

# Kernel Regression Utilizing Heterogeneous Datasets

Chi-Shian Dai<sup>1</sup> and Jun Shao<sup>2 1</sup>

<sup>1</sup>*Department of Statistics, University of Wisconsin-Madison*

<sup>2</sup>*School of Statistics, East China Normal University*

**Abstract:** Data analysis in modern scientific research and practice has shifted from analyzing a single dataset to coupling several datasets. We propose and study a kernel regression method that can handle the challenge of heterogeneous populations and greatly extends the constrained kernel regression (Dai and Shao, 2022) for homogeneous populations of different datasets. Asymptotic normality of proposed estimators is established under some conditions and simulation results are presented to confirm our theory and to quantify the improvements from datasets with heterogeneous populations.

*Key words and phrases:* Conditional expectation, constraints, data coupling and integration, external data, heterogeneous populations, kernel estimation.

# 1 Introduction

With advanced technologies in data collection and storage, in modern statistical analyses we have not only a primary random sample from a population of interest, which results in a dataset referred to as the internal dataset, but also some independent external datasets from sources such as past investigations and publicly available administrative datasets. In this paper, we consider nonparametric kernel regression (Bierens, 1987; Wand and Jones, 1994; Wasserman, 2006) between a univariate response  $Y$  and a covariate vector  $\mathbf{U}$  from a sampled subject, using the internal dataset with the help from independent external datasets. Specifically, we consider kernel estimation of the conditional expectation (regression function) of  $Y$  given  $\mathbf{U} = \mathbf{u}$  under internal data population,

$$\mu_1(\mathbf{u}) = E(Y \mid \mathbf{U} = \mathbf{u}, D = 1), \quad (1)$$

where  $D = 1$  indicates internal population and  $\mathbf{u}$  is a fixed point in  $\mathcal{U}$ , the range of  $\mathbf{U}$ . The subscript 1 in  $\mu_1(\mathbf{u})$  emphasizes that it is for internal data population ( $D = 1$ ), which may be different from  $\mu(\mathbf{u}) = E(Y \mid \mathbf{U} = \mathbf{u})$ , a mixture of quantities from the internal and external data populations.

When the external dataset also have measurements  $Y$  and  $\mathbf{U}$ , we may simply combine the internal and external datasets when the populations for internal and external data are identical (homogeneous). However, heterogeneity typically exists among populations for different datasets, especially when there are multiple external datasets collected in different ways and/or different time periods. In Section 2.1, we propose a method to handle heterogeneity among different populations and derive a kernel regression more efficient than the one using internal data alone. The result is also a building block for the more complicated

case in Section 2.2 where the external dataset contains fewer measured covariates.

In applications, it is often the case that the external dataset has measured  $Y$  and  $\mathbf{X}$  from each subject, where  $\mathbf{X}$  is a part of the vector  $\mathbf{U}$  with some components of  $\mathbf{U}$  not measured due to the high measurement cost or the progress of new technology and/or new scientific relevance for measuring some components of  $\mathbf{U}$ . With some unmeasured components of  $\mathbf{U}$ , the external datasets cannot be directly used to estimate  $\mu_1(\mathbf{u})$  in (1), since conditioning on the entire  $\mathbf{U}$  is involved. To solve this problem, Dai and Shao (2022) proposes a two-step kernel regression using external information as a constraint to improve kernel regression based on internal data along, following the idea of using constraints in Chatterjee et al. (2016) and Zhang et al. (2020). However, these three cited papers mainly assume that the internal and external datasets share the same population, which may be unrealistic. The challenge in dealing with the heterogeneity among different populations is similar to that in handling nonignorable missing data problems if unmeasured components of  $\mathbf{U}$  is treated as missing data, although in missing data problems we usually want to estimate  $\mu(\mathbf{u}) = E(Y \mid \mathbf{U} = \mathbf{u}) \neq \mu_1(\mathbf{u})$  in (1).

In Section 2.2, we develop a methodology to handle population heterogeneity for internal and external datasets, which extends the procedure in Dai and Shao (2022) to heterogeneous populations and greatly widens its application scope.

Under each scenario, we derive asymptotic normality in Section 3 for the proposed kernel estimators and obtain explicitly the asymptotic variances, which is important for large sample inference. Some simulation results are presented in Section 4 to compare finite sample performance of several estimators. Discussions on extensions and handling high dimension covariates are given in Section 5. All technical details are in the Appendix.

Our research fits into a general framework of data integration (Kim et al., 2021; Lohr and Raghunathan, 2017; Merkouris, 2004; Rao, 2021; Yang and Kim, 2020; Zhang et al., 2017; Zieschang, 1990).

## 2 Methodology

The internal dataset contains observations  $(Y_i, \mathbf{U}_i)$ ,  $i = 1, \dots, n$ , independent and identically distributed (iid) from  $\mathcal{P}_1$ , the internal population of  $(Y, \mathbf{U})$ , where  $Y$  is the response and  $\mathbf{U}$  is a  $p$ -dimensional covariate vector associated with  $Y$ . We are interested in the estimation of conditional expectation  $\mu_1(\mathbf{u})$  in (1). The standard kernel regression estimator of  $\mu_1(\mathbf{u})$  based on the internal dataset alone is

$$\hat{\mu}_1(\mathbf{u}) = \sum_{i=1}^n Y_i \kappa_b(\mathbf{u} - \mathbf{U}_i) / \sum_{i=1}^n \kappa_b(\mathbf{u} - \mathbf{U}_i), \quad (2)$$

where  $\kappa_b(\mathbf{a}) = b^{-p} \kappa(\mathbf{a}/b)$ ,  $\mathbf{a} \in \mathcal{U}$ ,  $\kappa(\mathbf{u})$  is a given kernel function, and  $b > 0$  is a bandwidth depending on  $n$ . We assume that  $\mathbf{U}$  is standardized so that the same bandwidth  $b$  is used for every component of  $\mathbf{U}$  in kernel regression. Because of the well-known curse of dimensionality for kernel-type methods, we focus on a low dimension  $p$  not varying with  $n$ . A discussion of handling a large dimensional  $\mathbf{U}$  is given in Section 5.

We consider the case with one external dataset, independent of the internal dataset. Extension to multiple external datasets is straightforward and discussed in Section 5.

### 2.1 The Case of No Unmeasured Covariates

In this subsection we consider the situation where the external dataset contains iid observations  $(Y_i, \mathbf{U}_i)$ ,  $i = n + 1, \dots, N$ , from  $\mathcal{P}_0$ , the external population of  $(Y, \mathbf{U})$ .

If we assume that the two populations  $\mathcal{P}_1$  and  $\mathcal{P}_0$  are identical, then we can simply combine two datasets to obtain the kernel estimator

$$\hat{\mu}_1^{E1}(\mathbf{u}) = \sum_{i=1}^N Y_i \kappa_b(\mathbf{u} - \mathbf{U}_i) \bigg/ \sum_{i=1}^N \kappa_b(\mathbf{u} - \mathbf{U}_i), \quad (3)$$

which is obviously more efficient than  $\hat{\mu}_1(\mathbf{u})$  in (2) as the sample size is increased to  $N > n$ . The estimator  $\hat{\mu}_1^{E1}(\mathbf{u})$  in (3), however, is not correct (i.e., it is biased) when populations  $\mathcal{P}_1$  and  $\mathcal{P}_0$  are different, because  $E(Y | \mathbf{U} = \mathbf{u}, D = 0)$  for external population may be different from  $\mu_1(\mathbf{u}) = E(Y | \mathbf{U} = \mathbf{u}, D = 1)$  for internal population.

We now derive a kernel estimator using two datasets and is asymptotically correct regardless of whether  $\mathcal{P}_1$  and  $\mathcal{P}_0$  are the same or not. Let  $f(y|\mathbf{u}, D)$  be the conditional density of  $Y$  given  $\mathbf{U} = \mathbf{u}$  and  $D = 1$  or  $0$  (for internal or external population). Then

$$\mu_1(\mathbf{x}) = E(Y|\mathbf{U} = \mathbf{u}, D = 1) = E \left\{ Y \frac{f(Y|\mathbf{u}, D = 1)}{f(Y|\mathbf{u}, D = 0)} \mid \mathbf{U} = \mathbf{u}, D = 0 \right\}. \quad (4)$$

The ratio  $f(Y|\mathbf{u}, D = 1)/f(Y|\mathbf{u}, D = 0)$  links internal and external populations so that we can overcome the difficulty in utilizing the external data under heterogeneous populations.

If we can construct an estimator  $\hat{f}(y|\mathbf{u}, D)$  of  $f(y|\mathbf{u}, D)$  for every  $y$ ,  $\mathbf{u}$ , and  $D = 0$  or  $1$ , then we can modify the estimator in (3) by replacing every  $Y_i$  with  $i > n$  by constructed response  $\hat{Y}_i = Y_i \hat{f}(Y_i|\mathbf{U}_i, D = 1)/\hat{f}(Y_i|\mathbf{U}_i, D = 0)$ . The resulting kernel estimator is

$$\hat{\mu}_1^{E2}(\mathbf{u}) = \left\{ \sum_{i=1}^n Y_i \kappa_b(\mathbf{u} - \mathbf{U}_i) + \sum_{i=n+1}^N \hat{Y}_i \kappa_b(\mathbf{u} - \mathbf{U}_i) \right\} \bigg/ \sum_{i=1}^N \kappa_b(\mathbf{u} - \mathbf{U}_i). \quad (5)$$

Note that we use internal data  $(Y_i, \mathbf{U}_i)$ ,  $i = 1, \dots, n$ , to obtain estimator  $\hat{f}(Y_i|\mathbf{U}_i, D = 1)$  and external data  $(Y_i, \mathbf{U}_i)$ ,  $i = n+1, \dots, N$ , to construct estimator  $\hat{f}(Y_i|\mathbf{U}_i, D = 0)$ . Applying kernel estimation, we obtain that

$$\begin{aligned}
\widehat{f}(y|\mathbf{U} = \mathbf{u}, D = 1) &= \sum_{i=1}^n \tilde{\kappa}_{\tilde{b}}(y - Y_i, \mathbf{u} - \mathbf{U}_i) \bigg/ \sum_{i=1}^n \bar{\kappa}_{\bar{b}}(\mathbf{u} - \mathbf{U}_i), \\
\widehat{f}(y|\mathbf{U} = \mathbf{u}, D = 0) &= \sum_{i=n+1}^N \tilde{\kappa}_{\tilde{b}}(y - Y_i, \mathbf{u} - \mathbf{U}_i) \bigg/ \sum_{i=n+1}^N \bar{\kappa}_{\bar{b}}(\mathbf{u} - \mathbf{U}_i),
\end{aligned} \tag{6}$$

where  $\tilde{\kappa}$  and  $\bar{\kappa}$  are kernels with bandwidths  $\tilde{b}$  and  $\bar{b}$ , respectively. The estimator in (5) is asymptotically valid under some regularity conditions for kernel and bandwidth, summarized in Theorem 1 of Section 3.

If additional information exists, then this approach can be improved. Assume that the internal and external datasets are formed according to a random binary indicator  $D$  such that  $(Y_i, \mathbf{U}_i, D_i)$ ,  $i = 1, \dots, N$ , are iid distributed as  $(Y, \mathbf{U}, D)$ , where  $Y_i$  and  $\mathbf{U}_i$  are observed internal data when  $D_i = 1$ ,  $Y_i$  and  $\mathbf{U}_i$  are observed external data when  $D_i = 0$ , and  $N$  is still the known total sample size for internal and external data. In this situation, the internal and external sample sizes are  $n = \sum_{i=1}^N D_i$  and  $N - n$ , respectively, both of which are random. In most applications, this assumption is not substantial. From the identity

$$\frac{f(Y|\mathbf{u}, D = 1)}{f(Y|\mathbf{u}, D = 0)} = \frac{P(D = 1|\mathbf{U} = \mathbf{u}, Y)}{P(D = 0|\mathbf{U} = \mathbf{u}, Y)} \frac{P(D = 0|\mathbf{U} = \mathbf{u})}{P(D = 1|\mathbf{U} = \mathbf{u})}. \tag{7}$$

we just need to estimate  $P(D = 1|\mathbf{U} = \mathbf{u}, Y)$  and  $P(D = 1|\mathbf{U} = \mathbf{u})$  for every  $\mathbf{u}$ , constructed using for example the nonparametric estimators in Fan et al. (1998) for binary response. For each estimator, both internal and external data on  $(Y, \mathbf{U})$  and the indicator  $D$  are used.

A further improved can be made if the following semi-parametric model holds,

$$\frac{P(D = 0 | \mathbf{U}, Y)}{P(D = 1 | \mathbf{U}, Y)} = \exp\{\alpha(\mathbf{U}) + \gamma Y\}, \tag{8}$$

where  $\alpha(\cdot)$  is an unspecified unknown function and  $\gamma$  is an unknown parameter. It follows from (7)-(8) that

$$\frac{f(Y|\mathbf{u}, D=1)}{f(Y|\mathbf{u}, D=0)} = e^{-\gamma Y} E(e^{\gamma Y} | \mathbf{U} = \mathbf{u}, D=1). \quad (9)$$

If  $\gamma = 0$ , then  $f(Y|\mathbf{u}, D=1) = f(Y|\mathbf{u}, D=0)$  and the estimator in (3) is correct. Under (9) with  $\gamma \neq 0$ , we just need to derive an estimator  $\hat{\gamma}$  of  $\gamma$  and apply kernel estimation to estimate  $E(e^{\hat{\gamma} Y} | \mathbf{U} = \mathbf{u}, D=1)$  as a function of  $\mathbf{u}$ . Note that we do not need to estimate the unspecified function  $\alpha(\cdot)$  in (8), which is a nice feature of semi-parametric model (8).

We now derive an estimator  $\hat{\gamma}$ . Applying (7)-(8) to (4), we obtain that

$$\begin{aligned} \mu_1(\mathbf{u}) &= E \left\{ Y \frac{P(D=1|\mathbf{U}=\mathbf{u}, Y)}{P(D=0|\mathbf{U}=\mathbf{u}, Y)} \middle| \mathbf{U} = \mathbf{u}, D=0 \right\} \frac{P(D=0|\mathbf{U}=\mathbf{u})}{P(D=1|\mathbf{U}=\mathbf{u})} \\ &= E \left( Y e^{-\alpha(\mathbf{u}) - \gamma Y} | \mathbf{U} = \mathbf{u}, D=0 \right) \frac{E\{P(D=0|\mathbf{U}=\mathbf{u}, Y)|\mathbf{U}=\mathbf{u}\}}{P(D=1|\mathbf{U}=\mathbf{u})} \\ &= e^{-\alpha(\mathbf{u})} E \left( Y e^{-\gamma Y} | \mathbf{U} = \mathbf{u}, D=0 \right) \frac{E\{e^{\alpha(\mathbf{u}) + \gamma Y} P(D=1|\mathbf{U}=\mathbf{u}, Y)|\mathbf{U}=\mathbf{u}\}}{P(D=1|\mathbf{U}=\mathbf{u})} \\ &= E \left( Y e^{-\gamma Y} | \mathbf{U} = \mathbf{u}, D=0 \right) \frac{E\{e^{\gamma Y} E(D|\mathbf{U}=\mathbf{u}, Y)|\mathbf{U}=\mathbf{u}\}}{P(D=1|\mathbf{U}=\mathbf{u})} \\ &= E \left( Y e^{-\gamma Y} | \mathbf{U} = \mathbf{u}, D=0 \right) \frac{E(e^{\gamma Y} D | \mathbf{U} = \mathbf{u})}{P(D=1|\mathbf{U}=\mathbf{u})} \\ &= E \left( Y e^{-\gamma Y} | \mathbf{U} = \mathbf{u}, D=0 \right) E \left( e^{\gamma Y} | \mathbf{U} = \mathbf{u}, D=1 \right), \end{aligned}$$

where the second and third equalities follow from (8).

I change the first and second to the second and third – Chi

For every real number  $t$ , define

$$h(\mathbf{u}, t) = E(Y e^{-tY} | \mathbf{U} = \mathbf{u}, D=0) E(e^{tY} | \mathbf{U} = \mathbf{u}, D=1).$$

Its estimator by kernel regression is

$$\hat{h}(\mathbf{u}, t) = \frac{\sum_{i=1}^N (1 - D_i) \check{\kappa}_{\check{b}}(\mathbf{u} - \mathbf{U}_i) Y_i e^{-tY_i}}{\sum_{i=1}^N (1 - D_i) \check{\kappa}_{\check{b}}(\mathbf{u} - \mathbf{U}_i)} \frac{\sum_{i=1}^N D_i \check{\kappa}_{\check{b}}(\mathbf{u} - \mathbf{U}_i) e^{tY_i}}{\sum_{i=1}^N D_i \check{\kappa}_{\check{b}}(\mathbf{u} - \mathbf{U}_i)} \quad (10)$$

where  $\check{\kappa}$  is a kernel and  $\check{b}$  is a bandwidth. Then, we estimate  $\gamma$  by

$$\hat{\gamma} = \arg \min_t \frac{1}{N} \sum_{i=1}^N D_i \{Y_i - \hat{h}(\mathbf{U}_i, t)\}^2, \quad (11)$$

motivated by the fact that the objective function for minimization in (11) approximates  $E[D\{Y - h(\mathbf{U}, t)\}^2 | D = 1]$  and, for any  $t$ ,

$$E[D\{Y - h(\mathbf{U}, \gamma)\}^2 | D = 1] \leq E[D\{Y - h(\mathbf{U}, t)\}^2 | D = 1]$$

because  $h(\mathbf{u}, \gamma) = \mu_1(\mathbf{u})$ .

Once  $\hat{\gamma}$  is obtained, our estimator of  $\mu_1(\mathbf{u})$  is

$$\hat{\mu}_1^{E3}(\mathbf{u}) = \left\{ \sum_{i=1}^N D_i Y_i \kappa_b(\mathbf{u} - \mathbf{U}_i) + \sum_{i=1}^N (1 - D_i) \hat{Y}_i \kappa_b(\mathbf{u} - \mathbf{U}_i) \right\} / \sum_{i=1}^N \kappa_b(\mathbf{u} - \mathbf{U}_i) \quad (12)$$

with

$$\hat{Y}_i = Y_i e^{-\hat{\gamma} Y_i} \sum_{j=1}^n e^{\hat{\gamma} Y_j} \check{\kappa}_{\check{b}}(\mathbf{U}_i - \mathbf{U}_j) / \sum_{j=1}^n \check{\kappa}_{\check{b}}(\mathbf{U}_i - \mathbf{U}_j),$$

in view of (9).

In applications, we need to choose bandwidths with given sample sizes  $n$  and  $N - n$ . We can apply the  $k$ -fold cross-validation as described in Györfi et al. (2002). Requirements on the rates of bandwidths are described in theorems in Section 3.

## 2.2 Constrained Kernel Regression with Unmeasured Covariates

We still consider the case with one external dataset, independent of the internal dataset. In this subsection, the external dataset contains iid observations  $(Y_i, \mathbf{X}_i)$ ,  $i = n + 1, \dots, N$ , from



the external population  $\mathcal{P}_0$ , where  $\mathbf{X}$  is a  $q$ -dimensional sub-vector of  $\mathbf{U}$  with  $q < p$ .

Since the external dataset has only  $\mathbf{X}$ , not the entire  $\mathbf{U}$ , we cannot apply the method in Section 2.1 when  $q < p$ . Instead, we consider kernel regression using external information in a constraint. First, we consider the estimation of the  $n$ -dimensional vector  $\boldsymbol{\mu}_1 = (\mu_1(\mathbf{U}_1), \dots, \mu_1(\mathbf{U}_n))^\top$ , where  $\mathbf{A}^\top$  denotes the transpose of vector or matrix  $\mathbf{A}$  throughout. Note that the standard kernel regression (2) actually estimates  $\boldsymbol{\mu}_1$  by

$$\hat{\boldsymbol{\mu}}_1 = \arg \min_{\mu_1, \dots, \mu_n} \sum_{i=1}^n \sum_{j=1}^n \kappa_b(\mathbf{U}_i - \mathbf{U}_j) (Y_j - \mu_i)^2 \bigg/ \sum_{k=1}^n \kappa_b(\mathbf{U}_i - \mathbf{U}_k). \quad (13)$$

We improve  $\hat{\boldsymbol{\mu}}_1$  by the following constrained minimization,

$$\hat{\boldsymbol{\mu}}_1^C = \arg \min_{\mu_1, \dots, \mu_n} \sum_{i=1}^n \sum_{j=1}^n \kappa_l(\mathbf{U}_i - \mathbf{U}_j) (Y_j - \mu_i)^2 \bigg/ \sum_{k=1}^n \kappa_l(\mathbf{U}_i - \mathbf{U}_k) \quad (14)$$

$$\text{subject to } \sum_{i=1}^n \{\mu_i - \hat{h}_1(\mathbf{X}_i)\} \mathbf{g}(\mathbf{X}_i)^\top = 0, \quad (15)$$

where  $\mathbf{g}(\mathbf{x})^\top = (1, \mathbf{x}^\top)$ ,  $l$  in (14) is a bandwidth that may be different from  $b$  in (2) or (13), and  $\hat{h}_1(\mathbf{x})$  is the kernel estimator of  $h_1(\mathbf{x}) = E(Y \mid \mathbf{X} = \mathbf{x}, D = 1)$  using one of the methods described in Section 2.1 with  $\mathbf{U}$  replaced by  $\mathbf{X}$ . Note that this can be done since both internal and external datasets have measured  $\mathbf{X}_i$ 's. More details of  $\hat{h}_1(\mathbf{x})$  are given in the end of this subsection.

Constraint (15) is an empirical analog (based on internal data) of

$$E \left[ \left\{ \mu_1(\mathbf{U}) - h_1(\mathbf{X}) \right\} \mathbf{g}(\mathbf{X})^\top \middle| D = 1 \right] = 0,$$

as  $E \{ E(Y \mid \mathbf{U}, D = 1) \mid \mathbf{X}, D = 1 \} = E(Y \mid \mathbf{X}, D = 1) = h_1(\mathbf{X})$ . Thus, if  $\hat{h}_1(\cdot)$  is a good estimator of  $h_1(\cdot)$ , then  $\hat{\boldsymbol{\mu}}_1^C$  in (14) is more accurate than the unconstrained  $\hat{\boldsymbol{\mu}}_1$  in (13).

It turns out that  $\hat{\boldsymbol{\mu}}_1^C$  in (14) has an explicit form  $\hat{\boldsymbol{\mu}}_1^C = \hat{\boldsymbol{\mu}}_1 + \mathbf{G}(\mathbf{G}^\top \mathbf{G})^{-1} \mathbf{G}^\top (\hat{\mathbf{h}}_1 - \hat{\boldsymbol{\mu}}_1)$ ,

where  $\mathbf{G}$  is the  $n \times n$  matrix whose  $i$ th row is  $\mathbf{g}(\mathbf{X}_i)^\top$  and  $\hat{\mathbf{h}}_1$  is the  $n$ -dimensional vector whose  $i$ th component is  $\hat{h}_1(\mathbf{X}_i)$ .

To obtain an improved estimator of the entire regression function  $\mu_1(\cdot)$  in (1), not just the function at  $\mathbf{u} = \mathbf{U}_i$ ,  $i = 1, \dots, n$ , we apply the standard kernel regression with response vector  $(Y_1, \dots, Y_n)^\top$  replaced by  $\hat{\boldsymbol{\mu}}_1^C$  in (14), which results in an estimator of  $\mu_1(\mathbf{u})$  for any  $\mathbf{u} \in \mathcal{U}$  as

$$\hat{\mu}_1^C(\mathbf{u}) = \sum_{i=1}^n \hat{\mu}_i \kappa_b(\mathbf{u} - \mathbf{U}_i) / \sum_{i=1}^n \kappa_b(\mathbf{u} - \mathbf{U}_i), \quad (16)$$

where  $\hat{\mu}_i$  is the  $i$ th component of  $\hat{\boldsymbol{\mu}}_1^C$  and  $b$  is the same bandwidth in (2).

With  $\mathbf{U}$  replaced by  $\mathbf{X}$ , there are essentially three methods described in Section 2.1 for constructing the kernel estimator  $\hat{h}_1(\mathbf{x})$  of  $h_1(\mathbf{x}) = E(Y \mid \mathbf{X} = \mathbf{x}, D = 1)$ . The first one is the simplest estimator (3) (with  $\mathbf{U}$  replaced by  $\mathbf{X}$ ), which may be incorrect unless populations  $\mathcal{P}_1$  and  $\mathcal{P}_0$  are identical. We still denote the resulting estimator in (16) as  $\hat{\mu}_1^{C1}(\mathbf{u})$ .

The second one is (5) with  $\mathbf{U}$  replaced by  $\mathbf{X}$  and  $\hat{f}(Y_i \mid \mathbf{X}_i, D = 1) / \hat{f}(Y_i \mid \mathbf{X}_i, D = 0)$  constructed either by the ratio of two kernel estimators given by (6) or kernel estimators for binary responses through (7), with  $\mathbf{U}$  replaced by  $\mathbf{X}$ . The resulting estimator in (16) is denoted as  $\hat{\mu}_1^{C2}(\mathbf{u})$ .

The third and last one uses  $\hat{h}(\mathbf{X}_i, \hat{\gamma})$  for  $\hat{h}_1(\mathbf{X}_i)$  in constraint (15), where  $\hat{\gamma}$  and  $h(\mathbf{X}_i, t)$  are defined by (11) and (10), respectively, with  $\mathbf{U}$  replaced by  $\mathbf{X}$ . The resulting estimator in (16) is denoted as  $\hat{\mu}_1^{C3}(\mathbf{u})$ .

### 3 Theory

We now establish the asymptotic normality of  $\hat{\mu}_1^E(\mathbf{u})$  and  $\hat{\mu}_1^C(\mathbf{u})$  for a fixed  $\mathbf{u}$ , as the sample size of the internal dataset increases to infinity. All technical proofs are given in the

Appendix.

The first result is about  $\hat{\mu}_1^{E2}(\mathbf{u})$  in (5). The result is also applicable to  $\hat{\mu}_1^{E1}(\mathbf{u})$  in (3) with an added condition that  $\mathcal{P}_1 = \mathcal{P}_0$ .

**Theorem 1.** *Assume the following conditions.*

- (B1) *The densities  $f_1(\mathbf{u})$  and  $f_0(\mathbf{u})$  for  $\mathbf{U}$  respectively under internal and external populations have continuous and bounded first and second order derivatives.*
- (B2)  *$\mu_1^2(\mathbf{u})f_k(\mathbf{u})$ ,  $\sigma_k^2(\mathbf{u})f_k(\mathbf{u})$ , and the first and second order derivatives of  $\mu_1(\mathbf{u})f_k(\mathbf{u})$  are continuous and bounded, where  $\sigma_1^2(\mathbf{u}) = E[\{Y - \mu_1(\mathbf{U})\}^2 \mid \mathbf{U} = \mathbf{u}, D = 1]$ ,  $\sigma_0^2(\mathbf{u}) = E[\{\tilde{Y} - \mu_1(\mathbf{U})\}^2 \mid \mathbf{U} = \mathbf{u}, D = 0]$ , and  $\tilde{Y} = Yf(Y|\mathbf{U}, D = 1)/f(Y|\mathbf{U}, D = 0)$ . Also,  $E(|Y|^s|\mathbf{U} = \mathbf{u}, D = 1)f_1(\mathbf{u})$  and  $E(|\tilde{Y}|^s|\mathbf{U} = \mathbf{u}, D = 0)f_0(\mathbf{u})$  are bounded for a constant  $s > 2$ .*
- (B3) *The kernel  $\kappa$  is second order, i.e.,  $\int \mathbf{u} \kappa(\mathbf{u}) d\mathbf{u} = 0$  and  $0 < \int \mathbf{u}^\top \mathbf{u} \kappa(\mathbf{u}) d\mathbf{u} < \infty$ .*
- (B4) *The bandwidth  $b$  satisfies  $b \rightarrow 0$  and  $(a + 1)nb^{p+4} \rightarrow c \in [0, \infty)$ , where  $a = \lim_{n \rightarrow \infty} (N - n)/n$  (assumed to exist without loss of generality).*
- (B5) *The kernels  $\tilde{\kappa}$  and  $\bar{\kappa}$  in (6) have bounded supports and orders  $\tilde{m} > 2 + 2/p$  and  $\bar{m} > 2$ , respectively, as defined by Bierens (1987),  $f(y, \mathbf{u}|D = 1)$ ,  $f(y, \mathbf{u}|D = 0)$  are  $\tilde{m}$ th order continuously differentiable with bounded derivatives, and  $f_1(\mathbf{u})$  and  $f_0(\mathbf{u})$  are  $\bar{m}$ th order continuously differentiable with bounded derivatives. Functions  $f(y, \mathbf{u}|D = 0)$  and  $f_1(\mathbf{u})$  are bounded away from zero. The bandwidths  $\tilde{b}$  and  $\bar{b}$  satisfy  $n\tilde{b}^{p+1}/\log(n) \rightarrow \infty$  and  $n\bar{b}^p/\log(n) \rightarrow \infty$ .*

Then, for any fixed  $\mathbf{u}$  with  $f_0(\mathbf{u}) > 0$  and  $f_1(\mathbf{u}) > 0$  and  $\hat{\mu}_1^{E2}$  in (5),

$$\sqrt{nb^p}\{\widehat{\mu}_1^{E2}(\mathbf{u}) - \mu_1(\mathbf{u})\} \xrightarrow{d} N(B_a(\mathbf{u}), V_a(\mathbf{u})), \quad (17)$$

where  $\xrightarrow{d}$  denotes convergence in distribution as  $n \rightarrow \infty$ ,

$$\begin{aligned} B_a(\mathbf{u}) &= \frac{c^{1/2}\{f_1(\mathbf{u})A_1(\mathbf{u}) + af_0(\mathbf{u})A_0(\mathbf{u})\}}{(a+1)^{1/2}\{f_1(\mathbf{u}) + af_0(\mathbf{u})\}}, \\ A_1(\mathbf{u}) &= \int \kappa(\mathbf{v}) \left\{ \frac{1}{2}\mathbf{v}^\top \nabla^2 \mu_1(\mathbf{u})\mathbf{v} + \mathbf{v}^\top \nabla \log f_1(\mathbf{u}) \nabla \mu_1(\mathbf{u})^\top \mathbf{v} \right\} d\mathbf{v}, \\ A_0(\mathbf{u}) &= \int \kappa(\mathbf{v}) \left\{ \frac{1}{2}\mathbf{v}^\top \nabla^2 \mu_1(\mathbf{u})\mathbf{v} + \mathbf{v}^\top \nabla \log f_0(\mathbf{u}) \nabla \mu_1(\mathbf{u})^\top \mathbf{v} \right\} d\mathbf{v}, \\ V_a(\mathbf{u}) &= \frac{f_1(\mathbf{u})\sigma_1^2(\mathbf{u}) + af_0(\mathbf{u})\sigma_0^2(\mathbf{u})}{\{f_1(\mathbf{u}) + af_0(\mathbf{u})\}^2} \int \kappa(\mathbf{v})^2 d\mathbf{v}. \end{aligned}$$

Conditions (B1)-(B4) are typically assumed for kernel estimation (Bierens, 1987). Condition (B5) is a sufficient condition for

$$\max_{i=n+1, \dots, N} \left| \frac{\widehat{f}(Y_i|\mathbf{U} = \mathbf{U}_i, D = 1)}{\widehat{f}(Y_i|\mathbf{U} = \mathbf{U}_i, D = 0)} - \frac{f(Y_i|\mathbf{U} = \mathbf{U}_i, D = 1)}{f(Y_i|\mathbf{U} = \mathbf{U}_i, D = 0)} \right| = \frac{o_p(1)}{\sqrt{nb^p}} \quad (18)$$

(Lemma 8.10 in Newey and McFadden (1994)), where  $o_p(1)$  denotes a term tending to 0 in probability. Result (18) implies that the estimation of ratio  $f(Y|\mathbf{U}, D = 1)/f(Y|\mathbf{U}, D = 0)$  does not affect the asymptotic distribution of  $\widehat{\mu}_1^{E2}(\mathbf{u})$  in (5).

Note that both the bias  $B_a(\mathbf{u})$  and variance  $V_a(\mathbf{u})$  in (17) are decreasing in the limit  $a = \lim_{n \rightarrow \infty} (N - n)/n$ , a quantity reflecting how many external data we have. In the extreme case of  $a = 0$ , i.e., the size of the external dataset is negligible compared with the size of the internal dataset, result (17) reduces to the well-known asymptotic normality for the standard kernel estimator  $\widehat{\mu}_1(\mathbf{u})$  in (2) (Bierens, 1987). In the other extreme case of  $a = \infty$ , on the other hand,  $B_a(\mathbf{u}) = V_a(\mathbf{u}) = 0$  and, hence,  $\widehat{\mu}_1^{E2}(\mathbf{u})$  has a convergence rate tending to 0 faster than  $1/\sqrt{nb^p}$ , the convergence rate of the standard kernel estimator  $\widehat{\mu}_1(\mathbf{u})$ .

The next result is about  $\hat{\mu}_1^{C2}(\mathbf{u})$  in (16) as described in the end of Section 2.2.

**Theorem 2.** *Assume (B1)-(B5) with  $\mathbf{U}$  and  $p$  replaced by  $\mathbf{X}$  and  $q$ , respectively, and the following conditions, where  $f_k(\mathbf{u})$  and  $\sigma_k^2(\mathbf{u})$ ,  $k = 0, 1$ , are defined in (B1)-(B2).*

(C1) *The range  $\mathcal{U}$  of  $\mathbf{U}$  is a compact set in the  $p$ -dimensional Euclidean space and  $f_1(\mathbf{u})$  is bounded away from infinity and zero on  $\mathcal{U}$ ;  $f_1(\mathbf{u})$  and  $f_0(\mathbf{u})$  have continuous and bounded first and second order derivatives.*

(C2) *Functions  $\mu_1(\mathbf{u}) = E(Y|\mathbf{U} = \mathbf{u})$  and  $\sigma_1^2(\mathbf{u})$  are Lipschitz continuous;  $\mu_1(\mathbf{u})$  has bounded third-order derivatives;  $h_1(\mathbf{x}) = E(Y | \mathbf{X} = \mathbf{x}, D = 1)$  has bounded first and second order derivatives; and  $E(|Y|^s | \mathbf{U} = \mathbf{u}, D = 1)$  is bounded with  $s > 2 + p/2$ .*

(C3) *All kernel functions are positive, bounded, and Lipschitz continuous with mean zero and finite sixth moments.*

(C4)  *$a = \lim_{n \rightarrow \infty} (N - n)/n > 0$  and the bandwidths  $b$  in (2) and  $l$  in (14) satisfy  $b \rightarrow 0$ ,  $l \rightarrow 0$ ,  $l/b \rightarrow r \in (0, \infty)$ ,  $nb^p \rightarrow \infty$ , and  $nb^{4+p} \rightarrow c \in [0, \infty)$ , as  $n \rightarrow \infty$ .*

(C5) *The densities  $f_{X0}(\mathbf{x})$  and  $f_{X1}(\mathbf{x})$  for  $\mathbf{X}$  respectively under internal and external populations are bounded away from zero. There exists a constant  $s > 4$  such that  $E(|Y|^s | D = 1)$  and  $E(|\tilde{Y}|^s | D = 0)$  are finite,  $E(|Y|^s | \mathbf{X} = \mathbf{x}, D = 1)f_{X1}(\mathbf{x})$  and  $E(|\tilde{Y}|^s | \mathbf{X} = \mathbf{x}, D = 0)f_{X0}(\mathbf{x})$  are bounded, and the bandwidth  $b_h$  for  $\hat{h}_1$  satisfies  $n^{1-2/s}b_h^q/\log(n) \rightarrow \infty$ . (Should we use something different for  $b$ ?)*

Use  $b_h$ ?? – Chi

Then, for any fixed  $\mathbf{u} \in \mathcal{U}$  and  $\hat{\mu}_1^{C2}(\mathbf{u})$  in (16),

$$\sqrt{nb^p}\{\widehat{\mu}_1^{C2}(\mathbf{u}) - \mu_1(\mathbf{u})\} \xrightarrow{d} N(B_r(\mathbf{u}), V_r(\mathbf{u})), \quad (19)$$

where

$$\begin{aligned} B_r(\mathbf{u}) &= c^{1/2}[(1+r^2)A_1(\mathbf{u}) - r^2\mathbf{g}(\mathbf{x})^\top \Sigma_g^{-1} E\{\mathbf{g}(\mathbf{X})A_1(\mathbf{U})|D=1\}], \\ A_1(\mathbf{u}) &= \int \kappa(\mathbf{v}) \left\{ \frac{1}{2}\mathbf{v}^\top \nabla^2 \mu_1(\mathbf{u})\mathbf{v} + \mathbf{v}^\top \nabla \log f_1(\mathbf{u}) \nabla \mu_1(\mathbf{u})^\top \mathbf{v} \right\} d\mathbf{v}, \\ V_r(\mathbf{u}) &= \frac{\sigma_1^2(\mathbf{u})}{f_1(\mathbf{u})} \int \left\{ \int \kappa(\mathbf{v} - r\mathbf{w})\kappa(\mathbf{w})d\mathbf{w} \right\}^2 d\mathbf{v}, \end{aligned}$$

and  $\Sigma = E\{\mathbf{g}(\mathbf{X})\mathbf{g}(\mathbf{X})^\top \mid D=1\}$  is assumed to be positive definite without loss of generality.

The next result is about  $\widehat{\gamma}$  in (11).

**Theorem 3.** *Suppose that (8) holds for binary random  $D$  indicating internal and external data. Assume also the following conditions.*

(D1) *The kernel  $\kappa$  in (10) is Lipschitz continuous, satisfies  $\int \kappa(\mathbf{u})d\mathbf{u} = 1$ , has a bounded support, and has order  $d > \max\{(p+4)/2, p\}$ .*

(D2) *The bandwidth  $\ell$  in (10) satisfies  $N\ell^{2q}/(\log N)^2 \rightarrow \infty$  and  $N\ell^{2d} \rightarrow 0$  as the total sample size of internal and external datasets  $N \rightarrow \infty$ , where  $d$  is given in (D1).*

(D3)  *$\gamma$  in (8) is an interior point of a compact domain  $\Gamma$  and it is the unique solution to  $h_1(\cdot) = h(\cdot, t)$ ,  $t \in \Gamma$ . For any  $\mathbf{u}$ ,  $h(\mathbf{u}, t)$  is second-order continuously differentiable in  $t$  and  $h$ ,  $\nabla_t h$ ,  $\nabla_t^2 h$  are bounded over  $t$  and  $\mathbf{u}$ . As  $t \rightarrow \gamma$ ,  $h(\cdot, t)$ ,  $\nabla_t h(\cdot, t)$ , and  $\nabla_t^2 h(\cdot, t)$  convergence uniformly.*

(D4)  *$\sup_{t \in \Gamma} E\|\mathbf{W}_t\|^4 < \infty$  and  $\sup_{t \in \Gamma} E[\|\mathbf{W}_t\|^4 | \mathbf{U}] f_U(\mathbf{U})$  is bounded, where  $\|\mathbf{a}\|^2 = \mathbf{a}^\top \mathbf{a}$ ,*

$$\mathbf{W}_t = (De^{tY}, (1-D)Ye^{-tY}, D, (1-D), DY e^{tY}, (1-D)Y^2 e^{-tY}, DY^2 e^{tY}, (1-D)Y^3 e^{-tY})^\top,$$

and  $f_U$  is the density of  $\mathbf{U}$ . Furthermore, there is a function  $\tau(Y, D)$  with  $E\{\tau(Y, D)\} < \infty$  such that  $\|\mathbf{W}_t - \mathbf{W}_{t'}\| < \tau(Y, D)|t - t'|$ .

(D5) The function  $\boldsymbol{\omega}_t(\mathbf{u}) = E(\mathbf{W}_t | \mathbf{U} = \mathbf{u})f_U(\mathbf{u})$  is bounded away from zero, and it is  $d$ th-order continuously differentiable with bounded derivatives on an open set containing the support of  $\mathbf{U}$ . There is a functional  $G(Y, D, \boldsymbol{\omega})$  linear in  $\boldsymbol{\omega}$  such that  $|G(Y, D, \boldsymbol{\omega})| \leq \iota(Y, D)\|\boldsymbol{\omega}\|_\infty$  and, for small enough  $\|\boldsymbol{\omega} - \boldsymbol{\omega}_\gamma\|_\infty$ ,  $|\psi(Y, D, \boldsymbol{\omega}) - \psi(Y, D, \boldsymbol{\omega}_\gamma) - G(Y, D, \boldsymbol{\omega} - \boldsymbol{\omega}_\gamma)| \leq \iota(Y, D)\|\boldsymbol{\omega} - \boldsymbol{\omega}_\gamma\|_\infty^2$ , where  $\iota(Y, D)$  is a function with  $E\{\iota(Y, D)\} < \infty$ ,  $\psi(Y, D, \boldsymbol{\omega}) = -2D \left( Y - \frac{\omega_1 \omega_2}{\omega_3 \omega_4} \right) \left( \frac{\omega_2 \omega_5 - \omega_1 \omega_6}{\omega_3 \omega_4} \right)$ ,  $\omega_j$  is the  $j$ th component of  $\boldsymbol{\omega}$ ,  $\|\boldsymbol{\omega}\|_\infty = \sup_{\mathbf{u} \in \mathcal{U}} \|\boldsymbol{\omega}(\mathbf{u})\|$ ,  $\|\boldsymbol{\omega} - \boldsymbol{\omega}_\gamma\|_\infty = \sup_{\mathbf{u} \in \mathcal{U}} \|\boldsymbol{\omega}(\mathbf{u}) - \boldsymbol{\omega}_\gamma(\mathbf{u})\|$ , and  $\mathcal{U}$  is the range of  $\mathbf{U}$ . Also, there exists an almost everywhere continuous 8-dimensional function  $\boldsymbol{\nu}(\mathbf{U})$  with  $\int \|\boldsymbol{\nu}(\mathbf{u})\| d\mathbf{u} < \infty$  and  $E\{\sup_{\|\boldsymbol{\delta}\| \leq \epsilon} \|\boldsymbol{\nu}(\mathbf{U} + \boldsymbol{\delta})\|^4\} < \infty$  for some  $\epsilon > 0$  such that  $E\{G(Y, D, \boldsymbol{\omega})\} = \int \boldsymbol{\nu}(\mathbf{u})^\top \boldsymbol{\omega}(\mathbf{u}) d\mathbf{u}$  for all  $\|\boldsymbol{\omega}\|_\infty < \infty$ .

Then, as the total sample size of internal and external datasets  $N \rightarrow \infty$ ,

$$\sqrt{N}(\hat{\gamma} - \gamma) \xrightarrow{d} N(0, \sigma_\gamma^2), \quad (20)$$

where  $\sigma_\gamma^2 = [2E\{D\nabla_\gamma h(\mathbf{U}, \gamma)\}^2]^{-1} \text{Var}[\psi(Y, D, \boldsymbol{\omega}_\gamma) + \boldsymbol{\nu}(\mathbf{U})^\top \mathbf{W}_\gamma - E\{\boldsymbol{\nu}(\mathbf{U})^\top \mathbf{W}_\gamma\}]$ .

Conditions (D1)-(D5) are technical assumptions discussed in Lemmas 8.11 and 8.12 in Newey and McFadden (1994). As discussed by Newey and McFadden (1994), the condition that  $\tilde{\kappa}$  has a bounded support can be relaxed, as it is imposed for a simple proof.

Combining Theorems 1-3, we obtain the following result for  $\hat{\mu}_1^{E3}(\mathbf{u})$  in (12) or  $\hat{\mu}_1^{C3}(\mathbf{u})$  in (16).

**Corollary 1.** Suppose that (8) holds for the binary random  $D$  indicating internal and ex-

ternal data.

(i) Under (B1)-(B4) and (D1)-(D5), result (17) holds with  $\hat{\mu}_1^{E2}(\mathbf{u})$  replaced by  $\hat{\mu}_1^{E3}(\mathbf{u})$  in (12).

(ii) Under (C1)-(C4) and (D1)-(D5) with  $\mathbf{U}$  and  $p$  replaced by  $\mathbf{X}$  and  $q$ , respectively, result (19) holds with  $\hat{\mu}_1^{C2}(\mathbf{u})$  replaced by  $\hat{\mu}_1^{C3}(\mathbf{u})$ .

## 4 Simulation Results

In this section, we present simulation results to examine and compare the performance of the standard kernel estimator (2) without using external information and our proposed estimator (16) with three variations,  $\hat{\mu}_1^{C1}$ ,  $\hat{\mu}_1^{C2}$ , and  $\hat{\mu}_1^{C3}$ , as described in the end of Section 2.2.

We consider  $\mathbf{U} = (X, Z)^\top$  with univariate covariate  $\mathbf{X} = X$  and unmeasured  $Z$  ( $p = 2$  and  $q = 1$ ) in two cases:

- (i) normal covariates:  $(X, Z)^\top$  is bivariate normal with means 0, variances 1, and correlation 0.5;
- (ii) bounded covariates:  $X = BW_1 + (1 - B)W_2$  and  $Z = BW_1 + (1 - B)W_3$ , where  $W_1$ ,  $W_2$ , and  $W_3$  are identically distributed as uniform on  $[-1, 1]$ ,  $B$  is uniform on  $[0, 1]$ , and  $W_1$ ,  $W_2$ ,  $W_3$ , and  $B$  are independent.

Conditioned on  $(X, Z)^\top$ , the response  $Y$  is normal with mean  $\mu(X, Z)$  and variance 1, where  $\mu(X, Z)$  follows one of the following four models:

- M1.  $\mu(X, Z) = X/2 - Z^2/4$ ;
- M2.  $\mu(X, Z) = \cos(2X)/2 + \sin(Z)$ ;
- M3.  $\mu(X, Z) = \cos(2XZ)/2 + \sin(Z)$ ;



M4.  $\mu(X, Z) = X/2 - Z^2/4 + \cos(XZ)/4$ .

Note that all four models are nonlinear in  $(X, Z)^\top$ ; M1-M2 are additive models, while M3-M4 are non-additive.

A total of  $N = 1,200$  data are generated from the population of  $(Y, X, Z)$  as previously described. A data point is treated as internal or external according to a random binary  $D$  with conditional probability  $P(D = 1 \mid Y, X, Z) = 1/\exp(1 + 2|X| + \gamma Y)$ , where  $\gamma = 0$  or  $1/2$ . Under this setting, the unconditional  $P(D = 1)$  is between 10% and 15%.

The simulation studies performance of kernel estimators in terms of mean integrated square error (MISE). The following measure is calculated by simulation with  $S$  replications:

$$\text{MISE} = \frac{1}{S} \sum_{s=1}^S \frac{1}{T} \sum_{t=1}^T \{\hat{\mu}_1^*(\mathbf{U}_{s,t}) - \mu_1(\mathbf{U}_{s,t})\}^2, \quad (21)$$

where  $\{\mathbf{U}_{s,t} : t = 1, \dots, T\}$  are test data for each simulation replication  $s$ , the simulation is repeated independently for  $s = 1, \dots, S$ , and  $\hat{\mu}_1^*(\cdot)$  is an estimator of  $\mu_1(\cdot)$  using a method described previously based on internal and external data, independent of test data. We consider two ways of generating test data  $\mathbf{U}_{s,t}$ 's. The first one is to use  $T = 121$  fixed grid points on  $[-1, 1] \times [-1, 1]$  with equal space. The second one is to take a random sample of  $T = 121$  without replacement from the covariate  $\mathbf{U}$ 's of the internal dataset, for each fixed  $s = 1, \dots, S$  and independently across  $s$ .

To show the benefit of using external information, we calculate the improvement in efficiency defined as follows:

$$\text{IMP} = 1 - \frac{\min\{\text{MISE}(\hat{\mu}_1^{Cj}) \text{ over } j = 1, 2, 3\}}{\text{MISE}(\hat{\mu}_1)}. \quad (22)$$

In all cases, we use the Gaussian kernel. The bandwidths  $b$  and  $l$  affect the performance of kernel methods. We consider two types of bandwidths in the simulation. The first one is “the best bandwidth”; for each method, we evaluate MISE in a pool of bandwidths and display the one that has the minimal MISE. This shows the best we can achieve in terms of bandwidth, but it cannot be used in applications. The second one is to select bandwidth from a pool of bandwidths via 10-fold cross validation (Györfi et al., 2002), which produces a decent bandwidth that can be applied to real data.

The simulated MISE based on  $S = 200$  replications is shown in Table 1 for the case of  $\gamma = 0$  and in Table 2 for the case of  $\gamma = 1/2$ .

Consider first the results in Table 1. Since  $\gamma = 0$ , all three estimators,  $\hat{\mu}_1^{C1}$ ,  $\hat{\mu}_1^{C2}$ , and  $\hat{\mu}_1^{C3}$ , are correct and more efficient than the standard estimator  $\hat{\mu}_1$  in (2) without using external information. The estimator  $\hat{\mu}_1^{C1}$  is the best, as it uses the correct information that populations are homogeneous ( $\gamma = 0$ ).

Next, the results in Table 2 for  $\gamma = 1/2$  indicate that the estimator  $\hat{\mu}_1^{C2}$  or  $\hat{\mu}_1^{C3}$  using a correct constraint is better than the estimator  $\hat{\mu}_1^{C1}$  using an incorrect constraint or the estimator  $\hat{\mu}_1$  without using external information. Since  $\hat{\mu}_1^{C3}$  uses more information, it is in general better than  $\hat{\mu}_1^{C2}$ . With an incorrect constraint,  $\hat{\mu}_1^{C1}$  can be much worse than  $\hat{\mu}_1$  without using external information.

## 5 Discussion

Curse of dimensionality is a well-known problem for nonparametric methods. Thus, the proposed method in Section 2 is intended for low dimensional covariate  $\mathbf{U}$ , i.e.,  $p$  is small. If  $p$  is not small, then we should reduce the dimension of  $\mathbf{U}$  prior to applying the CK, or

any kernel methods. For example, consider a single index model assumption (Li, 1991), i.e.,  $\mu_1(\mathbf{U})$  in (1) is assumed to be

$$\mu_1(\mathbf{U}) = \mu_1(\boldsymbol{\eta}^\top \mathbf{U}), \quad (23)$$

where  $\boldsymbol{\eta}$  is an unknown  $p$ -dimensional vector. The well-known SIR technique (Li, 1991) can be applied to obtain a consistent and asymptotically normal estimator  $\hat{\boldsymbol{\eta}}$  of  $\boldsymbol{\eta}$  in (23). Once  $\boldsymbol{\eta}$  is replaced by  $\hat{\boldsymbol{\eta}}$ , the kernel method can be applied with  $\mathbf{U}$  replaced by the one-dimensional “covariate”  $\hat{\boldsymbol{\eta}}^\top \mathbf{U}$ . We can also apply other dimension reduction techniques developed under assumptions weaker than (23) (Cook and Weisberg, 1991; Li and Wang, 2007; Shao et al., 2007; Xia et al., 2002; Ma and Zhu, 2012).

We turn to the dimension of  $\mathbf{X}$  in the external dataset. When the dimension of  $\mathbf{X}$  is high, we may consider the following approach. Instead of using constraint (15), we use component-wise constraints

$$\sum_{i=1}^n \{\mu_i - \hat{h}_1^{(k)}(X_i^{(k)})\} \mathbf{g}_k(X_i^{(k)})^\top = 0, \quad k = 1, \dots, q, \quad (24)$$

where  $X_i^{(k)}$  is the  $k$ th component of  $\mathbf{X}_i$ ,  $\mathbf{g}_k(X^{(k)}) = (1, X^{(k)})^\top$ , and  $\hat{h}_1^{(k)}(X_i^{(k)})$  is an estimator of  $h_1^{(k)}(X^{(k)}) = E(Y \mid X^{(k)}, D = 1)$  using methods described in Section 2. More constraints are involved in (24), but estimation only involves one dimensional  $X^{(k)}$ ,  $k = 1, \dots, q$ .

The kernel  $\kappa$  we adopted in (2) and (16) is the second order kernel so that the convergence rate of  $\hat{\mu}_1^E(\mathbf{u}) - \mu_1(\mathbf{u})$  is  $n^{-2/(4+p)}$ . An  $m$ th order kernel with  $m > 2$  as defined by Bierens (1987) may be used to achieve convergence rate  $n^{-m/(2m+p)}$ . Alternatively, we may apply other nonparametric smoothing techniques such as the local polynomial (Fan et al., 1997) to achieve convergence rate  $n^{-m/(2m+p)}$  with  $m \geq 2$ .

Our results can be extended to the scenarios where several external datasets are available. Since each external source may provide different covariate variables, we may need to apply component-wise constraints (24) by estimating  $\widehat{h}_1^{(k)}$  via combining all the external sources that collects covariate  $X^{(k)}$ . If populations of external datasets are different, then we may have to apply a combination of the methods described in Section 2.

## Appendix

**Proof of Theorem 1.** Let  $\tilde{\mu}_1(\mathbf{u}) = \widehat{p}(\mathbf{u})\widehat{\mu}_1(\mathbf{u}) + \{1 - \widehat{p}(\mathbf{u})\}\widehat{\mu}_0(\mathbf{u})$ , where  $\widehat{\mu}_1(\mathbf{u}) = \sum_{i=1}^n \kappa_b(\mathbf{u} - \mathbf{U}_i)Y_i / \sum_{i=1}^n \kappa_b(\mathbf{u} - \mathbf{U}_i)$ ,  $\widehat{\mu}_0(\mathbf{u}) = \sum_{i=n+1}^N \kappa_b(\mathbf{u} - \mathbf{U}_i)\tilde{Y}_i / \sum_{i=n+1}^N \kappa_b(\mathbf{u} - \mathbf{U}_i)$ , and  $\widehat{p}(\mathbf{u}) = \sum_{i=1}^n \kappa_b(\mathbf{u} - \mathbf{U}_i) / \sum_{i=1}^N \kappa_b(\mathbf{u} - \mathbf{U}_i)$ . Under (B3)-(B4), Theorem 2 in Nadaraya (1964) shows that  $\widehat{p}(\mathbf{u})$  converges to  $P(D = 1 | \mathbf{U} = \mathbf{u})$  in probability. Under (B1)-(B4),  $\sqrt{nb^p}\{\widehat{\mu}_1(\mathbf{u}) - \mu_1(\mathbf{u})\} \xrightarrow{d} N(B_1(\mathbf{u}), V_1(\mathbf{u}))$ ,  $B_1(\mathbf{u}) = c^{1/2}A_1(\mathbf{u})$ ,  $V_1(\mathbf{u}) = \frac{\sigma_1^2(\mathbf{u})}{f_1(\mathbf{u})} \int \kappa(\mathbf{v})^2 d\mathbf{v}$ , and  $\sqrt{n/(N-n)}\sqrt{(N-n)b^p}\{\widehat{\mu}_0(\mathbf{u}) - \mu_1(\mathbf{u})\} \xrightarrow{d} N(B_0(\mathbf{u}), V_0(\mathbf{u}))$ ,  $B_0(\mathbf{u}) = c^{1/2}A_0(\mathbf{u})$ ,  $V_0(\mathbf{u}) = \frac{\sigma_0^2(\mathbf{u})}{af_0(\mathbf{u})} \int \kappa(\mathbf{v})^2 d\mathbf{v}$ . Then (17) holds for  $\tilde{\mu}_1(\mathbf{u})$ , by Slutsky's theorem, the independence between  $\widehat{\mu}_1$  and  $\widehat{\mu}_0$ , and the definition of  $a$ . The desired result (17) follows from the fact that  $|\widehat{\mu}_1^{E2}(\mathbf{u}) - \tilde{\mu}_1(\mathbf{u})|$  is bounded by

$$\{1 - \widehat{p}(\mathbf{u})\} \max_{i > n} \left| \frac{\widehat{f}(Y_i | \mathbf{U} = \mathbf{U}_i, D = 1)}{\widehat{f}(Y_i | \mathbf{U} = \mathbf{U}_i, D = 0)} - \frac{f(Y_i | \mathbf{U} = \mathbf{U}_i, D = 1)}{f(Y_i | \mathbf{U} = \mathbf{U}_i, D = 0)} \right| \frac{\sum_{i=n+1}^N |Y_i| \kappa_b(\mathbf{u} - \mathbf{U}_i)}{\sum_{i=n+1}^N \kappa_b(\mathbf{u} - \mathbf{U}_i)}, \quad (25)$$

which is  $o_p(1/\sqrt{nb^p})$  by result (18) under condition (B5).

**Proof of Theorem 2.** Write

$$\sqrt{nb^p}\{\widehat{\mu}_1^{C2}(\mathbf{u}) - \mu_1(\mathbf{u})\} = T_1 + \cdots + T_6, \quad (26)$$

where  $T_1 = n^{-1/2}b^{p/2}\boldsymbol{\delta}_b(\mathbf{u})^\top (\mathbf{I}_n - \mathbf{P})\mathbf{B}_l^{-1}\boldsymbol{\Delta}_l\boldsymbol{\epsilon}/\widehat{f}_b(\mathbf{u})$ ,  $T_2 = n^{-1/2}b^{p/2}\boldsymbol{\delta}_b(\mathbf{u})^\top \{\boldsymbol{\mu}_1 - \mu_1(\mathbf{u})\mathbf{1}_n\}/\widehat{f}_b(\mathbf{u})$ ,

$T_3 = n^{-1/2} b^{p/2} \boldsymbol{\delta}_b(\mathbf{u})^\top (\mathbf{B}_l^{-1} \boldsymbol{\Delta}_l \boldsymbol{\mu}_1 - \boldsymbol{\mu}_1) / \widehat{f}_b(\mathbf{u})$ ,  $T_4 = -n^{-1/2} b^{p/2} \boldsymbol{\delta}_b(\mathbf{u})^\top \mathbf{P}(\mathbf{B}_l^{-1} \boldsymbol{\Delta}_l \boldsymbol{\mu}_1 - \boldsymbol{\mu}_1) / \widehat{f}_b(\mathbf{u})$ ,  
 $T_5 = n^{-1/2} b^{p/2} \boldsymbol{\delta}_b(\mathbf{u})^\top \mathbf{P}(\widehat{\mathbf{h}}_1 - \mathbf{h}_1) / \widehat{f}_b(\mathbf{u})$ ,  $T_6 = n^{-1/2} b^{p/2} \boldsymbol{\delta}_b(\mathbf{u})^\top \mathbf{P}(\mathbf{h}_1 - \boldsymbol{\mu}_1) / \widehat{f}_b(\mathbf{u})$ ,  $\widehat{f}_b(\mathbf{u}) = \sum_{i=1}^n \kappa_b(\mathbf{u} - \mathbf{U}_i) / n$ ,  $\boldsymbol{\delta}_b(\mathbf{u}) = (\kappa_b(\mathbf{u} - \mathbf{U}_1), \dots, \kappa_b(\mathbf{u} - \mathbf{U}_n))^\top$ ,  $\mathbf{I}_n$  is the identity matrix of order  $n$ ,  $\mathbf{1}_n$  is the  $n$ -vector with all components being 1,  $\mathbf{B}_l$  is the  $n \times n$  diagonal matrix whose  $i$ th diagonal element is  $\widehat{f}_l(\mathbf{U}_i)$ ,  $\boldsymbol{\Delta}_l$  is the  $n \times n$  matrix whose  $(i, j)$ th entry is  $\kappa_l(\mathbf{U}_i - \mathbf{U}_j) / n$ ,  $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^\top$  with  $\epsilon_i = Y_i - \mu_1(\mathbf{U}_i)$ ,  $\mathbf{h}_1$  is the  $n$ -dimensional vector whose  $i$ th component is  $h_1(\mathbf{X}_i)$ ,  $\mathbf{P} = \mathbf{G}(\mathbf{G}^\top \mathbf{G})^{-1} \mathbf{G}^\top$ , and  $\mathbf{G}$ ,  $\widehat{\mathbf{h}}_1$ , and  $\boldsymbol{\mu}_1$  are defined in Section 2.

We first show that  $T_1$  in (26) is asymptotically normal with mean 0 and variance  $V_r(\mathbf{u})$  defined in Theorem 2. Consider a further decomposition  $T_1 = \sqrt{n} V + T_{11} + T_{12} + T_{13}$ , where

$$V = \frac{1}{n^2} \sum_{j=1}^n \sum_{i=1}^n S(\mathbf{U}_i, \epsilon_i, \mathbf{U}_j, \epsilon_j)$$

is a V-statistic with

$$S(\mathbf{U}_i, \epsilon_i, \mathbf{U}_j, \epsilon_j) = \frac{b^{p/2}}{2f_1(\mathbf{u})} \left\{ \frac{\kappa_b(\mathbf{u} - \mathbf{U}_i) \kappa_l(\mathbf{U}_i - \mathbf{U}_j) \epsilon_j}{f_1(\mathbf{U}_i)} + \frac{\kappa_b(\mathbf{u} - \mathbf{U}_j) \kappa_l(\mathbf{U}_j - \mathbf{U}_i) \epsilon_i}{f_1(\mathbf{U}_j)} \right\},$$

$$T_{11} = \frac{b^{p/2}}{n^{3/2}} \sum_{i=1}^n \frac{\kappa_b(\mathbf{u} - \mathbf{U}_i) \kappa_l(0) \epsilon_i}{f_1(\mathbf{u}) f_1(\mathbf{U}_i)},$$

$$T_{12} = \frac{b^{p/2}}{n^{3/2}} \sum_{j=1}^n \sum_{i=1}^n \frac{\kappa_b(\mathbf{u} - \mathbf{U}_i) \kappa_l(\mathbf{U}_i - \mathbf{U}_j)}{f_1(\mathbf{u}) f_1(\mathbf{U}_i)} \left\{ \frac{f_1(\mathbf{u}) f_1(\mathbf{U}_i)}{\widehat{f}_b(\mathbf{u}) \widehat{f}_l(\mathbf{U}_i)} - 1 \right\} \epsilon_j,$$

and  $T_{13} = -n^{-1/2} b^{p/2} \boldsymbol{\delta}_b(\mathbf{u})^\top \mathbf{P} \mathbf{B}_l^{-1} \boldsymbol{\Delta}_l \boldsymbol{\epsilon} / \widehat{f}_b(\mathbf{u})$ . Note that

$$S_1(\mathbf{U}_1, \epsilon_1) = E\{S(\mathbf{U}_1, \epsilon_1, \mathbf{U}_2, \epsilon_2) \mid \mathbf{U}_1, \epsilon_1\} = \frac{b^{p/2}}{2f_1(\mathbf{u})} \left\{ \int \kappa_l(\mathbf{u}_2 - \mathbf{U}_1) \kappa_b(\mathbf{u} - \mathbf{u}_2) d\mathbf{u}_2 \right\} \epsilon_1$$

having variance

$$\begin{aligned}
 \text{Var}\{S_1(\mathbf{U}_1, \epsilon_1)\} &= \frac{b^{p/2}}{4f_1^2(\mathbf{u})} \int f_1(\mathbf{u}_1) \sigma_1^2(\mathbf{u}_1) \left\{ \int \kappa_l(\mathbf{u}_2 - \mathbf{u}_1) \kappa_b(\mathbf{u} - \mathbf{u}_2) d\mathbf{u}_2 \right\}^2 d\mathbf{u}_1 \\
 &= \frac{b^{p/2}}{4f_1^2(\mathbf{u})} \int f_1(\mathbf{u}_1) \sigma_1^2(\mathbf{u}_1) \left\{ \int \kappa_l(\mathbf{v}) \kappa_b(\mathbf{u} - \mathbf{u}_1 - \mathbf{lv}) d\mathbf{v} \right\}^2 d\mathbf{u}_1
 \end{aligned}$$

$$= \frac{1}{4f_1^2(\mathbf{u})} \int f_1(\mathbf{u} - b\mathbf{w}) \sigma_1^2(\mathbf{u} - b\mathbf{w}) \left\{ \int \kappa(\mathbf{v}) \kappa\left(\mathbf{w} - \nu \frac{l}{b}\right) d\nu \right\}^2 d\mathbf{w},$$

where  $\sigma_1^2(\cdot)$  is given in condition (C2), the second and third equalities follow from changing variables  $\mathbf{u}_2 - \mathbf{u}_1 = l\nu$  and  $\mathbf{u} - \mathbf{u}_1 = b\mathbf{w}$ , respectively. From the continuity of  $f_1(\cdot)$  and  $\sigma_1^2(\cdot)$ ,  $\text{Var}\{S_1(\mathbf{u}_1, \epsilon_1)\}$  converges to  $V_r(\mathbf{u})$ . Therefore, by the theory for asymptotic normality of V-statistics (e.g., Theorem 3.16 in Shao (2003)),  $\sqrt{n}V \xrightarrow{d} N(0, V_r(\mathbf{u}))$ .

Conditioned on  $\mathbf{U}_1, \dots, \mathbf{U}_n$ ,  $T_{11}$  has mean 0 and variance

$$\begin{aligned} \text{Var}(T_{11} | \mathbf{U}_1, \dots, \mathbf{U}_n) &= \frac{b^p}{4f_1^2(\mathbf{u})n^3} \sum_{i=1}^n \frac{\kappa_b(\mathbf{u} - \mathbf{U}_i)^2 \kappa_l(0)^2 \sigma_1^2(\mathbf{U}_i)}{f_1(\mathbf{U}_i)} \\ &\leq \frac{\sup_{\mathbf{u} \in \mathcal{U}} \kappa(\mathbf{u})^3}{4f_1^2(\mathbf{u})n^3 b^{2p}} \sum_{i=1}^n \frac{\kappa_b(\mathbf{u} - \mathbf{U}_i) \sigma_1^2(\mathbf{U}_i)}{f_1(\mathbf{U}_i)} = o_p(1). \end{aligned}$$

This proves that  $T_{11} = o_p(1)$ . Note that  $E(T_{12} | \mathbf{U}_1, \dots, \mathbf{U}_n) = 0$  and  $\text{Var}(T_{12} | \mathbf{U}_1, \dots, \mathbf{U}_n)$  is bounded by

$$\max \left\{ \frac{1}{f_1^2(\mathbf{u})}, \frac{1}{\widehat{f}_b^2(\mathbf{u})} \right\} \max_{i=1, \dots, n} \left| \frac{f_1(\mathbf{U}_i)}{\widehat{f}_l(\mathbf{U}_i)} - 1 \right|^2 \text{Var}(\sqrt{n}V + T_{11} | \mathbf{U}_1, \dots, \mathbf{U}_n).$$

Therefore, under the assumed condition that  $f_1$  is bounded away from zero, Lemma 3 in Dai and Shao (2022) implies  $T_{12} = o_p(1)$ . Note that

$$T_{13} = \frac{b^{p/2}}{n^{1/2}} \sum_{j=1}^n W_j(\mathbf{u}) \epsilon_j, \quad W_j(\mathbf{u}) = \frac{1}{n} \sum_{i=1}^n \frac{\kappa_b(\mathbf{u} - \mathbf{U}_i) \mathbf{g}(\mathbf{X}_i)^\top}{\widehat{f}_b(\mathbf{u})} (\mathbf{G}^\top \mathbf{G})^{-1} \sum_{i=1}^n \frac{\kappa_l(\mathbf{U}_i - \mathbf{U}_j) \mathbf{g}(\mathbf{X}_i)}{\widehat{f}_l(\mathbf{U}_i)}.$$

Conditioned on  $\mathbf{U}_1, \dots, \mathbf{U}_n$ ,  $T_{13}$  has mean 0 and variance

$$\text{Var}(T_{13} | \mathbf{U}_1, \dots, \mathbf{U}_n) = \frac{b^p}{n} \sum_{j=1}^n W_j^2(\mathbf{u}) \sigma_1^2(\mathbf{U}_j) = O_p(b^p) = o_p(1),$$

because, under the assumed condition that  $f_1$  is bounded away from zero, Lemma 3 in Dai

and Shao (2022) implies  $\max_{j=1,\dots,n} |W_j(\mathbf{u}) - \mathbf{g}(\mathbf{u})^\top \boldsymbol{\Sigma}_g^{-1} \mathbf{g}(\mathbf{X}_j)| = o_p(1)$ . Thus,  $T_{13} = o_p(1)$ .

Consequently,  $T_1$  has the same asymptotic distribution as  $\sqrt{n}V$ , the claimed result.

From Lemma 4 in Dai and Shao (2022) and (C4),  $T_2 = \sqrt{c}A_1(\mathbf{u})\{1 + o_p(1)\}$ . Note that

$$\begin{aligned} T_3 &= \frac{\sqrt{nb^p}l^2}{n\widehat{f}_b(\mathbf{u})} \sum_{j=1}^n \kappa_b(\mathbf{u} - \mathbf{U}_j) \left[ \frac{1}{nl^2\widehat{f}_b(\mathbf{U}_j)} \sum_{i=1}^n \kappa_l(\mathbf{u} - \mathbf{U}_i) \{\mu_1(\mathbf{U}_i) - \mu_1(\mathbf{U}_j)\} \right] \\ &= \left\{ \frac{\sqrt{cr^2}}{n\widehat{f}_b(\mathbf{u})} \sum_{j=1}^n \kappa_b(\mathbf{u} - \mathbf{U}_j) A_1(\mathbf{U}_j) \right\} \{1 + o_p(1)\} = \sqrt{cr^2} A_1(\mathbf{u}) \{1 + o_p(1)\}, \end{aligned}$$

where the second equality follows from (A4) and Lemmas 3-4 in Dai and Shao (2022), and the last equality follows from Lemma 2 in Dai and Shao (2022) and continuity of  $A_1(\cdot)$ . Also,

$$\begin{aligned} -\frac{n^{1/2}T_4}{b^{p/2}} &= \frac{1}{n} \sum_{i=1}^n \frac{\kappa_b(\mathbf{u} - \mathbf{U}_i) \mathbf{g}(\mathbf{X}_i)^\top}{\widehat{f}_b(\mathbf{u})} (\mathbf{G}^\top \mathbf{G})^{-1} \sum_{j=1}^n \frac{\mathbf{g}(\mathbf{X}_j)}{n\widehat{f}_b(\mathbf{U}_j)} \sum_{i=1}^n \kappa_l(\mathbf{u} - \mathbf{U}_i) \{\mu_1(\mathbf{U}_i) - \mu_1(\mathbf{U}_j)\} \\ &= \left\{ \mathbf{g}(\mathbf{x})^\top \boldsymbol{\Sigma}_g^{-1} \frac{1}{n} \sum_{j=1}^n \frac{\mathbf{g}(\mathbf{X}_j)}{n\widehat{f}_b(\mathbf{U}_j)} \sum_{i=1}^n \kappa_l(\mathbf{u} - \mathbf{U}_i) \{\mu_1(\mathbf{U}_i) - \mu_1(\mathbf{U}_j)\} \right\} \{1 + o_p(1)\} \\ &= \left\{ \mathbf{g}(\mathbf{x})^\top \boldsymbol{\Sigma}_g^{-1} \frac{l^{2/p}}{n} \sum_{j=1}^n \mathbf{g}(\mathbf{X}_j) A_1(\mathbf{U}_j) \right\} \{1 + o_p(1)\} \\ &= \sqrt{cr^2} \mathbf{g}(\mathbf{x})^\top \boldsymbol{\Sigma}_g^{-1} E\{\mathbf{g}(\mathbf{X}) A_1(\mathbf{U})\} \{1 + o_p(1)\}, \end{aligned}$$

where the first equality follows from Lemma 3 in Dai and Shao (2022) and the law of large numbers, the second equality follows from Lemma 4 in Dai and Shao (2022), and the last equality follows from the law of large numbers. Similarly, is equal to

$$\begin{aligned} \frac{n^{1/2}T_5}{b^{p/2}} &= \frac{1}{n} \sum_{i=1}^n \frac{\kappa_b(\mathbf{u} - \mathbf{U}_i) \mathbf{g}(\mathbf{X}_i)^\top}{\widehat{f}_b(\mathbf{u})} (\mathbf{G}^\top \mathbf{G})^{-1} \sum_{i=1}^n \mathbf{g}(\mathbf{X}_i) \{\widehat{h}_1(\mathbf{X}_i) - h_1(\mathbf{X}_i)\} \\ &= \left[ \mathbf{g}(\mathbf{x})^\top \boldsymbol{\Sigma}_g^{-1} \frac{1}{n} \sum_{i=1}^n \mathbf{g}(\mathbf{X}_i) \{\widehat{h}_1(\mathbf{X}_i) - h_1(\mathbf{X}_i)\} \right] \{1 + o_p(1)\} \end{aligned}$$

$$\leq \{1 + o_p(1)\} O_p(1) \max_{j=1, \dots, n} |\widehat{h}_1(\mathbf{X}_j) - h_1(\mathbf{X}_j)|$$

the equality follows from Lemma 3 in Dai and Shao (2022). Under (B1)-(B5) with  $\mathbf{U}$  and  $p$  replaced by  $\mathbf{X}$  and  $q$ , respectively, and (C5), Lemma 8.10 in Newey and McFadden (1994) implies that

$$\max |\widehat{h}_1(\mathbf{X}_i) - h_1(\mathbf{X}_1)| = O_p(\sqrt{\log(n)} n^{-2/(q+4)}), \quad (27)$$

which is  $o_p(1/\sqrt{nb^p}) = o_p(n^{-2/(p+4)})$  and, hence,  $T_5 = o_p(1)$ . From Lemma 3 in Dai and Shao (2022) and the Central Limit Theorem,

$$T_6 = \frac{b^{p/2}}{n^{1/2}} \sum_{i=1}^n \frac{\kappa_b(\mathbf{u} - \mathbf{U}_i) \mathbf{g}(\mathbf{X}_i)^\top}{\widehat{f}_b(\mathbf{u})} (\mathbf{G}^\top \mathbf{G})^{-1} \sum_{i=1}^n \mathbf{g}(\mathbf{X}_i) \{h_1(\mathbf{X}_i) - \mu_1(\mathbf{U}_i)\} = O_p(b^{p/2}) = o_p(1).$$

Combining these results, we obtain that  $T_2 + \dots + T_6 = B_r(\mathbf{u}) + o_p(1)$ . This completes the proof.

### Proof of Theorem 3.

**I change  $\mathbf{X}$  to  $\mathbf{U}$ . – Chi**

Define

$$\begin{aligned} \widehat{\omega}_{t1} &= \frac{1}{N} \sum_{i=1}^N R_i \check{\kappa}_b(\mathbf{u} - \mathbf{U}_i) e^{tY_i}, & \widehat{\omega}_{t2} &= \frac{1}{N} \sum_{i=1}^N (1 - R_i) \check{\kappa}_b(\mathbf{u} - \mathbf{U}_i) Y_i e^{-tY_i}, \\ \widehat{\omega}_{t3} &= \frac{1}{N} \sum_{i=1}^N R_i \check{\kappa}_b(\mathbf{u} - \mathbf{U}_i), & \widehat{\omega}_{t4} &= \frac{1}{N} \sum_{i=1}^N (1 - R_i) \check{\kappa}_b(\mathbf{u} - \mathbf{U}_i), \\ \widehat{\omega}_{t5} &= \frac{1}{N} \sum_{i=1}^N R_i \check{\kappa}_b(\mathbf{u} - \mathbf{U}_i) Y_i e^{tY_i}, & \widehat{\omega}_{t6} &= \frac{1}{N} \sum_{i=1}^N (1 - R_i) \check{\kappa}_b(\mathbf{u} - \mathbf{U}_i) Y_i^2 e^{-tY_i}, \\ \widehat{\omega}_{t7} &= \frac{1}{N} \sum_{i=1}^N R_i \check{\kappa}_b(\mathbf{u} - \mathbf{U}_i) Y_i^2 e^{tY_i}, & \widehat{\omega}_{t8} &= \frac{1}{N} \sum_{i=1}^N (1 - R_i) \check{\kappa}_b(\mathbf{u} - \mathbf{U}_i) Y_i^3 e^{-tY_i}. \end{aligned}$$



Then,  $\hat{h}(\mathbf{u}, t) = \hat{\omega}_{t1}\hat{\omega}_{t2}/\hat{\omega}_{t3}\hat{\omega}_{t4}$ ,  $\nabla_t \hat{h}(\mathbf{u}, t) = (\hat{\omega}_{t2}\hat{\omega}_{t5} - \hat{\omega}_{t1}\hat{\omega}_{t6})/\hat{\omega}_{t3}\hat{\omega}_{t4}$ , and  $\nabla_t^2 \hat{h}(\mathbf{u}, t) = (\hat{\omega}_{t1}\hat{\omega}_{t8} - 2\hat{\omega}_{t5}\hat{\omega}_{t6} + \hat{\omega}_{t2}\hat{\omega}_{t7})/\hat{\omega}_{t3}\hat{\omega}_{t4}$ . Let  $L(t) = E[R\{Y - h(\mathbf{U}, t)\}^2]$ ,  $\hat{L}_n(t) = N^{-1} \sum_{i=1}^N R_i\{Y_i - \hat{h}(\mathbf{U}_i, t)\}^2$ , and  $L_n(t) = N^{-1} \sum_{i=1}^N R_i\{Y_i - h(\mathbf{U}_i, t)\}^2$ . Taking derivatives with respect to  $t$ , we obtain

$$\nabla_t \hat{L}_n(t) = \frac{1}{N} \sum_{i=1}^N -2R_i\{Y_i - \hat{h}(\mathbf{U}_i, t)\} \nabla_t \hat{h}(\mathbf{U}_i, t) = \frac{1}{N} \sum_{i=1}^N \psi\{Y_i, R_i, \hat{\omega}_t(\mathbf{U}_i)\},$$

$$\nabla_t L_n(t) = \frac{1}{N} \sum_{i=1}^N -2R_i\{Y_i - h(\mathbf{U}_i, t)\} \nabla_t h(\mathbf{U}_i, t) = \frac{1}{N} \sum_{i=1}^N \psi\{Y_i, R_i, \omega_t(\mathbf{U}_i)\},$$

and

$$\nabla_t L(t) = -2E[R\{Y - h(\mathbf{u}, t)\} \nabla_t h(\mathbf{u}, t)] = E[\psi\{Y, R, \omega_t(\mathbf{U})\}],$$

where  $\psi$  is given in (D5). Note that  $\nabla_t L(\gamma) = 0$  and  $\nabla_t^2 L(\gamma) = 2E[\{\nabla_t h(\mathbf{U}, \gamma)\}^2 R] = \nu_\gamma \geq 0$ .

We establish the asymptotic normality of  $\hat{\gamma}$  in the following four steps.

**Step 1:** Since  $\gamma$  is the unique minimizer of  $L(t)$ , from Theorem 2.1 in Newey and McFadden (1994), it suffices to prove that  $\sup_{t \in \Gamma} |\nabla_t \hat{L}_n(t) - \nabla_t L(t)| \xrightarrow{p} 0$ . Note that

$$\begin{aligned} \sup_{t \in \Gamma} |\nabla_t \hat{L}_n(t) - \nabla_t L(t)| &\leq \sup_{t \in \Gamma} |\nabla_t L_n(t) - \nabla_t L(t)| + \sup_{t \in \Gamma} |\nabla_t \hat{L}_n(t) - \nabla_t L_n(t)| \\ &\leq \sup_{t \in \Gamma} |\nabla_t L_n(t) - \nabla_t L(t)| \\ &\quad + \frac{2}{n} \sum_{i=1}^n R_i |Y_i| \left\{ \sup_{t \in \Gamma, \mathbf{u} \in \mathcal{U}} |\nabla_t \hat{h}(\mathbf{u}, t) - \nabla_t h(\mathbf{u}, t)| \right. \\ &\quad \left. + \sup_{t \in \Gamma, \mathbf{u} \in \mathcal{U}} |\hat{h}(\mathbf{u}, t) \nabla_t \hat{h}(\mathbf{u}, t) - h(\mathbf{u}, t) \nabla_t h(\mathbf{u}, t)| \right\} \end{aligned}$$

From (D3),  $|2R\{Y - h(\mathbf{u}, t)\} \nabla_t h(\mathbf{u}, t)|$  is bounded by  $c|Y|$  for a constant  $c$  and hence

Lemma 2.4 in Newey and McFadden (1994) implies that  $\sup_{t \in \Gamma} |\nabla_t L_n(t) - \nabla_t L(t)| = o_p(1)$ .

Based on Lemma B.3 in Newey (1994), conditions (D1)-(D4) imply that  $\sup_{\mathbf{u} \in \mathcal{U}} |\hat{\omega}_t(\mathbf{u}) - \omega_t(\mathbf{u})| \rightarrow 0$  for all  $t \in \Gamma$ . As a result, by a similar argument of the proof of Lemma B.3 in

Newey (1994), we obtain that  $\sup_{t \in \Gamma, \mathbf{u} \in \mathcal{U}} |\hat{\omega}_t(\mathbf{u}) - \omega_t(\mathbf{u})| \xrightarrow{p} 0$ . Since  $\omega_t$  is bounded away from zero and  $h(\cdot, t)$  and  $\nabla_t h(\cdot, t)$  are Lipschitz continuous functions with respect to  $\omega_t$ ,  $\sup_{t \in \Gamma, \mathbf{u} \in \mathcal{U}} |\hat{h}(\mathbf{u}, t) - h(\mathbf{u}, t)| \xrightarrow{p} 0$  and  $\sup_{t \in \Gamma, \mathbf{u} \in \mathcal{U}} |\nabla_t \hat{h}(\mathbf{u}, t) - \nabla_t h(\mathbf{u}, t)| \xrightarrow{p} 0$ . These results together with the previous inequality implies that  $\hat{\gamma} \xrightarrow{p} \gamma$ .

**Step 2:** Conditions (D1)-(D5) ensure that Lemma 8.11 in Newey and McFadden (1994)

holds and hence  $\sqrt{N} \nabla_t \hat{L}_n(\gamma) \xrightarrow{d} N(0, \sigma_L^2)$  with  $\sigma_L^2 = \text{Var}\{m(Y, R, \mathbf{U}, \omega_\gamma) + \tau(Y, R, \mathbf{U}, \gamma)\}$ .

**Step 3:** Note that  $\nabla_t^2 L_n(t) = N^{-1} \sum_{i=1}^N -2R_i\{Y_i - h(\mathbf{U}_i, t)\} \nabla_t^2 h(\mathbf{U}_i, t) + 2R_i\{\nabla_t h(\mathbf{U}_i, t)\}^2$

and  $\sup_{|t-\gamma| \leq |\hat{\gamma}-\gamma|} |\nabla_t^2 \hat{L}_n(t) - \nabla_t^2 L(\gamma)| \leq A_1 + A_2 + A_3$ , where  $A_1 = |\nabla_t^2 L_n(\gamma) - \nabla_t^2 L(\gamma)|$ ,

$A_2 = \sup_{t \in \Gamma} |\nabla_t^2 \hat{L}_n(t) - \nabla_t^2 L_n(t)|$ , and the last term  $A_3 = \sup_{|t-\gamma| \leq |\hat{\gamma}-\gamma|} |\nabla_t^2 L_n(t) - \nabla_t^2 L_n(\gamma)|$ .

The law of large numbers guarantees that  $A_1 = o_p(1)$ . A similar argument in Step 1 shows that  $A_2 = o_p(1)$ . For  $A_3$ , we have

$$\begin{aligned} |\nabla_t^2 L_n(t) - \nabla_t^2 L_n(\gamma)| &\leq \frac{2}{N} \sum_{i=1}^N |\{\nabla_t h(\mathbf{U}_i, t)\}^2 - \{\nabla_t h(\mathbf{U}_i, \gamma)\}^2| \\ &\quad + \frac{2}{N} \sum_{i=1}^N |Y_i| |\nabla_t^2 h(\mathbf{U}_i, t) - \nabla_t^2 h(\mathbf{U}_i, \gamma)| \\ &\quad + \frac{2}{N} \sum_{i=1}^N |h(\mathbf{U}_i, t) \nabla_t^2 h(\mathbf{U}_i, t) - h(\mathbf{U}_i, \gamma) \nabla_t^2 h(\mathbf{U}_i, \gamma)|. \end{aligned}$$

Under (D3),  $h(\cdot, t)$ ,  $\nabla h(\cdot, t)$ , and  $\nabla^2 h(\cdot, t)$  converge uniformly for all  $\mathbf{x}$  as  $t \rightarrow \gamma$  and, thus,

the  $A_3 = o_p(1)$  because  $\hat{\gamma} \xrightarrow{p} \gamma$ . This shows that  $\sup_{|t-\gamma| \leq |\hat{\gamma}-\gamma|} |\nabla_t^2 \hat{L}_n(t) - \nabla_t^2 L(\gamma)| \xrightarrow{p} 0$ .

**Step 4:** By Taylor's expansion,  $\hat{L}_n(\hat{\gamma}) - \hat{L}_n(\gamma) = 0 - \hat{L}_n(\gamma) = \nabla_t \hat{L}_n(\xi)(\hat{\gamma} - \gamma)$  for some

$\xi \in (\gamma, \hat{\gamma})$ . From the results in Steps 1-3,  $\sqrt{N}(\hat{\gamma} - \gamma) \xrightarrow{d} N(0, [2E\{R \nabla_\gamma h(\mathbf{U}, \gamma)\}^2]^{-1} \sigma_L^2)$ . This completes the proof of (20).

**Proof of Corollary 1.** (i) From Theorem 3, (20) shows that  $\hat{\gamma} - \gamma = O_p(1/\sqrt{N})$ . Further-

more, Lemma 8.10 in Newey and McFadden (1994) shows that

$$\max_i \left| \frac{e^{-\gamma Y_i} \sum_{j=1}^n e^{\gamma Y_j} \check{\kappa}_\ell(\mathbf{U}_i - \mathbf{U}_j)}{\sum_{j=1}^n \check{\kappa}_\ell(\mathbf{U}_i - \mathbf{U}_j)} - \frac{f(Y_i|\mathbf{U} = \mathbf{U}_i, D = 1)}{f(Y_i|\mathbf{U} = \mathbf{U}_i, D = 0)} \right| = O_p \left( \sqrt{\frac{\log N}{N \ell^p}} + \ell^d \right), \quad (28)$$

which is  $o_p(N^{-2/(p+4)}) = o_p(1)/\sqrt{nb^p}$  under the assumed condition  $d > \max\{(p+4)/2, p\}$  and  $N \ell^{2d} \rightarrow 0$ . Since  $\hat{\gamma} - \gamma$  converges faster than (28), (18) holds. As a result, (17) holds with  $\mu_1^{E2}(\mathbf{u})$  replaced by  $\mu_1^{E3}(\mathbf{u})$  under (B1)-(B4) and (D1)-(D5).

(ii) Under (D1)-(D5) with  $\mathbf{U}$  replaced by  $\mathbf{X}$  and  $p$  replaced by  $q$ , Lemma 8.10 in Newey and McFadden (1994) implies that

$$\sup_{\mathbf{x} \in \mathbb{X}} |\hat{h}(\mathbf{x}, \gamma) - h_1(\mathbf{x})| = O_p((\log N)^{1/2} (N \ell^q)^{-1/2} + \ell^d) = o_p(n^{-2/(p+4)}).$$

From the asymptotic normality of  $\hat{\gamma}$ ,  $\hat{\gamma} - \gamma = O_p(1/\sqrt{N})$ , which converges to 0 faster than  $\sup_{\mathbf{x} \in \mathcal{X}} |\hat{h}(\mathbf{x}, \gamma) - h_1(\mathbf{x})| \rightarrow 0$ . Hence (27) holds while  $\hat{h}_1$  is estimated by  $\hat{h}(\mathbf{X}, \hat{\gamma})$ . Then, the rest of proof of the second claims follows the argument in the proof of Theorem 2.

## Acknowledgements

Jun Shao's research was partially supported by the National Natural Science Foundation of China (11831008) and the U.S. National Science Foundation (DMS-1914411).

## References

- Bierens, H. J. (1987). Kernel estimators of regression functions. In *Advances in Econometrics: Fifth World Congress*, Volume 1, pp. 99–144.
- Chatterjee, N., Y. H. Chen, P. Maas, and R. J. Carroll (2016). Constrained maximum

- likelihood estimation for model calibration using summary-level information from external big data sources. *Journal of the American Statistical Association* 111(513), 107–117.
- Cook, R. D. and S. Weisberg (1991). Sliced inverse regression for dimension reduction: Comment. *Journal of the American Statistical Association* 86(414), 328–332.
- Dai, C.-S. and J. Shao (2022). Kernel regression utilizing external information as constraints. *Statistica Sinica* 33, under revision.
- Fan, J., M. Farmen, and I. Gijbels (1998). Local maximum likelihood estimation and inference. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 60(3), 591–608.
- Fan, J., T. Gasser, I. Gijbels, M. Brockmann, and J. Engel (1997). Local polynomial regression: optimal kernels and asymptotic minimax efficiency. *Annals of the Institute of Statistical Mathematics* 49(1), 79–99.
- Györfi, L., M. Kohler, A. Krzyżak, and H. Walk (2002). *A Distribution-Free Theory of Nonparametric Regression*. Springer, New York.
- Kim, H. J., Z. Wang, and J. K. Kim (2021). Survey data integration for regression analysis using model calibration. *arXiv* 2107.06448.
- Li, B. and S. Wang (2007). On directional regression for dimension reduction. *Journal of the American Statistical Association* 102(479), 997–1008.
- Li, K.-C. (1991). Sliced inverse regression for dimension reduction. *Journal of the American Statistical Association* 86(414), 316–327.

- Lohr, S. L. and T. E. Raghunathan (2017). Combining survey data with other data sources. *Statistical Science* 32(2), 293–312.
- Ma, Y. and L. Zhu (2012). A semiparametric approach to dimension reduction. *Journal of the American Statistical Association* 107(497), 168–179.
- Merkouris, T. (2004). Combining independent regression estimators from multiple surveys. *Journal of the American Statistical Association* 99(468), 1131–1139.
- Nadaraya, E. A. (1964). On estimating regression. *Theory of Probability & Its Applications* 9(1), 141–142.
- Newey, W. K. (1994). Kernel estimation of partial means and a general variance estimator. *Econometric Theory* 10, 233–253.
- Newey, W. K. and D. McFadden (1994). Large sample estimation and hypothesis testing. *Handbook of Econometrics* 4, 2111–2245.
- Rao, J. (2021). On making valid inferences by integrating data from surveys and other sources. *Sankhya B* 83(1), 242–272.
- Shao, J. (2003). *Mathematical Statistics* (2nd ed.). Springer, New York.
- Shao, Y., R. D. Cook, and S. Weisberg (2007). Marginal tests with sliced average variance estimation. *Biometrika* 94(2), 285–296.
- Wand, M. P. and M. C. Jones (1994, December). *Kernel Smoothing*. Number 60 in Chapman & Hall/CRC Monographs on Statistics & Applied Probability. Boca Raton, FL, U.S.: Chapman & Hall.

- Wasserman, L. (2006). *All of Nonparametric Statistics*. Springer, New York.
- Xia, Y., H. Tong, W. K. Li, and L.-X. Zhu (2002). An adaptive estimation of dimension reduction space. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)* 64(3), 363–410.
- Yang, S. and J. K. Kim (2020). Statistical data integration in survey sampling: a review. *Japanese Journal of Statistics and Data Science* 3(2), 625–650.
- Zhang, H., L. Deng, M. Schiffman, J. Qin, and K. Yu (2020). Generalized integration model for improved statistical inference by leveraging external summary data. *Biometrika* 107(3), 689–703.
- Zhang, Y., Z. Ouyang, and H. Zhao (2017). A statistical framework for data integration through graphical models with application to cancer genomics. *The Annals of Applied Statistics* 11(1), 161–184.
- Zieschang, K. D. (1990). Sample weighting methods and estimation of totals in the consumer expenditure survey. *Journal of the American Statistical Association* 85(412), 986–1001.

Table 1: Simulated MISE (21) and IMP (22) with  $S = 200$  under  $\gamma = 0$

Covariate	Model	Test data	$b, l$	Estimator				IMP %	Mean of $\hat{\gamma}$
				$\hat{\mu}_1$	$\hat{\mu}_1^{C1}$	$\hat{\mu}_1^{C2}$	$\hat{\mu}_1^{C3}$		
Normal	M1	Sample	Best	0.057	0.038	0.051	0.044	33.26	-0.003
			CV	0.073	0.046	0.060	0.055	36.66	-0.006
		Grid	Best	0.054	0.019	0.038	0.030	64.86	-0.012
			CV	0.063	0.036	0.054	0.047	42.10	-0.013
	M2	Sample	Best	0.080	0.073	0.081	0.078	8.75	-0.033
			CV	0.091	0.083	0.093	0.089	8.79	-0.033
		Grid	Best	0.093	0.085	0.095	0.090	8.60	-0.037
			CV	0.110	0.101	0.115	0.108	8.18	-0.043
	M3	Sample	Best	0.077	0.072	0.078	0.076	6.49	-0.016
			CV	0.087	0.080	0.088	0.085	6.97	-0.018
		Grid	Best	0.067	0.059	0.067	0.062	11.94	0.001
			CV	0.087	0.081	0.091	0.087	6.89	-0.006
	M4	Sample	Best	0.061	0.040	0.054	0.047	33.83	-0.004
			CV	0.076	0.051	0.064	0.059	32.94	-0.010
		Grid	Best	0.059	0.022	0.045	0.033	62.17	-0.000
			CV	0.064	0.040	0.060	0.054	36.77	-0.014
Bounded	M1	Sample	Best	0.022	0.004	0.014	0.014	80.43	-0.003
			CV	0.029	0.005	0.014	0.014	80.96	-0.007
		Grid	Best	0.059	0.019	0.034	0.040	69.79	0.008
			CV	0.090	0.042	0.061	0.065	53.57	-0.014
	M2	Sample	Best	0.039	0.032	0.041	0.037	17.94	-0.014
			CV	0.044	0.039	0.047	0.044	11.36	-0.016
		Grid	Best	0.124	0.119	0.126	0.120	16.93	-0.005
			CV	0.158	0.148	0.170	0.157	20.25	-0.032
	M3	Sample	Best	0.033	0.027	0.035	0.034	19.54	-0.006
			CV	0.039	0.033	0.042	0.041	13.63	-0.009
		Grid	Best	0.103	0.097	0.103	0.101	9.34	0.009
			CV	0.130	0.125	0.144	0.139	9.81	-0.017
	M4	Sample	Best	0.023	0.005	0.015	0.014	77.03	-0.003
			CV	0.029	0.006	0.015	0.015	78.59	-0.007
		Grid	Best	0.063	0.025	0.040	0.044	62.57	0.007
			CV	0.096	0.049	0.066	0.072	52.23	-0.015

$\hat{\mu}_1$ : the standard kernel estimator given by (2).

$\hat{\mu}_1^{C1}$ ,  $\hat{\mu}_1^{C2}$ , and  $\hat{\mu}_1^{C3}$ : the estimators given by (16) as described in the end of Section 2.2.

Simulation standard deviations of  $\hat{\gamma}$  for all cases are between 0.005 and 0.006.

Table 2: Simulated MISE (21) and IMP (22) with  $S = 200$  under  $\gamma = 1/2$

Covariate	Model	Test data	$b, l$	Estimator				IMP %	Mean of $\hat{\gamma}$
				$\hat{\mu}_1$	$\hat{\mu}_1^{C1}$	$\hat{\mu}_1^{C2}$	$\hat{\mu}_1^{C3}$		
Normal	M1	Sample	Best	0.054	0.318	0.048	0.040	25.26	0.453
			CV	0.068	0.341	0.057	0.051	17.08	0.458
		Grid	Best	0.049	0.190	0.035	0.028	42.18	0.449
			CV	0.056	0.259	0.047	0.040	28.00	0.441
	M2	Sample	Best	0.082	0.509	0.083	0.081	1.07	0.426
			CV	0.093	0.528	0.095	0.091	1.12	0.429
		Grid	Best	0.089	0.588	0.089	0.085	4.21	0.419
			CV	0.099	0.608	0.103	0.098	1.31	0.426
	M3	Sample	Best	0.084	0.560	0.085	0.084	-0.50	0.449
			CV	0.097	0.565	0.101	0.095	1.99	0.442
		Grid	Best	0.070	0.513	0.069	0.067	4.42	0.452
			CV	0.082	0.551	0.087	0.081	0.38	0.456
	M4	Sample	Best	0.063	0.335	0.056	0.051	18.90	0.442
			CV	0.078	0.358	0.064	0.058	18.90	0.432
		Grid	Best	0.053	0.189	0.040	0.033	37.11	0.439
			CV	0.062	0.290	0.060	0.053	14.89	0.440
Bounded	M1	Sample	Best	0.021	0.242	0.015	0.013	36.29	0.480
			CV	0.028	0.241	0.015	0.014	42.94	0.475
		Grid	Best	0.057	0.326	0.035	0.039	31.95	0.483
			CV	0.083	0.346	0.061	0.064	22.69	0.487
	M2	Sample	Best	0.040	0.338	0.042	0.040	0.76	0.478
			CV	0.049	0.350	0.051	0.050	-2.07	0.469
		Grid	Best	0.127	0.555	0.125	0.122	4.15	0.463
			CV	0.155	0.591	0.163	0.154	0.59	0.463
	M3	Sample	Best	0.033	0.328	0.035	0.034	-4.03	0.475
			CV	0.039	0.343	0.041	0.041	-6.37	0.486
		Grid	Best	0.105	0.484	0.102	0.101	4.04	0.483
			CV	0.126	0.519	0.130	0.126	-0.30	0.482
	M4	Sample	Best	0.021	0.243	0.016	0.014	31.91	0.476
			CV	0.030	0.238	0.017	0.015	42.21	0.477
		Grid	Best	0.064	0.344	0.043	0.044	30.67	0.486
			CV	0.086	0.370	0.068	0.068	20.44	0.489

$\hat{\mu}_1$ : the standard kernel estimator given by (2).

$\hat{\mu}_1^{C1}$ ,  $\hat{\mu}_1^{C2}$ , and  $\hat{\mu}_1^{C3}$ : the estimators given by (16) as described in the end of Section 2.2.

Simulation standard deviations of  $\hat{\gamma}$  for all cases are between 0.005 and 0.006.