}[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} \right) \\
&\quad + \lim_{n \rightarrow \infty} \mathbf{P}_\theta \left( \frac{\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1}) - \text{Bias}_{\theta,k}(\widehat{\psi}_1)}{\widehat{\mathbf{s.e.}}[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} < -\varsigma_k - \frac{\text{Bias}_{\theta,k}(\widehat{\psi}_1)}{\widehat{\mathbf{s.e.}}[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} - \delta \frac{\widehat{\mathbf{s.e.}}[\widehat{\psi}_1]}{\widehat{\mathbf{s.e.}}[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} \right) \\
&= \lim_{n \rightarrow \infty} \mathbf{P}_\theta \left( \frac{\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1}) - \text{Bias}_{\theta,k}(\widehat{\psi}_1)}{\mathbf{s.e.}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} (1 + o_{\mathbf{P}_\theta}(1)) > \varsigma_k - (\gamma - \delta) \frac{\mathbf{s.e.}_\theta[\widehat{\psi}_1]}{\mathbf{s.e.}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} (1 + o_{\mathbf{P}_\theta}(1)) \right) \\
&\quad + \lim_{n \rightarrow \infty} \mathbf{P}_\theta \left( \frac{\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1}) - \text{Bias}_{\theta,k}(\widehat{\psi}_1)}{\mathbf{s.e.}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} (1 + o_{\mathbf{P}_\theta}(1)) < -\varsigma_k - (\gamma + \delta) \frac{\mathbf{s.e.}_\theta[\widehat{\psi}_1]}{\mathbf{s.e.}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} (1 + o_{\mathbf{P}_\theta}(1)) \right) \\
&= 1 - \Phi \left( \varsigma_k - (\gamma - \delta) \frac{\mathbf{s.e.}_\theta[\widehat{\psi}_1]}{\mathbf{s.e.}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} \right) + \Phi \left( -\varsigma_k - (\gamma + \delta) \frac{\mathbf{s.e.}_\theta[\widehat{\psi}_1]}{\mathbf{s.e.}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} \right) \\
&= 2 - \Phi \left( \varsigma_k - (\gamma - \delta) \frac{\mathbf{s.e.}_\theta[\widehat{\psi}_1]}{\mathbf{s.e.}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} \right) - \Phi \left( \varsigma_k + (\gamma + \delta) \frac{\mathbf{s.e.}_\theta[\widehat{\psi}_1]}{\mathbf{s.e.}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})]} \right).
\end{aligned}$$

■



**A.5. Proof of Proposition 4.3.** As discussed in Section 4, since  $\Sigma_k^{-1}$  is generally unknown,  $\widehat{\mathbb{I}\mathbb{F}}_{22,k}(\Sigma_k^{-1})$  needs to be replaced by  $\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})$ . We consider the following statistic used in the test  $\widehat{\chi}_{3,k}(\widehat{\Sigma}_k^{-1}; \varsigma_k, \delta)$  after standardization:

(A.17)

$$\begin{aligned}
& \frac{\widehat{\text{s.e.}}[\widehat{\psi}_1]}{\widehat{\text{s.e.}}[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]} \left( \frac{\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})}{\widehat{\text{s.e.}}[\widehat{\psi}_1]} - \delta \right) \\
&= \left( \frac{\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1}) - \text{Bias}_{\theta,k}(\widehat{\psi}_1)}{\widehat{\text{s.e.}}[\widehat{\psi}_1]} \frac{\widehat{\text{s.e.}}[\widehat{\psi}_1]}{\widehat{\text{s.e.}}[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]} \right) + \frac{\widehat{\text{s.e.}}[\widehat{\psi}_1]}{\widehat{\text{s.e.}}[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]} \left( \frac{\text{Bias}_{\theta,k}(\widehat{\psi}_1)}{\widehat{\text{s.e.}}[\widehat{\psi}_1]} - \delta \right) \\
&= \left( \frac{\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1}) - \mathbb{E}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]}{\widehat{\text{s.e.}}[\widehat{\psi}_1]} \frac{\widehat{\text{s.e.}}[\widehat{\psi}_1]}{\widehat{\text{s.e.}}[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]} \right) + \frac{\text{EB}_{\theta,3,k}(\widehat{\Sigma}_k^{-1})}{\widehat{\text{s.e.}}[\widehat{\psi}_1]} \frac{\widehat{\text{s.e.}}[\widehat{\psi}_1]}{\widehat{\text{s.e.}}[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]} \\
&\quad + \frac{\widehat{\text{s.e.}}[\widehat{\psi}_1]}{\widehat{\text{s.e.}}[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]} \left( \frac{\text{Bias}_{\theta,k}(\widehat{\psi}_1)}{\widehat{\text{s.e.}}[\widehat{\psi}_1]} - \delta \right) \\
&= \left\{ \left( \frac{\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1}) - \mathbb{E}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]}{\widehat{\text{s.e.}}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]} \right) + \left( \frac{\text{Bias}_{\theta,k}(\widehat{\psi}_1) + \text{EB}_{\theta,3,k}(\widehat{\Sigma}_k^{-1})}{\widehat{\text{s.e.}}_\theta[\widehat{\psi}_1]} - \delta \right) \frac{\widehat{\text{s.e.}}_\theta[\widehat{\psi}_1]}{\widehat{\text{s.e.}}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]} \right\} (1 + o_{\mathbf{P}_\theta}(1)) \\
&= \left\{ \left( \frac{\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1}) - \mathbb{E}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]}{\widehat{\text{s.e.}}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]} \right) + \left( \gamma + \frac{\text{EB}_{\theta,3,k}(\widehat{\Sigma}_k^{-1})}{\widehat{\text{s.e.}}_\theta[\widehat{\psi}_1]} - \delta \right) \frac{\widehat{\text{s.e.}}_\theta[\widehat{\psi}_1]}{\widehat{\text{s.e.}}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]} \right\} (1 + o_{\mathbf{P}_\theta}(1)) \\
&= \left\{ \underbrace{\left( \frac{\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1}) - \mathbb{E}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]}{\widehat{\text{s.e.}}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]} \right)}_A - (\delta - \gamma) \frac{\widehat{\text{s.e.}}_\theta[\widehat{\psi}_1]}{\widehat{\text{s.e.}}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]} + \underbrace{\frac{\text{EB}_{\theta,3,k}(\widehat{\Sigma}_k^{-1})}{\widehat{\text{s.e.}}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]}}_B \right\} (1 + o_{\mathbf{P}_\theta}(1)).
\end{aligned}$$

The effect of estimating  $\Sigma_k^{-1}$  on the asymptotic validity of the test  $\widehat{\chi}_{3,k}(\widehat{\Sigma}_k^{-1}; \varsigma_k, \delta)$  of  $\text{H}_{0,k}(\delta)$  thus depends on the orders of terms A and B. A has variance 1 and mean  $-(\delta - \gamma) \frac{\widehat{\text{s.e.}}_\theta[\widehat{\psi}_1]}{\widehat{\text{s.e.}}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]}$ . B depends on the estimation bias due to estimating  $\Sigma_k^{-1}$  by  $\widehat{\Sigma}_k^{-1}$ . Hence if we have:

- (1)  $\frac{\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1}) - \mathbb{E}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]}{\widehat{\text{s.e.}}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]}$  is asymptotically  $N(0, 1)$  conditional on the training sample;
- (2)  $\frac{\widehat{\text{s.e.}}_\theta[\widehat{\psi}_1]}{\widehat{\text{s.e.}}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]} / \frac{\widehat{\text{s.e.}}_\theta[\widehat{\psi}_1]}{\widehat{\text{s.e.}}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]} \rightarrow 1$ ;
- (3)  $B = o(1)$  under  $\text{H}_{0,k}(\delta)$  or fixed alternatives to  $\text{H}_{0,k}(\delta)$ ;
- (3')  $B \ll -(\delta - \gamma) \frac{\widehat{\text{s.e.}}_\theta[\widehat{\psi}_1]}{\widehat{\text{s.e.}}_\theta[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]}$  under diverging alternatives to  $\text{H}_{0,k}(\delta)$  i.e.  $\delta - \gamma = c$  for some  $c \rightarrow \infty$  (at any rate).

$\widehat{\chi}_{3,k}(\widehat{\Sigma}_k^{-1}; \varsigma_k, \delta)$  is an asymptotically level  $2 - 2\Phi(\varsigma_k)$  two-sided test for  $\text{H}_{0,k}(\delta)$  and rejects the null with probability approaching 1 under diverging alternatives to  $\text{H}_{0,k}(\delta)$ .

According to Proposition 4.2, both (1) and (2) hold for  $\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})$ . In terms of (3) and (3'), under the conditions of Proposition 4.3:

- Under  $H_{0,k}(\delta)$  or fixed alternatives to  $H_{0,k}(\delta)$ ,  $EB_{\theta,3,k}(\widehat{\Sigma}_k^{-1}) \lesssim \frac{k \log(k)}{n^{3/2}} \ll \frac{\sqrt{k}}{n} + \frac{1}{\sqrt{n}}$  and hence  $B = o(1)$ . Thus (3) is satisfied.
- Under diverging alternatives to  $H_{0,k}(\delta)$  i.e.  $\gamma - \delta = c \rightarrow \infty$ ,  $(\gamma - \delta) \text{s.e.}_{\theta}(\widehat{\psi}_1) \asymp \mathbb{L}_{\theta,2,\widehat{b},k} \mathbb{L}_{\theta,2,\widehat{p},k}$ ,

$$\begin{aligned} B &\lesssim \frac{\mathbb{L}_{\theta,2,\widehat{b},k} \mathbb{L}_{\theta,2,\widehat{p},k}}{\text{s.e.}_{\theta}[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]} \frac{k \log(k)}{n} \ll \frac{\mathbb{L}_{\theta,2,\widehat{b},k} \mathbb{L}_{\theta,2,\widehat{p},k}}{\text{s.e.}_{\theta}[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]} \\ &\asymp -(\delta - \gamma) \frac{\text{s.e.}_{\theta}[\widehat{\psi}_1]}{\text{s.e.}_{\theta}[\widehat{\mathbb{I}\mathbb{F}}_{22 \rightarrow 33,k}(\widehat{\Sigma}_k^{-1})]}. \end{aligned}$$

Thus (3') is satisfied.

In summary,  $\widehat{\chi}_{3,k}(\widehat{\Sigma}_k^{-1}; \varsigma_k, \delta)$  is an asymptotically valid level  $2 - 2\Phi(\varsigma_k)$  two-sided test for  $H_{0,k}(\delta)$  and rejects the null with probability approaching 1 under diverging alternatives to  $H_{0,k}(\delta)$ .

## APPENDIX B. SUPPLEMENTARY INFORMATION FOR SIMULATION STUDIES

The functions  $h_f$ ,  $h_b$  and  $h_p$  appeared in Section 6 are of the following forms:

$$(B.1) \quad h_f(x; s_f) \propto 1 + \exp \left\{ \frac{1}{2} \sum_{j \in \mathcal{J}, \ell \in \mathbb{Z}} 2^{-j(s_f + 0.25)} \omega_{j,\ell}(x) \right\},$$

$$(B.2) \quad h_b(x; s_b) = \sum_{j \in \mathcal{J}, \ell \in \mathbb{Z}} 2^{-j(s_b + 0.25)} \omega_{j,\ell}(x),$$

$$(B.3) \quad h_p(x; s_p) = \expit \left\{ -2 \sum_{j \in \mathcal{J}, \ell \in \mathbb{Z}} 2^{-j(s_p + 0.25)} \omega_{j,\ell}(x) \right\}$$

where  $\mathcal{J} = \{0, 3, 6, 9, 10, 16\}$  and  $\omega_{j,\ell}(\cdot)$  is the D12 (or equivalently db6) father wavelets function dilated at resolution  $j$ , shifted by  $\ell$  (Daubechies, 1992; Härdle et al., 1998; Mallat, 1999). Härdle et al. (1998)[Theorem 9.6] indeed implies that  $h_f(\cdot; s_f) \in \text{Hölder}(s_f)$ ,  $h_b(\cdot; s_b) \in \text{Hölder}(s_b)$  and  $h_p(\cdot; s_p) \in \text{Hölder}(s_p)$ . We fix  $s_f = 0.1$ . In simulation setup I we choose  $s_b = s_p = 0.25$  whereas in simulation setup II we choose  $s_b = s_p = 0.6$ .

In Table 5, we provide the numerical values for  $(\tau_{b,j}, \tau_{p,j})_{j=1}^8$  used in generating the simulation experiments in Section 6.

**B.1. Generating correlated multidimensional covariates  $X$  with fixed non-smooth marginal densities.** In the simulation study conducted in Section 6, one key step of generating the simulated datasets is to draw correlated multidimensional covariates  $X \in [0, 1]^d$  with fixed non-smooth marginal densities. First, we fix the marginal densities of  $X$  in each dimension proportional to  $h_f(\cdot)$  (equation (B.1)). Then we make  $2K$  independent draws of  $\widetilde{X}_{i,j}$ ,

$j$	$\tau_{b,j}$	$\tau_{p,j}$
1	-0.2819	0.09789
2	0.4876	0.08800
3	-0.1515	-0.4823
4	-0.1190	0.4588

TABLE 5. Coefficients used in constructing  $b$  and  $\pi$  in Section 6.

$i = 1, \dots, 2K$ , from  $h_f$  for every  $j = 1, \dots, d$  so  $\tilde{X} = (\tilde{X}_{1,\cdot}, \dots, \tilde{X}_{2K,\cdot})^\top \in [0, 1]^{2K \times d}$ . Next, to create correlations between different dimensions, we follow the strategy proposed in Baker (2008). First we group every two consecutive draws:  $(\tilde{X}_{1,\cdot}, \tilde{X}_{2,\cdot})^\top, (\tilde{X}_{3,\cdot}, \tilde{X}_{4,\cdot})^\top, \dots, (\tilde{X}_{2K-1,\cdot}, \tilde{X}_{2K,\cdot})^\top$ . Then for each pair  $(\tilde{X}_{2i-1,\cdot}, \tilde{X}_{2i,\cdot})^\top$  for  $i = 1, \dots, K$ , we form the following  $d$ -dimensional random vectors

$$U_i := (\max(\tilde{X}_{2i-1,1}, \tilde{X}_{2i,1}), \dots, \max(\tilde{X}_{2i-1,d}, \tilde{X}_{2i,d}))^\top,$$

$$V_i := (\min(\tilde{X}_{2i-1,1}, \tilde{X}_{2i,1}), \dots, \min(\tilde{X}_{2i-1,d}, \tilde{X}_{2i,d}))^\top.$$

Lastly, we construct  $K$  independent  $d$ -dimensional vectors  $X$  by the following rule: for each  $i = 1, \dots, K$ , we draw a Bernoulli random variable  $B_i$  with probability  $1/2$ , and if  $B_i = 0$ ,  $X_{i,\cdot} = U_i$ , otherwise  $X_{i,\cdot} = V_i$ . Following the above strategy, we conserve the marginal density of  $X_{\cdot,j}$  as that of  $\tilde{X}_{\cdot,j}$  but create dependence between different dimensions.