

# Minimum Hellinger Distance Estimators for Complex Survey Designs

David Kepplinger and Anand N. Vidyashankar

Department of Statistics, George Mason University, Fairfax, VA, USA

## Abstract

Reliable inference from complex survey samples can be derailed by outliers and high-leverage observations induced by unequal inclusion probabilities and calibration. We develop a minimum Hellinger distance estimator (MHDE) for parametric superpopulation models under complex designs, including Poisson PPS and fixed-size SRS/PPS without replacement, with possibly stochastic post-stratified or calibrated weights. Using a Horvitz–Thompson–adjusted kernel density plug-in, we show: (i)  $L^1$ –consistency of the KDE with explicit large-deviation tail bounds driven by a variance-adaptive effective sample size; (ii) uniform exponential bounds for the Hellinger affinity that yield MHDE consistency under mild identifiability; (iii) an asymptotic Normal distribution for the MHDE with covariance  $\mathbf{A}^{-1}\boldsymbol{\Sigma}\mathbf{A}^\top$  (and a finite-population correction under without-replacement designs); and (iv) robustness via the influence function and  $\alpha$ –influence curves in the Hellinger topology. Simulations under Gamma and lognormal superpopulation models quantify efficiency–robustness trade-offs relative to weighted MLE under independent and high-leverage contamination. An application to NHANES 2021–2023 total water consumption shows that the MHDE remains stable despite extreme responses that markedly bias the MLE. The estimator is simple to implement via quadrature over a fixed grid and is extensible to other divergence families.

*Keywords:* Hellinger distance, complex survey design, Horvitz–Thompson, probability-proportional-to-size (PPS), kernel density estimation, large deviations, asymptotic normality, robust estimation, influence function.

# 1 Introduction

Reliable estimation and inference in complex survey samples is a challenging problem, particularly when outliers can be present. In this work, we develop a robust estimator and asymptotic inference for survey samples from a finite population with possibly unequal inclusion probabilities which can be based on auxiliary information. Outliers, i.e., unusually small or large values, in the observed sample must be handled carefully to avoid biased and invalid inference from the sample survey. These outliers may be legitimate values, but can also be caused by data entry errors and other problems. Regardless of the legitimacy of these unusual values, inclusion probabilities from auxiliary information can drastically amplify the outliers' effects on the estimator. Borrowing from the terminology of high breakdown estimators for linear regression, we call outliers in units with low inclusion probability *high-leverage observations*.

In large-scale surveys, it is common to adjust the survey weights derived from the inclusion probabilities, to match certain characteristics to known totals for the entire population or within strata [1]. Post-stratification or calibration leads to stochastic survey weights even if the initial inclusion probabilities are deterministic. When such adjustments are applied, outliers and high-leverage observations can be further amplified and have an even more detrimental effect on an estimator.

We propose a reliable minimum Hellinger distance estimator (MHDE) for model parameters under complex survey designs, with potentially random survey weights. Minimum divergence estimators are known for their robustness toward outliers without sacrificing efficiency in clean samples in a wide range of models and settings [e.g., 2–6]. Recently, minimum phi-divergence estimators have been shown to achieve robustness and high efficiency for multinomial and polytomous logistic regression in complex survey designs [7, 8]. Following that line of research, we develop an MHDE for the parameters of a superpopulation model from a survey sample. We allow inclusion probabilities derived from auxiliary information, e.g., probability proportional to size (PPS) sampling [9], cluster sampling or stratified sampling, and stochastic survey weights adjusted by post-stratification or calibration.

In Section 2 we define our MHDE for complex survey designs and show in Section 3 that it is consistent under mild assumptions, is asymptotically Normal and robust under arbitrary contamination. The empirical studies in Section 4 demonstrate that the estimator is highly efficient and yields valid inference, even in the presence of outliers and high-leverage observations. We further apply our estimator to the National Health and Nutrition Examination Survey (NHANES) [10], where we show that our MHDE is much less affected by unusual values than the maximum likelihood estimator.

## 1.1 Background

For each  $\gamma \in \mathbb{N}$  we consider a finite population of  $N_\gamma$  units  $\mathcal{U}_\gamma = \{1, \dots, N_\gamma\}$ . We observe i.i.d. draws  $\{(Y_{\gamma i}, Z_{\gamma i}) : i \in \mathcal{U}_\gamma\}$  from a superpopulation law on  $\mathbb{R}^{d+1} \times (0, \infty)$ . The  $Y_{\gamma i}$  are the characteristics of interest with unknown but measurable and integrable density  $g \in L^1(\mathbb{R}^d)$ . The auxiliary variable  $Z_{\gamma i}$  can be used to derive inclusion probabilities and, if used, is assumed to be known and greater than 0 with probability 1. From  $\mathcal{U}_\gamma$ , a sample of size  $n_\gamma$  is drawn according to pre-specified, potentially unequal, inclusion probabilities  $\pi_{\gamma i} > 0$ ,  $i \in \mathcal{U}_\gamma$ . The units included in the random sample are denoted by  $\mathcal{S}_\gamma \subset \mathcal{U}_\gamma$ .

For simple designs, such as fixed-size simple random sampling (SRS) with or without replacement, the inclusion probabilities are equal for all units. However, in many survey samples, the design is more complicated. In this work, we focus on the probability proportional-to-size (PPS) design with random size (Poisson-PPS), where the inclusion probabilities depend on an auxiliary variable,  $Z_{\gamma i}$ ,  $i \in \mathcal{U}_\gamma$ , with  $\pi_{\gamma i} = \frac{n_\gamma Z_{\gamma i}}{\sum_{k \in \mathcal{U}_\gamma} Z_{\gamma k}}$  and hence  $\sum_i \pi_{\gamma i} = n_\gamma$ . In the PPS design,  $Z_{\gamma i}$  is known prior to sampling for all units in  $\mathcal{U}_\gamma$ , e.g., the earnings of business entities in the previous year(s) or the taxable income of households. If the auxiliary variable,  $Z_{\gamma i}$ , is correlated with  $Y_{\gamma i}$ , PPS sampling can reduce the sampling variance of an estimator. Other sampling designs also use auxiliary information to derive inclusion probabilities, such as cluster sampling or stratified sampling based on geographic location or school districts, to name just a few. Unequal inclusion probabilities yield non-identically distributed observations, as the unit-level density becomes  $\tilde{g}_{\gamma, i}(y) = \pi_{\gamma i} g(y)$ .

Based on the inclusion probabilities, we define the sample weight for each unit in the sample as  $w_{\gamma i} = 1/\pi_{\gamma i}$ ,  $i \in \mathcal{S}_\gamma$ . These sample weights need to be considered in the estimation to achieve consistency and reduce bias, e.g., using the Horvitz-Thompson (HT) adjustment [11]. However, with post-stratification or calibration, the sample weights may be further adjusted to  $\omega_{\gamma i} = \zeta_{\gamma i} w_{\gamma i}$ , where  $\zeta_{\gamma i}$  is a positive random factor with  $\mathbb{P}(\zeta_{\gamma i} > 0) = 1$ .

In this paper our goal is to find a parametric density from a family  $\mathcal{F} := \{f_\theta : \theta \in \Theta \subseteq \mathbb{R}^p\}$  which is “closest” to the true superpopulation distribution in the topology defined by divergence  $D$ . Hence, we seek  $f_{\hat{\theta}}$  where  $\hat{\theta}_\gamma = \arg \min_{\theta \in \Theta} D(f_\theta \| g)$ .

The theoretical and empirical properties of  $\hat{\theta}$  are intricately linked to the divergence,  $D$ . Information divergences between probability density functions are a rich family of measures, but not all are suitable under model misspecification, i.e.,  $g \notin \mathcal{F}$ , for example the Tukey-Huber  $\varepsilon$ -contamination model Tukey [12]. The Kullback-Leibler divergence, for example, yields the maximum likelihood estimate [13] but can lead to arbitrarily biased estimates under contamination. In this paper, we therefore focus on the more robust (squared) Hellinger distance:

$$H^2(f, g) = \frac{1}{2} \int_{\Omega} \left( \sqrt{f(y)} - \sqrt{g(y)} \right)^2 dy = 1 - \int_{\Omega} \sqrt{f(y)g(y)} dy. \quad (1)$$

The Hellinger divergence is known to yield estimates that are robust towards model misspecification [2], yet achieve high efficiency if  $g \in \mathcal{F}$  [5]. Important for large-scale surveys, the MHDE can be quickly computed using numerical integration if the dimension  $Y$  is reasonably low. In the following Section 2 we describe the MHDE for complex survey designs based on a Horvitz-Thompson adjusted Kernel Density Estimator.

## 1.2 Notation

Throughout we denote by  $\delta_{\gamma i} \in \{0, 1\}$  whether a unit in the finite population  $\mathcal{U}_\gamma$  is included in the sample  $\mathcal{S}_\gamma$  or not, i.e.,  $\delta_{\gamma i} = 1$  if  $i \in \mathcal{S}_\gamma$  and  $\delta_{\gamma i} = 0$  otherwise. The first-order inclusion probabilities are thus  $\pi_{\gamma i} = \mathbb{P}(\delta_{\gamma i} = 1)$ . We let the “effective” sample size be  $n_{\text{eff}, \gamma} := N_\gamma^2 / \sum_{i \in \mathcal{U}_\gamma} \pi_{\gamma i}^{-1}$ , and the variance-adaptive effective sample size  $n_{\text{V-eff}, \gamma} := N_\gamma^2 / \sum_{i \in \mathcal{U}_\gamma} (1 - \pi_{\gamma i}) / \pi_{\gamma i}$ . For fixed-size designs, such as SRS-WOR or fixed-size PPS-WOR, we write  $\alpha_\gamma := n_\gamma / N_\gamma$ . The bandwidth of the kernel density estimator depends on the sample size  $n_\gamma$  and we denote it by  $h_\gamma > 0$ . We then write the normalized kernel as  $K_{h_\gamma}(x) = h_\gamma^{-d} K(x/h_\gamma)$ .

We denote the Hellinger affinity (Bhattacharyya coefficient) between the parametric density  $f_\theta$  and the KDE  $\hat{f}_\gamma$  or the (arbitrary) distribution  $F$  with density  $f$ , respectively, by

$$\Gamma_\gamma(\theta) := \int \sqrt{\hat{f}_\gamma(y) f_\theta(y)} dy, \quad \Gamma_F(\theta) := \int \sqrt{f(y) f_\theta(y)} dy.$$

We simply write  $\Gamma(\theta) := \Gamma_G(\theta)$  when referring to the true superpopulation distribution.

Finally, we define the score function as  $u_\theta(y) := \nabla_\theta \log f_\theta(y)$  and use

$$\phi_g(y) := \frac{1}{4} u_{\theta_0}(y) \sqrt{\frac{f_{\theta_0}(y)}{g(y)}}, \quad \Sigma := \mathbb{E}_G [\phi_g(Y) \phi_g(Y)^\top], \quad \mathbf{A} := -\nabla_\theta^2 \Gamma(\theta) \Big|_{\theta=\theta_0},$$

to denote the scaled score function, the expected information and the Hessian of the Hellinger affinity, respectively. Where obvious, we omit the subscript from the scaled score function and write  $\phi(y) := \phi_g(y)$ .

## 2 Methodology

Let the true superpopulation distribution of  $Y$  be  $G$  with density  $g$ . Introducing the parametric family  $\mathcal{F} = \{f_\theta : \theta \in \Theta\}$ , the population minimizer  $\theta_0 = \arg \min_{\theta \in \Theta} H^2(f_\theta, g)$  represents the “closest” parametric density in  $\mathcal{F}$  to  $G$  in the Hellinger topology. To estimate  $\theta_0$ , we minimize the Hellinger distance between the estimated density and the densities in the parametric family:

$$\hat{\theta}_\gamma := \arg \min_{\theta \in \Theta} H^2(f_\theta, \hat{f}_\gamma) = \arg \max_{\theta \in \Theta} \Gamma_\gamma(\theta). \quad (2)$$

To obtain a consistent estimate of  $\theta_0$ , we use the Horvitz-Thompson (HT) adjusted kernel

density estimator:

$$\hat{f}_\gamma(y) := \frac{1}{\sum_{i \in \mathcal{U}_\gamma} \delta_{\gamma i} / \pi_{\gamma i}} \sum_{i \in \mathcal{U}_\gamma} \frac{\delta_{\gamma i}}{\pi_{\gamma i}} K_{h_\gamma}(Y_{\gamma i} - y). \quad (3)$$

Underpinning the robustness properties of the MHDE defined in (2) is its continuity in the Hellinger topology, as shown in Proposition C.2 in the Appendix under mild conditions. As the proportion of contaminated observations decreases, the estimator converges to the maximizer of (2) without contamination.

## 2.1 A Note on Computation

Our software implementation maximizes  $\Gamma_\gamma(\theta)$  by Nelder-Mead [14]. The integral in  $\Gamma_\gamma$  is computed over a fixed grid  $\mathcal{G}_Y$  using the Gauss-Kronrod quadrature and a given number of subdivisions. Therefore,  $\hat{f}_\gamma$  must be evaluated only once for each  $y \in \mathcal{G}_Y$ . Particularly for large  $n_\gamma$ , this substantially eases the computational burden compared to adaptive quadrature. The grid is chosen to cover only the regions where  $\hat{f}_\gamma > 0$ , which can be quickly evaluated knowing the kernel and bandwidth.

## 3 Theory

We present three main results for the MHDE (2) in the finite population setting under the superpopulation model framework. We first show that the HT-adjusted KDE under PPS sampling converges in  $L_1$  to the superpopulation density  $f_\gamma$ , while the naïve KDE converges to a size-biased density. We then prove that the MHDE based on the HT-adjusted KDE is consistent for  $\theta_0$  and derive its limiting normal distribution under several sample designs. Finally, we obtain the influence function and demonstrate the robustness of the estimator.

In the following, we write the HT-adjusted KDE as  $\hat{f}_\gamma(y) = \frac{1}{S_\gamma} T_\gamma(y)$ , with

$$S_\gamma := \frac{1}{N_\gamma} \sum_{i \in \mathcal{U}_\gamma} \frac{\delta_{\gamma i}}{\pi_{\gamma i}}, \quad T_\gamma(y) := \frac{1}{N_\gamma} \sum_{i \in \mathcal{U}_\gamma} \frac{\delta_{\gamma i}}{\pi_{\gamma i}} K_{h_\gamma}(y - Y_{\gamma i}).$$

### 3.1 Consistency of the HT-adjusted KDE

For consistency to hold, we assume that the kernel function  $K$  is smooth and that the bandwidth decreases at a prescribed rate. We also make concrete our assumptions about the superpopulation model and the regularity of the design.

**A1** (Smoothness of the kernel). The kernel  $K \in L^1 \cap L^2$  is bounded, non-negative, Lipschitz ( $\nabla K \in L^1$ ) and integrates to one,  $\int_{\mathbb{R}^d} K(x) dx = 1$ .

**A2** (Bandwidth and growth). The bandwidth  $h_\gamma \rightarrow 0$  such that  $n_{\text{eff},\gamma} h_\gamma^d \rightarrow \infty$  and  $N_\gamma h_\gamma^d \rightarrow \infty$ . Moreover,  $\alpha_\gamma \rightarrow \alpha \in (0, 1]$ .

**A3** (Superpopulation model).  $(Y_{\gamma i}, Z_{\gamma i})$  are i.i.d. across  $i \in \mathcal{U}_\gamma$  with  $Y_{\gamma i} \sim g \in L^1(\mathbb{R}^d)$ . The design may depend on  $Z$  but not directly on  $Y$  given  $Z$  (PPS).

**A4** (Design regularity). There exists  $0 < c_0 < \infty$  such that

$$\lim_{\gamma \rightarrow \infty} \mathbb{P} \left( \max_{i \in \mathcal{U}_\gamma} \pi_{\gamma i}^{-1} \leq c_0 / \alpha_\gamma \right) = 1.$$

Equivalently, we write  $\max_i \pi_{\gamma i}^{-1} = O_p(1/\alpha_\gamma)$ . To satisfy this assumption in applications, extremely large inverse inclusion weights can be truncated.

Lemma A.1 in the Appendix shows that under these assumptions  $\hat{f}_\gamma(y)$  is self-normalizing, i.e., integrates to 1 for every sample. The following theorem shows that the HT-adjusted KDE converges to the true density  $g$  in  $L^1$ .

**Theorem 3.1** (Large-deviation-based  $L_1$ -consistency of HT-adjusted KDE). *Under Assumptions A1–A4,*

$$\|\hat{f}_\gamma - g\|_{L^1} \xrightarrow{\mathbb{P}} 0.$$

Moreover, there exist constants  $C_1, C_2, C_3 > 0$ , depending only on  $K$  and  $c_0$ , such that for all  $\tau \in (0, 1]$ ,

$$\mathbb{P} \left( \|\hat{f}_\gamma - g\|_1 > \tau \right) \leq C_1 \exp \left\{ -C_2 n_{\text{eff}, \gamma} h_\gamma^d \tau^2 \right\} + C_1 \exp \left\{ -C_3 N_\gamma h_\gamma^d \tau^2 \right\} + o(1).$$

If in addition  $n_{\text{eff}, \gamma} h_\gamma^d / \log(1/h_\gamma) \rightarrow \infty$ , then  $\|\hat{f}_\gamma - g\|_1 \rightarrow 0$  almost surely.

A key ingredient in the proof of the  $L^1$  consistency is the following proposition about the large-deviation bounds for the design term.

**Proposition 3.2** (Direct large-deviation bounds for the design term). *Under assumptions A1 and A4, and letting*

$$\bar{f}_{\gamma, h}(y) := \frac{1}{N_\gamma} \sum_{i \in \mathcal{U}_\gamma} K_{h_\gamma}(y - Y_{\gamma i}),$$

there exist constants  $c, C > 0$  (depending only on  $c_0, K, d$ ) such that for all  $\tau \in (0, 1]$ ,

$$\mathbb{P} \left( \|T_\gamma - \bar{f}_{\gamma, h}\|_{L^1} > \tau \mid \{Y, Z\} \right) \leq C \exp \left\{ -c n_{\text{eff}, \gamma} h_\gamma^d \min(\tau^2, \tau) \right\} + C \exp \left\{ -c N_\gamma h_\gamma^d \right\}.$$

Under sampling without replacement (rejective), the first exponent is multiplied by  $(1 - \alpha_\gamma)$ .

The proof of the proposition as well as the consistency of the HT-adjusted KDE are given in Appendix A.

*Remark 3.3* (Rates under smoothness). If  $g$  is  $\beta$ -Hölder and  $K$  has order  $\beta$ , the three-way decomposition in the large-deviation bound yields

$$\|\hat{f}_\gamma - g\|_1 = O_{\mathbb{P}}(h_\gamma^\beta) + O_{\mathbb{P}}\left((n_{\text{eff},\gamma} h_\gamma^d)^{-1/2}\right) + O_{\mathbb{P}}\left((N_\gamma h_\gamma^d)^{-1/2}\right).$$

Since  $N_\gamma \geq n_{\text{eff},\gamma}$ , the last term is dominated by the middle one. Choosing  $h_\gamma \asymp n_{\text{eff},\gamma}^{-1/(2\beta+d)}$  balances the (design) variance and bias, giving

$$\|\hat{f}_\gamma - g\|_1 = O_{\mathbb{P}}\left(n_{\text{eff},\gamma}^{-\beta/(2\beta+d)}\right).$$

**Corollary 3.4** (Simple random sampling). *If  $\pi_{\gamma i} \equiv n_\gamma/N_\gamma$  (i.e., simple random sampling with replacement), then  $n_{\text{eff},\gamma} \asymp n_\gamma$  and Theorem 3.1 recovers the well-known triangular-array SRS result: if  $h_\gamma \downarrow 0$  and  $n_\gamma h_\gamma^d \rightarrow \infty$ , then  $\|\hat{f}_\gamma - g\|_1 \rightarrow 0$  in probability.*

### 3.2 Consistency of the MHDE

The MHDE with HT-adjusted KDE plug-in (2) is equivalent to any measurable maximizer  $\hat{\theta}_\gamma \in \arg \max_{\theta \in \Theta} \Gamma_\gamma(\theta)$ , and we make the following identifiability assumption:

**A5** (Identifiability).  $\theta_0$  uniquely maximizes  $\Gamma(\theta)$  and, for each  $\varepsilon > 0$ ,  $\sup_{\|\theta - \theta_0\| \geq \varepsilon} \Gamma(\theta) \leq \Gamma(\theta_0) - \Delta(\varepsilon)$  for some  $\Delta(\varepsilon) > 0$ .

The following proposition establishes the tail bounds for the MHDE deviations and is proven in Appendix A.5.

**Proposition 3.5** (Exponential tail bounds for uniform MHDE deviation). *Under Assumptions A1–A4, for all  $t \in (0, 1]$ ,*

$$\mathbb{P}\left(\sup_{\theta} |\Gamma_\gamma(\theta) - \Gamma(\theta)| > t\right) \leq C \exp\{c n_{\text{eff},\gamma} h_\gamma^d \min\{t^4, t^2\}\} + C \exp\{-c N_\gamma h_\gamma^d\}. \quad (4)$$

*For Poisson-PPS without replacement, multiply the first exponent by  $(1 - \alpha_\gamma)$ .*

**Theorem 3.6** (Consistency of MHDE with HT plug-in). *Under Assumptions A1–A5, let  $\hat{\theta}_\gamma$  be any sequence of  $\varepsilon_\gamma$ -maximizers, i.e.,*

$$\Gamma_\gamma(\hat{\theta}_\gamma) \geq \sup_{\theta} \Gamma_\gamma(\theta) - \varepsilon_\gamma$$

*with  $\varepsilon_\gamma \rightarrow 0$ . Then  $\hat{\theta}_\gamma \rightarrow \theta_0$  in probability. Moreover, for every  $\varepsilon > 0$ ,*

$$\mathbb{P}\left(\|\hat{\theta}_\gamma - \theta_0\| \geq \varepsilon\right) \leq \mathbb{P}\left(\sup_{\theta} |\Gamma_\gamma(\theta) - \Gamma(\theta)| > \frac{1}{3}\Delta(\varepsilon)\right) + \mathbf{1}\{\varepsilon_\gamma > \frac{1}{3}\Delta(\varepsilon)\}.$$

*In particular, using (4), the RHS decays at the exponential rate in  $n_{\text{eff},\gamma} h_\gamma^d$ .*

*Proof.* Fix  $\varepsilon > 0$  and write  $\Delta = \Delta(\varepsilon)$  from Assumption **A5**. For the event

$$\mathcal{E}_\gamma(\varepsilon) := \left\{ \sup_{\theta} |\Gamma_\gamma(\theta) - \Gamma(\theta)| \leq \frac{1}{3}\Delta \right\} \cap \{\varepsilon_\gamma \leq \frac{1}{3}\Delta\},$$

we claim that every  $\varepsilon_\gamma$ -maximizer  $\hat{\theta}_\gamma$  lies in  $B(\theta_0, \varepsilon)$ . In fact, for any  $\theta$  with  $\|\theta - \theta_0\| \geq \varepsilon$ ,

$$\Gamma_\gamma(\theta) \leq \Gamma(\theta) + \frac{1}{3}\Delta \leq \Gamma(\theta_0) - \frac{2}{3}\Delta,$$

while at  $\theta_0$ ,  $\Gamma_\gamma(\theta_0) \geq \Gamma(\theta_0) - \frac{1}{3}\Delta$ . Thus  $\sup_{\|\theta - \theta_0\| \geq \varepsilon} \Gamma_\gamma(\theta) \leq \Gamma_\gamma(\theta_0) - \frac{1}{3}\Delta$ . If  $\hat{\theta}_\gamma$  is an  $\varepsilon_\gamma$ -maximizer with  $\varepsilon_\gamma \leq \Delta/3$ , then

$$\Gamma_\gamma(\hat{\theta}_\gamma) \geq \sup_{\theta} \Gamma_\gamma(\theta) - \varepsilon_\gamma > \Gamma_\gamma(\theta_0) - \frac{1}{3}\Delta,$$

so  $\hat{\theta}_\gamma$  cannot lie outside  $B(\theta_0, \varepsilon)$ . Therefore

$$\mathbb{P}\left(\|\hat{\theta}_\gamma - \theta_0\| \geq \varepsilon\right) \leq \mathbb{P}(\mathcal{E}_\gamma(\varepsilon)^c) \leq \mathbb{P}\left(\sup_{\theta} |\Gamma_\gamma(\theta) - \Gamma(\theta)| > \frac{1}{3}\Delta\right) + \mathbf{1}\{\varepsilon_\gamma > \Delta/3\}.$$

Applying Proposition 3.5 yields the RHS  $\rightarrow 0$  as  $\gamma \rightarrow \infty$ .  $\square$

*Remark 3.7.* The argument uses only uniform (in  $\theta$ ) control via Lemma A.14 and separation in Assumption **A5**. Compactness of  $\Theta$  or upper semicontinuity of  $\Gamma$  are not required.

*Remark 3.8.* All statements hold verbatim under Poisson–PPS without replacement. In the tail bound (4), multiply the first exponent by  $(1 - \alpha_\gamma)$  (finite-population correction).

*Remark 3.9.* In case  $g \in \mathcal{F}$  such that  $g \equiv f_0$ ,  $\hat{\theta}_\gamma$  converges to the true parameter  $\theta_0$ .

### 3.3 Asymptotic normality

To derive the central limit theorem for  $\hat{\theta}_\gamma$  we need the following additional assumptions.

**A6** (Design) For Poisson–PPS,  $\sqrt{n_{V\text{-eff},\gamma}(S_\gamma - 1)} = O_{\mathbb{P}}(1)$  and

$$\max_{i \in \mathcal{U}_\gamma} \frac{(1 - \pi_{\gamma i})/\pi_{\gamma i}}{\sum_{j \in \mathcal{U}_\gamma} (1 - \pi_{\gamma j})/\pi_{\gamma j}} \xrightarrow{\mathbb{P}} 0.$$

For fixed-size SRS–WOR, with  $n_{V\text{-eff},\gamma} \equiv n_\gamma/(1 - \alpha_\gamma)$  and  $S_\gamma \equiv 1$ , these are satisfied if the finite-population correction (FPC) factor  $(1 - \alpha_\gamma) \rightarrow (1 - \alpha)$  with  $\alpha \in [0, 1)$ .

**A7** (Kernel approximation) Either  $|1 - K(\xi)| \lesssim |\xi|^\beta$  as  $\xi \rightarrow 0$  for some  $\beta > 0$ ; or  $K$  has order  $\beta > 0$  (i.e., vanishing moments up to  $\lfloor \beta \rfloor$ ) and  $\sqrt{g} \in \mathcal{C}^\beta$  on compacta. In both cases,  $K_h * \phi_g \rightarrow \phi_g$  in  $L^2(g)$  and  $g * K_h - g = O(h_\gamma^\beta)$  in the weighted sense used in the theorem below.



**A8** (Model smoothness and identifiability)  $\theta_0$  is the unique maximizer of  $\Gamma(\theta)$  with positive definite Hessian **A**. Moreover,  $\phi \in L^2(g; \mathbb{R}^p)$  and  $\theta \mapsto \sqrt{f_\theta(y)}$  is twice continuously differentiable in a neighborhood of  $\theta_0$ , with dominated derivatives allowing differentiation under the integral for  $\Gamma$  and  $\Gamma_\gamma$ .

**A9** (Bandwidth regime) The bandwidth  $h_\gamma \downarrow 0$  and

$$\sqrt{n_{V\text{-eff},\gamma}} h_\gamma^{2\beta} \rightarrow 0, \quad \frac{1}{\sqrt{n_{V\text{-eff},\gamma}} h_\gamma^d} \rightarrow 0, \quad \frac{\sqrt{n_{V\text{-eff},\gamma}}}{N_\gamma h_\gamma^d} \rightarrow 0.$$

**A10** (Localized risk control) Fix an exhaustion by compacta  $A_R \uparrow \mathbb{R}^d$  with  $c_R := \inf_{A_R} g > 0$  (automatic if  $g$  is continuous and strictly positive). For each  $R$ , there exist  $h_0(R) > 0$ ,  $C_R < \infty$ , and a tail remainder  $\tau_R(h) \downarrow 0$  (as  $R \uparrow \infty$  uniformly on  $h \leq h_0(R)$ ) such that for all  $0 < h \leq h_0(R)$ ,

$$\sup_{t \in \mathbb{R}^d} \int_{A_R} \frac{K_h^2(y-t)}{g(y)} dy \leq C_R h^{-d}, \quad \sup_{0 < h \leq h_0(R)} \int_{A_R^c} \frac{(K_h^2 * g)(y)}{g(y)} dy \leq \tau_R(h).$$

In addition, Assumption **A7** implies the bias bound on compacta

$$\int_{A_R} \frac{(g * K_h - g)^2}{g} dy \lesssim_R h^{2\beta}, \text{ and } \sup_{0 < h \leq h_0(R)} \int_{A_R^c} \frac{(g * K_h - g)^2}{g} dy \leq \tau_R(h).$$

**A11** (Lindeberg/no dominant unit) Let  $\psi_{h_\gamma} = K_{h_\gamma} * \phi_g$  and define

$$X_{\gamma i} := \frac{1}{N_\gamma} \left( \frac{\delta_{\gamma i}}{\pi_{\gamma i}} - 1 \right) \psi_{h_\gamma}(Y_{\gamma i}) \in \mathbb{R}^p.$$

Assume the Lindeberg condition for triangular-arrays holds conditional on  $\{Y_{\gamma i}\}$ , i.e., for every  $\varepsilon > 0$ ,

$$\frac{\sum_{i \in \mathcal{U}_\gamma} \mathbb{E} \left[ \|X_{\gamma i}\|^2 \mathbf{1} \{ \|X_{\gamma i}\| > \varepsilon / \sqrt{n_{V\text{-eff},\gamma}} \} \mid \{Y_{\gamma i}\} \right]}{\sum_{i \in \mathcal{U}_\gamma} \text{Var}(X_{\gamma i} \mid \{Y_{\gamma i}\})} \xrightarrow{\mathbb{P}} 0,$$

and  $\sum_i \text{Var}(X_{\gamma i} \mid \{Y\})$  converges in probability to  $\Sigma / n_{V\text{-eff},\gamma}$  (see the variance limit below).

Assumption **A6** is standard and mild for Poisson-PPS and automatically satisfied for SRS-WOR. The assumptions **A7** and **A9** are standard for KDE in finite populations, while **A10** is substantially weaker than the usual assumption of  $\inf g > 0$  globally. A sufficient condition for **A11** to hold is  $\mathbb{E}_G \|\phi(Y)\|^{2+\eta} < \infty$  for some  $\eta > 0$  together with the no-dominant-unit condition in Assumption **A6**.

**Corollary 3.10.** *Under assumptions **A6**, **A7** and **A10** with  $\psi_{h_\gamma} = K_{h_\gamma} * \phi_g$  we have*

$$n_{V\text{-eff},\gamma} \text{Var} \left( \frac{1}{N_\gamma} \sum_{i \in \mathcal{U}_\gamma} \left( \frac{\delta_{\gamma i}}{\pi_{\gamma i}} - 1 \right) \psi_{h_\gamma}(Y_{\gamma i}) \middle| \{Y_{\gamma i}\} \right) \xrightarrow{\mathbb{P}} \Sigma.$$

For SRS-WOR, the variance is multiplied by the FPC  $(1 - \alpha_\gamma)$ .

**Theorem 3.11.** *Under assumptions **A6**–**A11** the asymptotic distribution of the MHDE  $\hat{\theta}_\gamma$  under Poisson-PPS is Gaussian:*

$$\sqrt{n_{V\text{-eff},\gamma}}(\hat{\theta}_\gamma - \theta_0) \Rightarrow N_p(\mathbf{0}, \mathbf{A}^{-1} \Sigma \mathbf{A}^{-\top}).$$

For fixed-size SRS-WOR, the covariance matrix must be multiplied by the FPC  $(1 - \alpha_\gamma)$ .

The proof borrows ideas from Cheng and Vidyashankar [15] and is given in Appendix B. Here we want to discuss a few important insights from the proof technique.

*Remark 3.12.* Our proof does not require a global lower bound on  $f_0$ . All variance and bias controls in assumption **A10** localized on  $A_R$  with a tail remainder  $\tau_R(h)$  that can be driven to 0 by taking  $R = R_\gamma \uparrow \infty$  slowly, uniformly over  $h \leq h_0(R)$ .

*Remark 3.13.* The bandwidth assumption **A9** is necessary to remove the kernel bias, the i.i.d. KDE smoothing noise of order  $1/\sqrt{N_\gamma h_\gamma^d}$  and also to leave the HT design fluctuation at the scale of  $\sqrt{n_{V\text{-eff},\gamma}}$ .

*Remark 3.14.* Under SRS the assumptions reduce to the classical conditions  $\sqrt{n_\gamma} h_\gamma^d \rightarrow \infty$  and  $\sqrt{n_\gamma} h_\gamma^{2\beta}$ , up to the FPC, as in i.i.d. MHDE analyses without global lower bound on  $f$ .

*Remark 3.15.* For fixed-size PPS-WOR, assumption **A6** would need to be replaced with the usual rejective design with  $\sum_{i \in \mathcal{U}_\gamma} \delta_{\gamma i} = n_\gamma$  and first-order  $\{\pi_{\gamma i}\}$ . The same FPC factor as with SRS-WOR appears asymptotically, and the rest of the statement is unchanged.

### 3.4 Robustness

Finally, we turn to the robustness of the MHDE against contaminated superpopulations. We define the estimator functional

$$T(G) \in \arg \max_{\theta \in \Theta} \Gamma_G(\theta),$$

and the gradient of the population-level  $\Gamma_G$  as  $S_G(\theta) := \nabla_\theta \Gamma_G(\theta) = \frac{1}{2} \int u_\theta(y) \sqrt{f_\theta(y)g(y)}$ . The MHDE (2) targets  $\theta_0 = T(G)$ . Note that the sampling design does not affect the functional, only the estimator. Hence, the design does not affect the robustness properties.

We denote the contaminated superpopulation distribution by  $G_\epsilon := (1 - \epsilon)G + \epsilon H$  with arbitrary contamination distribution  $H$ . We work in the Hellinger topology,  $H(f_1, f_2) := \|\sqrt{f_1} - \sqrt{f_2}\|_{L^2}$  and make the following assumptions about the model smoothness.

**A12** (Model smoothness and identifiability).

(i) For each  $y$ ,  $\theta \mapsto \sqrt{f_\theta(y)}$  is twice continuously differentiable in a neighborhood  $\mathcal{N}(\theta_0)$  of  $\theta_0 = T(G)$ .

(ii) There exists an envelope  $e \in L^2(g)$  such that for all  $\theta \in \mathcal{N}(\theta_0)$ ,

$$\begin{aligned} \|u_\theta \sqrt{f_\theta}\|_{L^2(g)} &\leq \|e\|_{L^2(g)}, \\ \|\partial_\theta u_\theta \sqrt{f_\theta}\|_{L^2(g)} &\leq \|e\|_{L^2(g)}. \end{aligned}$$

(iii)  $\theta_0$  is the unique maximizer of  $\Gamma$  and there is a separation margin:  $\forall \varepsilon > 0, \exists \Delta(\varepsilon) > 0$  s.t.,

$$\sup_{\|\theta - \theta_0\| \geq \varepsilon} \Gamma(\theta) \leq \Gamma(\theta_0) - \Delta(\varepsilon).$$

(iv) The Hessian  $-\nabla_\theta^2 \Gamma(\theta) \Big|_{\theta=\theta_0}$  exists and is nonsingular.

**A13** (Directional Gateaux derivative in  $F_0$ ). Let  $H$  be a finite signed measure on  $(\mathbb{R}^d, \mathcal{B})$  with  $\int \|u_{\theta_0}(y)\| \sqrt{f_{\theta_0}(y)} \frac{d|H|(y)}{\sqrt{g(y)}} < \infty$  (e.g.,  $H \ll G$  with density in  $L^2(g)$ , or  $H = \Delta_z$  with  $g(z) > 0$  and finite integrand). The Gateaux derivative of  $S_G$ ,

$$\dot{S}_G(\theta; H) := \lim_{\varepsilon \downarrow 0} \frac{S_{G_\varepsilon}(\theta) - S_G(\theta)}{\varepsilon}$$

exists for  $\theta$  near  $\theta_0$  and

$$\dot{S}_G(\theta; H) = \frac{1}{4} \int u_\theta(y) \sqrt{f_\theta(y)} \frac{dH(y)}{\sqrt{g(y)}}.$$

Based on these assumptions, we derive the influence function [16] and the  $\alpha$ -influence curve to describe the estimator's behavior under small levels of contamination. The proofs of the following theorem and corollary are given in the Appendix C.

**Theorem 3.16** (Influence function). *Under assumptions **A12**–**A13**, the functional  $T$  is Gateaux differentiable at  $G$  in direction  $H$  and has influence function*

$$IF(H; T, G) := \frac{d}{d\varepsilon} T(G_\varepsilon) \Big|_{\varepsilon=0} = -\mathbf{Q}^{-1} \dot{S}_G(\theta_0; H),$$

with

$$\mathbf{Q} := \nabla_\theta S_G(\theta) \Big|_{\theta=\theta_0} = \frac{1}{2} \int \left[ [\nabla_\theta u_\theta(y)] + \frac{1}{2} u_\theta(y) u_\theta(y)^\top \right]_{\theta=\theta_0} \sqrt{f_{\theta_0}(y) g(y)} dy.$$

In particular, for a point mass  $H = \Delta_z$  with  $g(z) > 0$ ,  $IF(z; T, G) = -\mathbf{Q}^{-1} \phi_g(z)$ .

**Corollary 3.17** ( $\alpha$ -influence curve). *For  $G_\varepsilon$  with small  $\varepsilon$ ,*

$$T(G_\varepsilon) = \theta_0 - \varepsilon \mathbf{Q}^{-1} \dot{S}_G(\theta_0; H) + O(\varepsilon^2).$$

*In particular, for point-contamination at  $z$ ,  $H = \Delta_z$ , with  $g(z) > 0$ ,  $T(G_\varepsilon) = \theta_0 - \varepsilon \mathbf{Q}^{-1} \phi_g(z) + O(\varepsilon^2)$ .*

*Remark 3.18.* All statements are made in the Hellinger topology. The influence function holds for any direction  $H$  satisfying Assumption **A13**. In particular, for point-mass contamination  $\Delta_z$ ,  $g(z) > 0$  to avoid division by zero in  $\sqrt{f_{\theta_0}(z)/g(z)}$ . For directions  $H \ll G$  with density  $h = dH/dG \in L^2(g)$ , Assumption **A13** is always satisfied since  $\int \|u_{\theta_0}\| \sqrt{f_{\theta_0}} h / \sqrt{g} dy = \int \|u_{\theta_0}\| \sqrt{f_{\theta_0}/g} h dG$  is finite under the  $L^2(g)$  envelope.

## 4 Empirical Studies

To bring the theoretical properties derived above into perspective and compare with the maximum likelihood estimator, we conduct a large simulation study. We then demonstrate the versatility of the MHDE (2) by applying it in the National Health and Nutrition Examination Survey (NHANES)[10].

### 4.1 Simulation study

We simulate data from a finite population of size  $N \in \{10^6, 10^{6.5}, 10^7, 10^{7.5}, 10^8\}$  with two different sampling ratios  $\alpha \in \{10^{-3}, 10^{-4}\}$ . The characteristic of interest follows a  $\Gamma$  superpopulation model,  $Y \sim \Gamma(2, 35000)$ . For Poisson-PPS we simulate a log-normal auxiliary variable  $Z$  using different correlations with  $Y$ ,  $\rho_{YZ} \in \{0.25, 0.75\}$ . In Section D.1 of the supplementary materials, we present the results with the survey weights calibrated to match known cluster totals. The conclusions from the calibrated survey weights are similar to what is presented here.

To inspect the robustness properties of the MHDE, we introduce point-mass contamination in a fraction of the sampled observations. Specifically, we replace  $\lfloor \varepsilon n \rfloor$  observations in the sample with draws from a Normal distribution with mean  $z \gg \mathbb{E}[Y]$  and variance  $10^{-2} \text{Var}(Y)$ . For “independent contamination,” the contaminated observations are chosen completely at random, while for “high-leverage contamination,” observations with higher sample weight are more likely to be contaminated,  $\mathbb{P}(\text{obs. } i \text{ is contaminated}) \propto (1 - \pi_{\gamma_i})^{-10}$ . The supplementary materials (Section D.1.1) contain results for a scenario where the contamination comes from a truncated  $t$  distribution with 3 degrees of freedom.

For each combination of simulation parameters, we present the relative absolute bias and the

relative root mean square error across  $R = 100$  replications:

$$\text{RelBias} := \frac{1}{\theta_0} \frac{1}{R} \sum_{r=1}^R (\hat{\theta}_r - \theta_0), \quad \text{RelRMSE} := \frac{1}{\theta_0} \sqrt{\frac{1}{R} \sum_{r=1}^R (\hat{\theta}_r - \theta_0)^2}.$$

We compare the MHDE with the weighted maximum likelihood estimate (MLE) for the Gamma model.

#### 4.1.1 Results

Figure 1 shows the relative bias of the MHDE and the MLE in the Gamma superpopulation model as the finite population size and the sample size increase. When  $N$  is sufficiently large, the bias is within  $\pm 1\%$  for each parameter with both MHDE and MLE. As expected, the variance of both estimators also decreases rapidly with increasing finite population size (Figure 2).

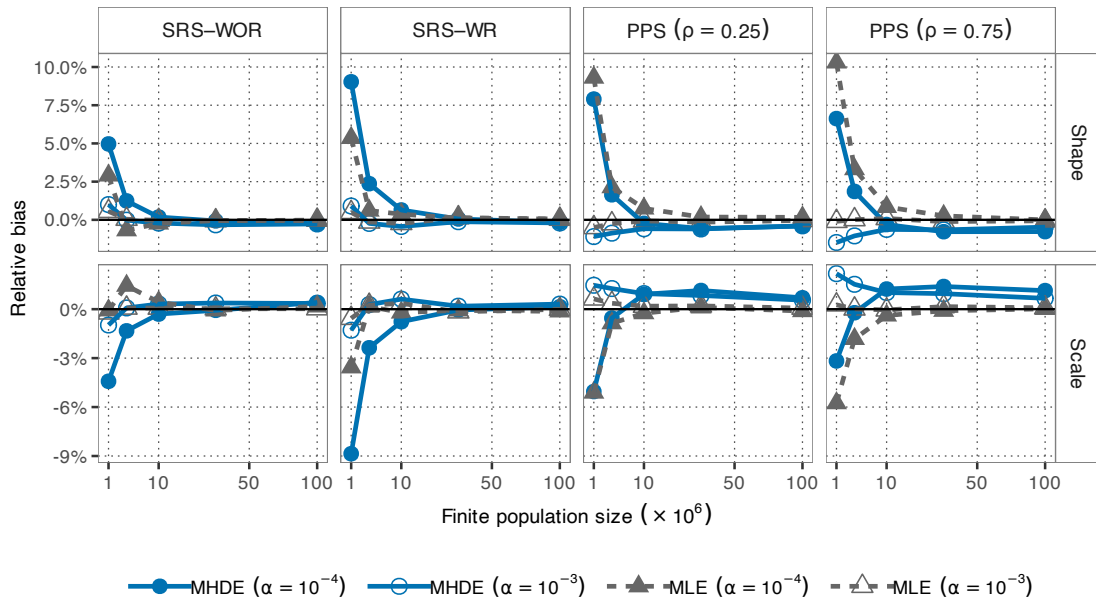
In the presence of contamination, the MHDE clearly shows its advantages over the MLE (Figure 3). Overall, the estimates for the scale parameter of the Gamma superpopulation model are much more affected by contamination than the shape parameter. Importantly, the influence function for the MHDE under independent and high-leverage contamination is bounded, whereas it is unbounded for the MLE. From the  $\alpha$ -influence curve, we can further see that the MHDE can withstand up to 30% high-leverage contamination before becoming unstable. In the presence of independent contamination, on the other hand, the bias of the MHDE remains bounded even when approaching 50% contamination.

We also verify the coverage of the asymptotic confidence intervals derived from Theorem 3.11 using 10 000 replications for  $N = 10^7$  and two sample sizes,  $n \in \{1\,000, 10\,000\}$ . Table 1 summarizes the coverage and width of the 95% confidence intervals for the different sampling strategies. The coverage rate for SRS with and without replacement is very close to the nominal level. For Poisson-PPS, on the other hand, the CI coverage is below the nominal level, likely due to the slightly higher bias observed also in Figure 1. However, this is not unique to the MHDE, but the MLE also suffers from the same issue in this setting.

In Section D.2 of the supplementary materials, we present a second simulation study using the log-normal distribution for  $Y$ . The conclusions align with the Gamma model presented here, but the CI coverage is close to the nominal level for all sampling schemes.

## 4.2 Application to NHANES

We now analyze the total daily water consumption by U.S. residents as collected through the National Health and Nutrition Examination Survey (NHANES) [10]. Over each 2-year period, NHANES surveys health, dietary, sociodemographic and other information from about 10,000 adults and children in the U.S. using several interviews, health assessments and other survey instruments spread over several days. NHANES uses a complex survey design, and calibrated



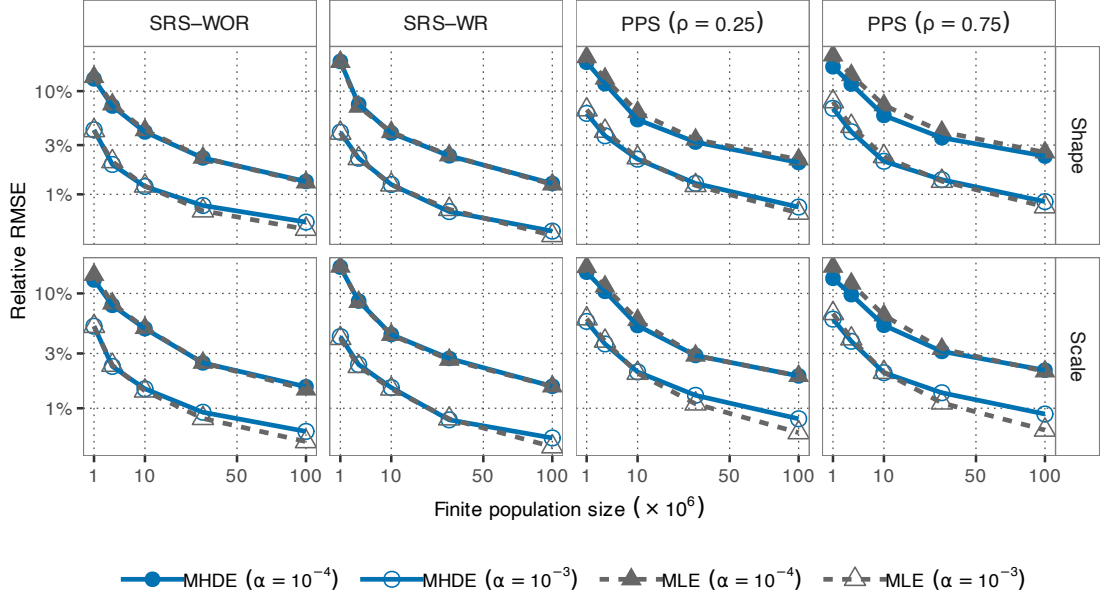
**Figure 1:** Relative bias of the MHDE (blue dots) and MLE (gray triangles) in a Gamma superpopulation model using various sample designs. The sample size in each simulation is determined by  $n = \alpha N$ , with  $\alpha \in \{10^{-3}, 10^{-4}\}$ .

Sample size	Sampling scheme	Shape		Scale	
		CI coverage	CI avg. width	CI coverage	CI avg. width
1 000	PPS ( $\rho = 0.75$ )	91.1%	9.9%	97.5%	11.5%
	SRS-WOR	95.8%	8.2%	94.3%	9.2%
	SRS-WR	95.4%	8.2%	94.0%	9.2%
10 000	PPS ( $\rho = 0.75$ )	86.1%	3.1%	92.6%	3.6%
	SRS-WOR	95.2%	2.6%	94.9%	2.9%
	SRS-WR	95.2%	2.6%	94.8%	2.9%

**Table 1:** Coverage and average (relative) width of 95% confidence intervals across 10 000 replicates of the MHDE estimates in the Gamma model with finite population size  $N = 10^7$ .

survey weights are reported separately for each part of the survey. Here, we analyze the dietary interview data from the 2021–2023 survey cycle, specifically the total daily water consumption. Each NHANES participant is eligible for two 24-hour dietary recall interviews. In the 2021–2023 survey cycle, both interviews were conducted by telephone, as opposed to the first interview being conducted in person as in earlier iterations of NHANES. This may decrease the reliability of the first interview compared to previous years. In fact, three and four respondents reported drinking more than 10 liters a day in the first interview and the second interview, respectively. These values are not only unusual, but can even lead to hyponatremia Adrogué and Madias [17].

We fit a Gamma model and a Weibull model to the survey data, estimating the parameters using the proposed MHDE and the reference MLE. Figure 4 shows the fitted densities for these two models for the second day. In both models, the MLE is shifted rightwards, apparently affected by the few unusually high values. There is a single response of 44.2 liters/day with a



**Figure 2:** Relative RMSE of the MHDE (blue dots) and MLE (gray triangles) under a Gamma superpopulation model using various sample designs. The sample size in each simulation is determined by  $n = \alpha N$ , with  $\alpha \in \{10^{-3}, 10^{-4}\}$ .

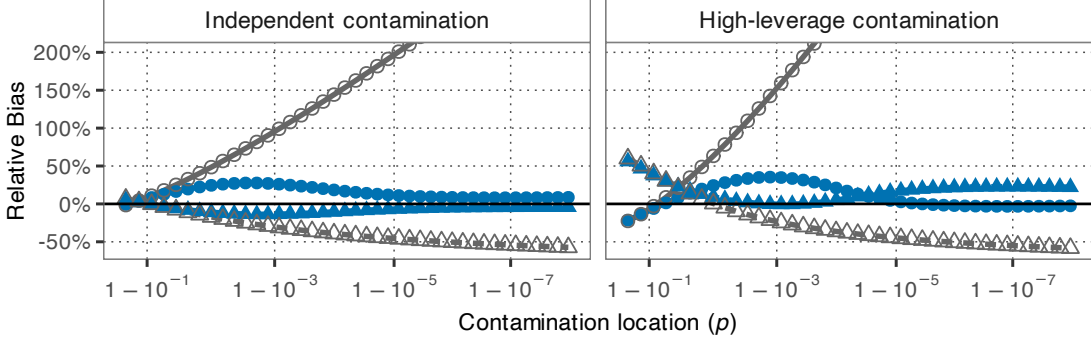
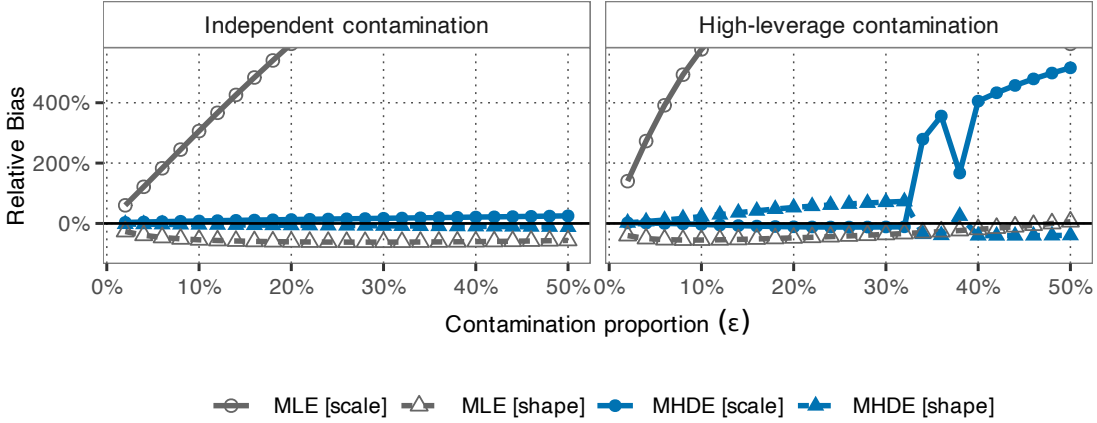
Model	Estimator	Mean [95% CI]	Median [95% CI]
<i>Non-parametric</i>		1.26 [1.21, 1.31]	1.01 [1.00, 1.08]
Gamma	MHDE	1.32 [1.28, 1.36]	0.99 [0.96, 1.02]
	MLE	1.45 [1.41, 1.48]	1.16 [1.13, 1.19]
Weibull	MHDE	1.32 [1.19, 1.35]	1.02 [0.97, 1.17]
	MLE	1.45 [1.37, 1.44]	1.17 [1.20, 1.27]

**Table 2:** Population-level estimates of average total daily water consumption (in liters) using a non-parametric estimator and the MHDE and MLE for two different parametric models. The 95% confidence intervals for the parametric estimates are Monte-Carlo approximations using 10 000 draws from the asymptotic distribution.

sampling weight in the 99th percentile, which can have a devastating effect on the MLE.

In sample surveys, the interest is often in population statistics, like population averages or totals. We can easily obtain these statistics and associated confidence intervals from the fitted superpopulation models. Here, we estimate the effective sample size according to Kish [18] by  $\widehat{n}_{\text{eff}} = (\sum_i w_i)^2 / (\sum_i w_i^2)$  since the inclusion probabilities are unknown. In Table 2 we can again see that the MLE is shifted upward, likely due to the bias from the unreasonable outliers in the data. The non-parametric estimates are computed using the weighted mean and median, with confidence intervals derived from the Taylor series expansion implemented in the survey R package Lumley et al. [19]. In general, the non-parametric estimates seem to agree with the MHDE estimates, with overlapping confidence intervals. The ML estimates, on the other hand, are substantially higher.

(a) Influence function

(b)  $\alpha$ -influence curve

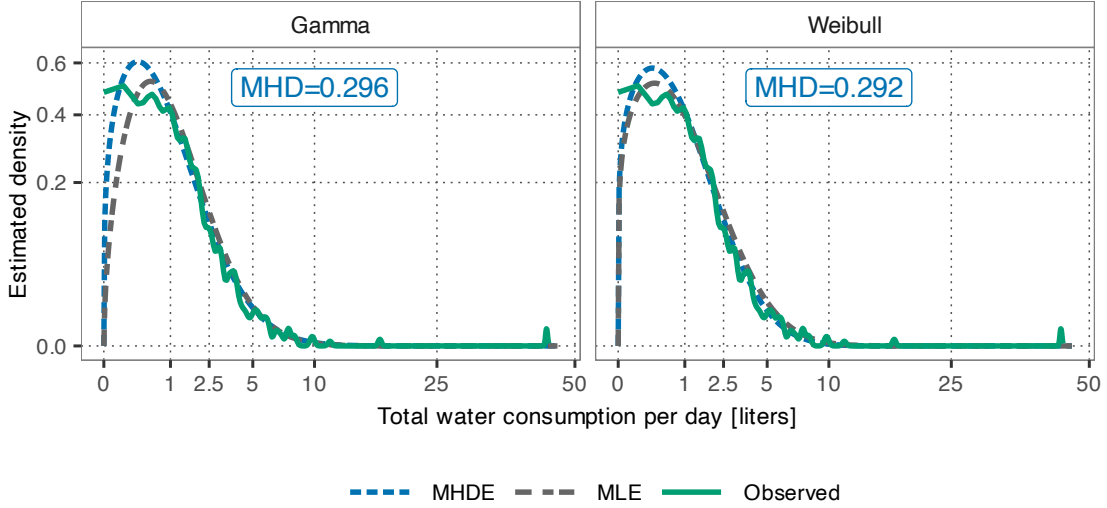
**Figure 3:** Influence functions (top) and alpha curves (bottom) for the MHDE and MLE of the scale ( $\circ$ ) and shape ( $\triangle$ ) parameters in the Gamma model. The contamination proportion in the influence function at the top is set to  $\varepsilon = 0.1$ . The horizontal axis shows the location of the point-mass contamination in terms of the quantile of the true superpopulation model, i.e.,  $z = G^{-1}(p)$ . For the alpha curve the point-mass contamination is located at  $z = G^{-1}(1 - 10^{-7}) \approx 669.193$ .

## 5 Discussion

In this paper, we develop the Minimum Hellinger Distance Estimator (MHDE) with Horvitz-Thompson adjusted kernel density estimator for finite populations under various sampling designs. In the superpopulation framework with potential model misspecification, we prove that the MHDE is consistent in the Hellinger topology and admits an asymptotic Normal distribution, with fully efficient covariance if the true distribution is in the parametric family. We further derive the influence function and the  $\alpha$ -influence curve, showing that the MHDE is highly robust against contamination, including high-leverage points. Our theory requires minimal assumptions on the true density and the sampling design, allowing for efficient estimation and valid inference even under post-stratification or calibration. Hence, the MHDE is as efficient as the MLE if the superpopulation assumption is correct, but much more reliable and stable if the model is misspecified or the sample contaminated.

The MHDE is easy to implement for a wide class of parametric families with minimal ad-





**Figure 4:** MHDE and MLE estimates for two different parametric models to describe the total daily water consumption in the NHANES survey. The minimum Hellinger distance (MHD) is achieved by the MHDE shown here.

justments. In the numerical experiments, we applied the MHDE for the Gamma and Weibull models, but other models, such as the log-normal, are equally straightforward to implement. The numerical experiments further underscore the utility of the MHDE in complex survey samples, particularly its stability under contamination and its versatility.

The simplicity and efficiency our HT-adjusted MHDE make it an ideal candidate for complex survey samples and a wide range of superpopulation models. While the focus in this paper is on the Hellinger distance, the techniques used in the proofs are expressively more general. With appropriate adjustments to the assumptions, our results can be generalized to broader classes of divergences, such as power divergences [20] or  $\phi$  divergences Pardo [21]. A better understanding of the theoretical properties of more general HT-adjusted minimum divergence estimators is crucial to choosing the best estimator under different sampling strategies and contamination expectations.

## Appendix A Proof of Consistency

### A.1 Technical Lemmas

**Lemma A.1** (Self-normalization). *Let*

$$\hat{f}_\gamma(y) = \frac{\sum_{i \in \mathcal{U}_\gamma} \frac{\delta_{\gamma i}}{\pi_{\gamma i}} K_{h_\gamma}(y - Y_{\gamma i})}{\sum_{i \in \mathcal{U}_\gamma} \frac{\delta_{\gamma i}}{\pi_{\gamma i}}} = \frac{T_\gamma(y)}{S_\gamma}.$$

*If  $S_\gamma > 0$ , then  $\int_{\mathbb{R}^d} \hat{f}_{HT,\gamma}(y) dy = 1$  almost surely.*

*Proof.* By Fubini and the change of variables  $u = (y - Y_{\gamma i})/h_{\gamma}$ ,

$$\begin{aligned} \int T_{\gamma}(y) dy &= \frac{1}{N_{\gamma}} \sum_{i=1}^{N_{\gamma}} \frac{\delta_{\gamma i}}{\pi_{\gamma i}} \int K_{h_{\gamma}}(y - Y_{\gamma i}) dy \\ &= \frac{1}{N_{\gamma}} \sum_{i=1}^{N_{\gamma}} \frac{\delta_{\gamma i}}{\pi_{\gamma i}} \int K(u) du = \frac{1}{N_{\gamma}} \sum_{i=1}^{N_{\gamma}} \frac{\delta_{\gamma i}}{\pi_{\gamma i}} = S_{\gamma}, \end{aligned}$$

since  $K$  integrates to one by assumption **A1**. Hence  $\int \hat{f}_{\gamma} = S_{\gamma}/S_{\gamma} = 1$  on  $\{S_{\gamma} > 0\}$ . Under Poisson-PPS with  $n_{\gamma} = \sum_i \pi_{\gamma i} \rightarrow \infty$ ,  $\mathbb{P}(S_{\gamma} = 0) \leq e^{-n_{\gamma}} \rightarrow 0$ ; under sampling without replacement,  $S_{\gamma} > 0$  deterministically.  $\square$

**Lemma A.2** (Bernstein, independent case). *Let  $X_1, \dots, X_m$  be independent,  $\mathbb{E}X_i = 0$ ,  $|X_i| \leq b$ , and  $\sum_{i=1}^m \text{Var}(X_i) \leq v$ . Then for all  $t > 0$ ,*

$$\mathbb{P}\left(\sum_{i=1}^m X_i \geq t\right) \leq \exp\left(-\frac{t^2}{2(v + bt/3)}\right), \quad \mathbb{P}\left(\left|\sum_{i=1}^m X_i\right| \geq t\right) \leq 2 \exp\left(-\frac{t^2}{2(v + bt/3)}\right).$$

The next lemma is a variant as applied to simple random sampling without replacement (SRSWOR).

**Lemma A.3** (Bernstein under WOR (Serfling-type)). *Let a finite population  $\{y_1, \dots, y_N\} \subset \mathbb{R}$  have mean  $\mu = \frac{1}{N} \sum_{i=1}^N y_i$ , range  $|y_i - \mu| \leq B$ , and population variance  $\sigma_N^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \mu)^2$ . Draw a sample of size  $m$  without replacement and let  $X_1, \dots, X_m$  be the sampled values in any order. Then, with  $v_{\text{wor}} := \frac{N-m}{N-1} m \sigma_N^2$ , for all  $t > 0$ ,*

$$\mathbb{P}\left(\sum_{k=1}^m (X_k - \mu) \geq t\right) < \exp\left(-\frac{t^2}{2(v_{\text{wor}} + Bt/3)}\right).$$

*Equivalently, compared to the independent (with-replacement) Bernstein bound with variance proxy  $v_{\text{ind}} := m \sigma_N^2$ , the WOR bound holds with the finite-population correction  $v_{\text{wor}} = \frac{N-m}{N-1} v_{\text{ind}}$ .*

*Proof.* Write  $Z_k := X_k - \mu$ . Let  $\mathcal{F}_k$  be the  $\sigma$ -field generated by the first  $k$  draws, and consider the Doob (Hájek) decomposition

$$M_k := \sum_{r=1}^k (Z_r - \mathbb{E}[Z_r \mid \mathcal{F}_{r-1}]), \quad k = 0, 1, \dots, m,$$

with  $M_0 = 0$ . Then  $(M_k)_{k \leq m}$  is a martingale and  $M_m = \sum_{k=1}^m Z_k$  because  $\sum_{k=1}^m \mathbb{E}[Z_k \mid \mathcal{F}_{k-1}] = 0$  (the remaining population is always centered around its current mean and those means telescope to 0; see Remark A.4 below).

**Bounded increments.** Each increment is bounded as

$$|M_k - M_{k-1}| = |Z_k - \mathbb{E}[Z_k \mid \mathcal{F}_{k-1}]| \leq |Z_k| + |\mathbb{E}[Z_k \mid \mathcal{F}_{k-1}]| \leq 2B =: b.$$

**Predictable quadratic variation.** Define the predictable quadratic variation

$$V_m := \sum_{k=1}^m \mathbb{E}[(M_k - M_{k-1})^2 \mid \mathcal{F}_{k-1}] = \sum_{k=1}^m \text{Var}(Z_k \mid \mathcal{F}_{k-1}).$$

Taking expectations and using the variance decomposition for martingales gives

$$\mathbb{E}[V_m] = \text{Var}\left(\sum_{k=1}^m Z_k\right).$$

Under simple random sampling without replacement, the variance of the sample *sum* is the classical finite-population formula

$$\text{Var}\left(\sum_{k=1}^m X_k\right) = \frac{N-m}{N-1} m \sigma_N^2,$$

hence  $\mathbb{E}[V_m] = v_{\text{wor}}$  as claimed.

**Freedman's inequality and optimization.** Freedman's inequality for martingales with bounded increments (e.g., Theorem 1.6 in Freedman [22]) states that for all  $t, v > 0$ ,

$$\mathbb{P}(M_m \geq t \cap V_m \leq v) \leq \exp\left\{-\frac{t^2}{2(v + bt/3)}\right\}.$$

We use the standard peeling argument on the random  $V_m$ :

$$\begin{aligned} \mathbb{P}(M_m \geq t) &= \sum_{j \geq 0} \mathbb{P}(M_m \geq t, V_m \in (2^{j-1}v_{\text{wor}}, 2^j v_{\text{wor}}]) \\ &\leq \sum_{j \geq 0} \exp\left\{-\frac{t^2}{2(2^j v_{\text{wor}} + bt/3)}\right\} \\ &\leq \exp\left\{-\frac{t^2}{2(v_{\text{wor}} + bt/3)}\right\}, \end{aligned}$$

where the last inequality uses that the series is dominated by its first term (geometric decay in  $j$  once  $t$  is fixed). Substituting  $b = 2B$  and  $v_{\text{wor}}$  yields the stated bound. The two-sided tail follows by symmetry. □

*Remark A.4* (Centering under WOR). At step  $k$ , conditional on  $\mathcal{F}_{k-1}$ ,  $X_k$  is uniformly distributed over the remaining  $N - k + 1$  units; its conditional mean equals the mean of the remaining values, which is  $-\frac{1}{N-k+1} \sum_{r=1}^{k-1} Z_r$ . Consequently,  $\sum_{k=1}^m \mathbb{E}[Z_k \mid \mathcal{F}_{k-1}] = 0$ .

*Remark A.5* (Asymptotics and the  $(1-\alpha)$  factor). Since  $\frac{N-m}{N-1} = (1-\alpha)(1-\frac{1}{N})^{-1}$  with  $\alpha = m/N$ , the variance proxy satisfies  $v_{\text{wor}} = (1-\alpha + o(1))m\sigma_N^2$  as  $N \rightarrow \infty$ . Thus, in triangular-array asymptotics with  $N \rightarrow \infty$ , replacing  $v$  by  $(1-\alpha)v$  in the independent Bernstein bound is correct up to a vanishing factor.

**Lemma A.6** (HT normalizer concentration). *Let  $S_\gamma = \frac{1}{N_\gamma} \sum_{i=1}^{N_\gamma} \frac{\delta_{\gamma i}}{\pi_{\gamma i}}$ . Under Poisson-PPS,*

$$\mathbb{P}(|S_\gamma - 1| > t) \leq 2 \exp \left( - \frac{N_\gamma^2 t^2}{2 \left( \sum_{i \in \mathcal{U}_\gamma} \frac{1-\pi_{\gamma i}}{\pi_{\gamma i}} + \frac{t N_\gamma}{3} \max_i \frac{1}{\pi_{\gamma i}} \right)} \right) \leq 2 \exp \left( - \frac{cn_{\text{eff},\gamma} t^2}{1 + c't} \right),$$

for constants  $c, c' > 0$  under the regularity  $\max_i \pi_{\gamma i}^{-1} \lesssim 1/\alpha_\gamma$ . For WOR, multiply the denominator's variance term by  $(1-\alpha_\gamma)$ .

*Proof.* Write  $S_\gamma - 1 = N_\gamma^{-1} \sum_i X_i$  with  $X_i = (\delta_{\gamma i}/\pi_{\gamma i} - 1)$ , so  $\mathbb{E}X_i = 0$ ,  $|X_i| \leq \max_i \pi_{\gamma i}^{-1}$ ,  $\sum \text{Var}(X_i) = \sum (1-\pi_{\gamma i})/\pi_{\gamma i}$ . Apply Lemma A.2 with  $t \leftarrow N_\gamma t$ ; for WOR use Lemma A.3.  $\square$

The next lemma is useful in establishing the first step of the proof of the Theorem 3.1.

**Lemma A.7** (Three-way  $L_1$  decomposition). *Let*

$$T_\gamma(y) = \frac{1}{N_\gamma} \sum_{i=1}^{N_\gamma} \frac{\delta_{\gamma i}}{\pi_{\gamma i}} K_{h_\gamma}(y - Y_{\gamma i}), \quad S_\gamma = \frac{1}{N_\gamma} \sum_{i=1}^{N_\gamma} \frac{\delta_{\gamma i}}{\pi_{\gamma i}}, \quad \hat{f}_\gamma = \frac{T_\gamma}{S_\gamma},$$

and let  $\bar{f}_{\gamma,h}(y) = \frac{1}{N_\gamma} \sum_{i \in \mathcal{U}_\gamma} K_{h_\gamma}(y - Y_{\gamma i})$  be the i.i.d. kernel average. Then

$$\|\hat{f}_\gamma - f\|_{L^1} \leq |S_\gamma^{-1} - 1| + \|T_\gamma - \bar{f}_{\gamma,h}\|_{L^1} + \|\bar{f}_{\gamma,h} - f\|_{L^1}. \quad (5)$$

*Proof.* Add and subtract  $\bar{f}_{\gamma,h}$  and use the triangle inequality:

$$\|\hat{f}_\gamma - f\|_1 \leq \|\hat{f}_\gamma - \bar{f}_{\gamma,h}\|_1 + \|\bar{f}_{\gamma,h} - f\|_1.$$

For the first term,

$$\hat{f}_\gamma - \bar{f}_{\gamma,h} = \frac{T_\gamma}{S_\gamma} - \bar{f}_{\gamma,h} = \left( \frac{1}{S_\gamma} - 1 \right) \bar{f}_{\gamma,h} + \frac{1}{S_\gamma} (T_\gamma - \bar{f}_{\gamma,h}),$$

so by subadditivity of  $\|\cdot\|_1$ ,

$$\|\hat{f}_\gamma - \bar{f}_{\gamma,h}\|_1 \leq \left| \frac{1}{S_\gamma} - 1 \right| \|\bar{f}_{\gamma,h}\|_1 + \frac{1}{S_\gamma} \|T_\gamma - \bar{f}_{\gamma,h}\|_1.$$

Now  $\|\bar{f}_{\gamma,h}\|_1 = \int \bar{f}_{\gamma,h}(y) dy = \frac{1}{N_\gamma} \sum_i \int K_{h_\gamma}(y - Y_{\gamma i}) dy = 1$  because  $\int K = 1$ . Also,  $S_\gamma > 0$  with probability  $1 - o(1)$  (Poisson-PPS) or deterministically (WOR), and on  $\{S_\gamma > 0\}$  we may write

$$\frac{1}{S_\gamma} \leq 1 + \left| \frac{1}{S_\gamma} - 1 \right|.$$

Hence,

$$\|\hat{f}_\gamma - \bar{f}_{\gamma,h}\|_1 \leq \left| \frac{1}{S_\gamma} - 1 \right| + \left( 1 + \left| \frac{1}{S_\gamma} - 1 \right| \right) \|T_\gamma - \bar{f}_{\gamma,h}\|_1.$$

Finally, since  $|1/S_\gamma - 1| \|T_\gamma - \bar{f}_{\gamma,h}\|_1 \geq 0$ , we may drop that product term to obtain the simpler bound

$$\|\hat{f}_\gamma - \bar{f}_{\gamma,h}\|_1 \leq \left| \frac{1}{S_\gamma} - 1 \right| + \|T_\gamma - \bar{f}_{\gamma,h}\|_1.$$

Combine with the first display to conclude (5).  $\square$

## A.2 Large-deviation bounds for the design term

To prove Proposition 3.2 we need the following technical lemmas. Throughout this proof we use the definition of  $\bar{f}_{\gamma,h}$  from Proposition 3.2, set  $\xi_{\gamma i} := \delta_{\gamma i}/\pi_{\gamma i} - 1$  and define the signed measure

$$\mu_\gamma := \frac{1}{N_\gamma} \sum_{i \in \mathcal{U}_\gamma} \xi_{\gamma i} \delta_{Y_{\gamma i}},$$

such that  $T_\gamma - \bar{f}_{\gamma,h} = K_{h_\gamma} * \mu_\gamma$ . We further partition  $\mathbb{R}^d$  into half-open cubes  $\{B_{\gamma j}\}_{j=1}^{M_\gamma}$  with sides of length  $h_\gamma$  and centers  $c_{\gamma j}$ .

**Lemma A.8** (Cellwise smoothing reduction). *Under assumption **A1** and defining*

$$S_\gamma(y) := (T_\gamma - \bar{f}_{\gamma,h})(y) = (K_{h_\gamma} * \mu_\gamma)(y),$$

*there exist constants  $C_K, c, C' > 0$  (depending only on  $K$  and  $d$ ) such that*

$$\int_{\mathbb{R}^d} |S_\gamma(y)| dy \leq C_K \sum_{j=1}^{M_\gamma} |\mu_\gamma(B_{\gamma j})| + R_\gamma, \text{ and}$$

$$\mathbb{P}(R_\gamma > \tau) \leq C' e^{-c N_\gamma h_\gamma^d \min(\tau^2, \tau)}.$$

*Proof sketch.* Decompose  $\mu_\gamma$  into its restrictions on the cells and replace each atom in  $B_{\gamma j}$  by a single atom at  $c_{\gamma j}$ ; integrate the Lipschitz translation error of  $K_{h_\gamma}$  over each cell. The sum of these errors concentrates with  $N_\gamma h_\gamma^d$  by Bernstein applied to i.i.d. occupancies of the cells (details as in Lemma A.2).  $\square$

**Lemma A.9** (Convolution reduction on a grid). *Under assumptions **A1** and **A2**, for any finite signed measure  $\mu$ ,*

$$\|K_h * \mu\|_{L^1} \leq \|K\|_{L^1} \sum_j |\mu(B_j)| + C_K \sum_j |r_j|(B_j),$$

$$\text{with } r_j := \mu \upharpoonright_{B_j} - \mu(B_j) \delta_{c_j},$$

where one may take  $C_K := \|\nabla K\|_{L^1} \frac{\sqrt{d}}{2}$  (so  $C_K$  does **not** depend on  $h$ ). In particular, since  $|r_j|(B_j) \leq 2|\mu|(B_j)$ ,

$$\|K_h * \mu\|_{L^1} \leq \|K\|_{L^1} \sum_j |\mu(B_j)| + 2C_K \sum_j |\mu|(B_j).$$

*Proof.* Write  $\mu = \sum_j \nu_j$  with  $\nu_j = \mu \upharpoonright_{B_j}$  and decompose  $\nu_j = \mu(B_j) \delta_{c_j} + r_j$  with  $r_j(B_j) = 0$ . Then

$$K_h * \mu = \sum_j \mu(B_j) K_h(\cdot - c_j) + \sum_j K_h * r_j.$$

Taking  $L^1$  norms gives the first term as  $\|K\|_{L^1} \sum_j |\mu(B_j)|$ . For the remainder, by Minkowski and the translation inequality valid for  $K \in W^{1,1}$ ,

$$\begin{aligned} \|K_h * r_j\|_{L^1} &\leq \int_{B_j} \|K_h(\cdot - t) - K_h(\cdot - c_j)\|_{L^1} |dr_j|(t) \\ &\leq \left( \sup_{t \in B_j} \|K_h(\cdot - t) - K_h(\cdot - c_j)\|_{L^1} \right) |r_j|(B_j). \end{aligned}$$

Now  $\|K_h(\cdot - t) - K_h(\cdot - c_j)\|_{L^1} \leq \|\nabla K\|_{L^1} \|t - c_j\|/h \leq \|\nabla K\|_{L^1} (\sqrt{d}/2)$  because  $\|t - c_j\| \leq (\sqrt{d}/2)h$  for  $t \in B_j$ . Summing over  $j$  yields the claim.  $\square$

**Lemma A.10** (Per-cell Bernstein/Serfling bound). *Let  $A_{\gamma ij} := \mathbf{1}\{Y_{\gamma i} \in B_{\gamma j}\}$ . Conditional on  $\{Y, Z\}$ ,*

$$\mu_\gamma(B_{\gamma j}) = \frac{1}{N_\gamma} \sum_{i \in \mathcal{U}_\gamma} \xi_{\gamma i} A_{\gamma ij}, \quad \mathbb{E}[\mu_\gamma(B_{\gamma j}) \mid \{Y, Z\}] = 0,$$

and

$$\text{Var}(\mu_\gamma(B_{\gamma j}) \mid \{Y, Z\}) = \frac{1}{N_\gamma^2} \sum_{i \in \mathcal{U}_\gamma} \frac{1 - \pi_{\gamma i}}{\pi_{\gamma i}} A_{\gamma ij} \leq \frac{1}{n_{\text{eff}, \gamma}}.$$

Moreover, if  $\max_i \pi_{\gamma i}^{-1} \leq c_0/\alpha_\gamma$ , then for all  $t > 0$ ,

$$\mathbb{P}(|\mu_\gamma(B_{\gamma j})| > t \mid \{Y, Z\}) \leq 2 \exp \left\{ -\frac{t^2}{2(v_\gamma + b_\gamma t/3)} \right\},$$

with  $v_\gamma \leq n_{\text{eff}, \gamma}^{-1}$  and  $b_\gamma \leq c_0/n_\gamma$ . Under sampling without replacement (rejective), replace  $v_\gamma$  by  $(1 - \alpha_\gamma)n_{\text{eff}, \gamma}^{-1}$ .

*Proof.* The variance identity follows immediately from independence (Poisson-PPS) of  $\delta_{\gamma i}$  given  $\{Y, Z\}$ ; the tail bound is Bernstein's inequality with the stated  $v_\gamma, b_\gamma$  (and Lemma A.3 for WOR).  $\square$

*Proof of Proposition 3.2.* By Lemma A.8,  $\|T_\gamma - \bar{f}_{\gamma, h}\|_1 \leq C_K \sum_{j=1}^{M_\gamma} |\mu_\gamma(B_{\gamma j})| + R_\gamma$ . Fix  $\tau \in (0, 1]$

and set  $u := \tau/(2C_K)$ ; then

$$\mathbb{P} \left( C_K \sum_j |\mu_\gamma(B_{\gamma j})| > \frac{\tau}{2} \mid \{Y, Z\} \right) \leq \mathbb{P} \left( \exists j : |\mu_\gamma(B_{\gamma j})| > \frac{u}{M_\gamma} \mid \{Y, Z\} \right).$$

By Lemma A.10 with  $t = u/M_\gamma$  and a union bound over  $M_\gamma \asymp h_\gamma^{-d}$  cells,

$$\begin{aligned} \mathbb{P} \left( \exists j : |\mu_\gamma(B_{\gamma j})| > \frac{u}{M_\gamma} \mid \{Y, Z\} \right) &\leq 2M_\gamma \exp \left\{ -\frac{t^2}{2(v_\gamma + b_\gamma t/3)} \right\} \\ &\leq C \exp \left\{ -c \frac{n_{\text{eff}, \gamma}}{M_\gamma} \min\{u^2, u\} \right\}, \end{aligned}$$

since  $v_\gamma \leq n_{\text{eff}, \gamma}^{-1}$  and  $b_\gamma \leq c_0/n_\gamma$ , and for large  $\gamma$  the  $b_\gamma t$  term is dominated by  $v_\gamma$  when  $t \lesssim M_\gamma^{-1}$ . Because  $M_\gamma \asymp h_\gamma^{-d}$ , we obtain

$$\mathbb{P} \left( C_K \sum_j |\mu_\gamma(B_{\gamma j})| > \frac{\tau}{2} \mid \{Y, Z\} \right) \leq C \exp \left\{ -c n_{\text{eff}, \gamma} h_\gamma^d \min(\tau^2, \tau) \right\}.$$

Finally, add the bound for the remainder  $\mathbb{P}(R_\gamma > \tau/2) \leq C e^{-c N_\gamma h_\gamma^d}$  from Lemma A.8. The WOR version follows by replacing  $v_\gamma$  by  $(1 - \alpha_\gamma) n_{\text{eff}, \gamma}^{-1}$  in Lemma A.10.  $\square$

### A.3 KDE large-deviation bounds

**Lemma A.11** (i.i.d. KDE  $L^1$  tail). *Under assumptions **A1** and **A3**, there exist constants  $C, c > 0$  (depending only on  $K, d$ ) such that for all  $\tau \in (0, 1]$ ,*

$$\mathbb{P} \left( \|\bar{f}_{\gamma, h} - g * K_{h_\gamma}\|_{L^1} > \tau \right) \leq C \exp \left\{ -c N_\gamma h_\gamma^d \min(\tau^2, \tau) \right\}. \quad (6)$$

*Proof sketch.* Partition  $\mathbb{R}^d$  into cubes  $\{B_{\gamma j}\}_{j=1}^{M_\gamma}$  of side  $h_\gamma$  with centers  $z_{\gamma j}$ , where  $M_\gamma \asymp h_\gamma^{-d}$ . For each cell center,

$$S_{\gamma j} := \frac{1}{N_\gamma} \sum_{i=1}^{N_\gamma} \{K_{h_\gamma}(z_{\gamma j} - Y_{\gamma i}) - \mathbb{E}[K_{h_\gamma}(z_{\gamma j} - Y)]\}$$

is a sum of independent centered bounded variables with variance  $\lesssim (N_\gamma h_\gamma^d)^{-1}$ ; Bernstein yields  $\mathbb{P}(|S_{\gamma j}| > u) \leq 2 \exp\{-c N_\gamma h_\gamma^d \min(u^2, u)\}$ . A union bound over  $M_\gamma \asymp h_\gamma^{-d}$  centers gives  $\max_j |S_{\gamma j}| \leq u$  with probability at least  $1 - C \exp\{-c N_\gamma h_\gamma^d \min(u^2, u)\}$ . Using the Lipschitz–translation inequality for  $K_{h_\gamma}$ , which is valid since  $K \in W^{1,1}$ ,

$$\int_{B_{\gamma j}} |(\bar{f}_{\gamma, h} - g * K_{h_\gamma})(y) - S_{\gamma j}| \, dy \leq C_K h_\gamma^d.$$

Summing over  $j$  yields  $\|\bar{f}_{\gamma, h} - g * K_{h_\gamma}\|_{L^1} \leq M_\gamma h_\gamma^d \max_j |S_{\gamma j}| + C_K M_\gamma h_\gamma^d = \max_j |S_{\gamma j}| + C'_K$ .

Choosing  $u \simeq \tau/2$  and noticing that  $C'_K$  can be absorbed into the  $\min(\tau^2, \tau)$  regime for  $\tau \in (0, 1]$  gives the desired bound (6). A full proof parallels the one of Proposition 3.2, with independence replacing design weighting.  $\square$

**Lemma A.12** (Approximate identity). *If  $K \in L^1(\mathbb{R}^d)$  with  $\int K = 1$ , then for every  $g \in L^1(\mathbb{R}^d)$ ,  $\|g * K_{h_\gamma} - g\|_{L^1} \rightarrow 0$  as  $h_\gamma \downarrow 0$ .*

## A.4 Consistency of the HT-adjusted KDE

*Proof of Theorem 3.1.* The proof uses several technical lemmas listed in Appendix A.1. With the definition of  $\bar{f}_{\gamma,h}(y)$  from Proposition 3.2 and using Lemma A.7 we get

$$\|\hat{f}_\gamma - g\|_1 \leq \underbrace{|S_\gamma^{-1} - 1|}_{A_\gamma} + \underbrace{\|T_\gamma - \bar{f}_{\gamma,h}\|_1}_{B_\gamma} + \underbrace{\|\bar{f}_{\gamma,h} - g\|_1}_{C_\gamma}. \quad (7)$$

**HT normalizer concentration.** Let  $X_{\gamma i} := N_\gamma^{-1}(\delta_{\gamma i}/\pi_{\gamma i} - 1)$ . Conditional on  $Z$ , the  $X_{\gamma i}$  are independent, mean zero, and  $|X_{\gamma i}| \leq c_0/n_\gamma$  by Assumption A4. Moreover

$$\text{Var}(S_\gamma \mid Z) = \frac{1}{n_{\text{V-eff},\gamma}} \leq \frac{1}{n_{\text{eff},\gamma}}.$$

Bernstein bounds from Lemmas A.2–A.3 yield constants  $c, c' > 0$  with

$$\mathbb{P}(|S_\gamma - 1| > t \mid Z) \leq 2 \exp \left\{ -\frac{cn_{\text{eff},\gamma}t^2}{1 + c't} \right\}.$$

For  $|S_\gamma - 1| \leq 1/2$ ,  $A_\gamma \leq 2|S_\gamma - 1|$ , hence  $A_\gamma = o_{\mathbb{P}}(1)$  with exponential tails.

**Design noise in the numerator.** Proposition 3.2 shows there exist constants  $c, C > 0$  such that

$$\mathbb{P}(B_\gamma > \tau \mid \{Y, Z\}) \leq C \exp \left\{ -cn_{\text{eff},\gamma}h_\gamma^d \min(\tau^2, \tau) \right\} + C \exp \left\{ -cN_\gamma h_\gamma^d \right\}.$$

**Smoothing noise and bias.** Decompose

$$C_\gamma \leq \|\bar{f}_{\gamma,h} - g * K_{h_\gamma}\|_{L^1} + \|g * K_{h_\gamma} - g\|_{L^1}.$$

Then, combining Lemmas A.11 and A.12 shows that  $C_\gamma \rightarrow 0$  in probability, with the tail

$$\mathbb{P}(\|\bar{f}_{\gamma,h} - g * K_{h_\gamma}\|_{L^1} > \tau) \leq C \exp \left\{ -cN_\gamma h_\gamma^d \min(\tau^2, \tau) \right\}.$$

for the stochastic part.



**Conclusion** Plugging the three steps above into the three-way decomposition (7) we obtain

$$\mathbb{P}\left(\|\hat{f}_\gamma - g\|_{L^1} > \tau\right) \leq C \exp\{-c n_{\text{eff},\gamma} h_\gamma^d \min(\tau^2, \tau)\} + C \exp\{-c N_\gamma h_\gamma^d\} + o(1),$$

and therefore  $\|\hat{f}_\gamma - g\|_{L^1} \xrightarrow{p} 0$  as  $\gamma \rightarrow \infty$  whenever  $n_{\text{eff},\gamma} h_\gamma^d \rightarrow \infty$ .  $\square$

*Remark A.13* (On the  $N_\gamma h_\gamma^d$  growth). Note that

$$n_{\text{eff},\gamma} = \frac{N_\gamma^2}{\sum_{i \in \mathcal{U}_\gamma} \pi_{\gamma i}^{-1}} \leq \frac{N_\gamma^2}{N_\gamma^2 / \sum_{i \in \mathcal{U}_\gamma} \pi_{\gamma i}} = \sum_{i \in \mathcal{U}_\gamma} \pi_{\gamma i} =: n_\gamma \leq N_\gamma,$$

so  $N_\gamma h_\gamma^d \geq n_{\text{eff},\gamma} h_\gamma^d$  for all  $\gamma$ . Therefore, the condition  $n_{\text{eff},\gamma} h_\gamma^d \rightarrow \infty$  implies  $N_\gamma h_\gamma^d \rightarrow \infty$ , without any additional assumptions on  $\alpha_\gamma$ .

## A.5 Consistency of the MHDE

**Lemma A.14** (Uniform Hellinger control). *For all  $\gamma$ ,*

$$\sup_{\theta \in \Theta} |\Gamma_\gamma(\theta) - \Gamma(\theta)| \leq \left\| \sqrt{\hat{f}_\gamma} - \sqrt{g} \right\|_{L^2}.$$

*Proof.* For any  $\theta$ , by Cauchy-Schwarz,

$$|\Gamma_\gamma(\theta) - \Gamma(\theta)| = \left| \int \sqrt{f_\theta} \left( \sqrt{\hat{f}_\gamma} - \sqrt{g} \right) \right| \leq \left\| \sqrt{f_\theta} \right\|_2 \left\| \sqrt{\hat{f}_\gamma} - \sqrt{g} \right\|_2 = \left\| \sqrt{\hat{f}_\gamma} - \sqrt{g} \right\|_2,$$

since  $\|\sqrt{f_\theta}\|_2 = (\int f_\theta)^{1/2} = 1$ . Taking the supremum over  $\theta$  gives the claim.  $\square$

**Lemma A.15** (Hellinger vs.  $L_1$ ). *For densities  $p, q$  on  $\mathbb{R}^d$ ,*

$$\|\sqrt{p} - \sqrt{q}\|_{L^2}^2 \leq \|p - q\|_{L^1}.$$

*Proof.* Pointwise for  $a, b \geq 0$  and w.l.o.g.  $a \geq b$ ,  $(\sqrt{a} - \sqrt{b})^2 = (a - b)/(\sqrt{a} + \sqrt{b}) \leq a - b$ . Integrate with  $a = p(y)$  and  $b = q(y)$ .  $\square$

*Proof of Proposition 3.5.* By Lemmas A.14 and A.15,

$$\sup_{\theta} |\Gamma_\gamma(\theta) - \Gamma(\theta)| \leq \left\| \sqrt{\hat{f}_\gamma} - \sqrt{g} \right\|_2 \leq \left\| \hat{f}_\gamma - g \right\|_1^{1/2}.$$

Hence, for  $t \in (0, 1]$ ,

$$\mathbb{P}\left(\sup_{\theta} |\Gamma_\gamma(\theta) - \Gamma(\theta)| > t\right) \leq \mathbb{P}\left(\|\hat{f}_\gamma - g\|_1 > t^2\right).$$

Apply the  $L_1$  large-deviation bound of Theorem 3.1 with  $\tau = t^2$ , which yields  $\exp\{-c n_{\text{eff},\gamma} h_\gamma^d \min(t^4, t^2)\}$  for the design term and  $\exp\{-c N_\gamma h_\gamma^d\}$  for the i.i.d. smoothing term. The WOR factor  $(1 - \alpha_\gamma)$

enters as in Theorem 3.1. □

## Appendix B Proof of the CLT

We first define the notation used throughout this proof. For any distribution  $H$  with density  $h$  we write

$$\nabla_{\theta} \Gamma_h(\theta) \Big|_{\theta=\theta_0} = \int \phi_g(y) (h(y) - g(y)) R_h(y) dy, \quad R_h(y) := \frac{2\sqrt{g(y)}}{\sqrt{h(y)} + \sqrt{g(y)}} \in [0, 2].$$

Since  $\theta_0$  maximizes  $\Gamma_G(\theta) = \int \sqrt{f_{\theta}} \sqrt{g}$ , we have  $\nabla_{\theta} \Gamma_G(\theta_0) = 0$ .

Let's further define  $\psi_{h_{\gamma}} := K_{h_{\gamma}} * \psi_g$ . We continue to use the notation  $\hat{f}_{\gamma} = S_{\gamma}^{-1} T_{\gamma}$  for the HT-adjusted KDE and  $\bar{f}_{\gamma,h}(y)$  for the unweighted KDE as in Proposition 3.2.

**Algebraic decomposition.** Add and subtract  $\bar{f}_{\gamma,h}$  and  $g * K_{h_{\gamma}}$ , and isolate the normalizer:

$$\begin{aligned} \nabla_{\theta} \Gamma_{\gamma}(\theta_0) &= \int \phi_g (S_{\gamma}^{-1} T_{\gamma} - g) R_{\gamma} = \underbrace{\int \phi_g (T_{\gamma} - \bar{f}_{\gamma,h}) R_{\gamma}}_{\mathbf{A}_{\gamma,1}} + \underbrace{\int \phi_g (\bar{f}_{\gamma,h} - g * K_{h_{\gamma}}) R_{\gamma}}_{\mathbf{A}_{\gamma,2}} \\ &\quad + \underbrace{\int \phi_g (g * K_{h_{\gamma}} - g) R_{\gamma}}_{\mathbf{B}_{\gamma}} + \underbrace{(S_{\gamma}^{-1} - 1) \int \phi_g T_{\gamma} R_{\gamma}}_{\mathbf{N}_{\gamma}^{(S)}}. \end{aligned}$$

Split  $A_{\gamma,1}$  and  $A_{\gamma,2}$  into a linear piece and a nonlinear remainder involving  $R_{\gamma} - 1$ :

$$\begin{aligned} \mathbf{A}_{\gamma,1} &= \underbrace{\int \phi_g (T_{\gamma} - \bar{f}_{\gamma,h})}_{\Xi_{\gamma}} + \underbrace{\int \phi_g (T_{\gamma} - \bar{f}_{\gamma,h}) (R_{\gamma} - 1)}_{\mathbf{R}_{\gamma,2}^{(1)}}, \\ \mathbf{A}_{\gamma,2} &= \underbrace{\int \phi_g (\bar{f}_{\gamma,h} - g * K_{h_{\gamma}})}_{\mathbf{U}_{\gamma}} + \underbrace{\int \phi_g (\bar{f}_{\gamma,h} - g * K_{h_{\gamma}}) (R_{\gamma} - 1)}_{\mathbf{R}_{\gamma,2}^{(2)}}. \end{aligned}$$

Set  $\mathbf{R}_{\gamma,2} := \mathbf{R}_{\gamma,2}^{(1)} + \mathbf{R}_{\gamma,2}^{(2)}$ .

**Identify the HT fluctuation.** By Fubini,

$$\Xi_{\gamma} = \int \phi_g (T_{\gamma} - \bar{f}_{\gamma,h}) = \frac{1}{N_{\gamma}} \sum_{i \in \mathcal{U}_{\gamma}} \left( \frac{\delta_{\gamma i}}{\pi_{\gamma i}} - 1 \right) \psi_{h_{\gamma}}(Y_{\gamma i}),$$

with  $\psi_{h_{\gamma}} = K_{h_{\gamma}} * \phi_g$ .

**Variance limit and CLT for  $\Xi_\gamma$ .** Define

$$\mathbf{V}_\gamma := n_{\text{V-eff},\gamma} \text{Var}(\Xi_\gamma \mid \{Y_{\gamma i}\}) = n_{\text{V-eff},\gamma} \frac{1}{N_\gamma^2} \sum_{i \in \mathcal{U}_\gamma} \frac{1 - \pi_{\gamma i}}{\pi_{\gamma i}} \psi_{h_\gamma}(Y_{\gamma i}) \psi_{h_\gamma}(Y_{\gamma i})^\top.$$

Let  $w_{\gamma i} := \frac{(1 - \pi_{\gamma i})/\pi_{\gamma i}}{\sum_{j \in \mathcal{U}_\gamma} (1 - \pi_{\gamma j})/\pi_{\gamma j}}$  so that  $\mathbf{V}_\gamma = \sum_{i \in \mathcal{U}_\gamma} w_{\gamma i} \psi_{h_\gamma}(Y_{\gamma i}) \psi_{h_\gamma}(Y_{\gamma i})^\top$  and  $\max_i w_{\gamma i} \rightarrow 0$  by the “no dominant unit” in Assumption **A11**. By Lemma B.2 (applied to each coordinate) and  $\mathbb{E}_g \|\psi_{h_\gamma}(Y)\| < \infty$  (from  $\phi_g \in L^2(g)$  and  $K_{h_\gamma} \in L^1$ ),

$$\mathbf{V}_\gamma \xrightarrow{\mathbb{P}} \mathbb{E}_g [\psi_{h_\gamma}(Y) \psi_{h_\gamma}(Y)^\top] =: \Sigma_{h_\gamma}.$$

By Lemma B.1,

$$\|\psi_{h_\gamma} - \phi_g\|_{L^2(g)} \rightarrow 0 \quad \Rightarrow \quad \Sigma_{h_\gamma} \rightarrow \Sigma := \mathbb{E}_g[\phi_g \phi_g^\top]$$

(since  $\|ab^\top - ac^\top\|_{\text{HS}} \leq \|a\|_2 \|b - c\|_2 + \|c\|_2 \|a - b\|_2$ ). Finally, the triangular-array Lindeberg condition in Assumption **A11** yields, conditionally on  $\{Y_{\gamma i}\}$ ,

$$\sqrt{n_{\text{V-eff},\gamma}} \Xi_\gamma \rightarrow \mathcal{N}_p(0, \Sigma),$$

and unconditioning preserves the limit.

**i.i.d. smoothing term  $U_\gamma$ .** Write

$$\begin{aligned} \mathbf{U}_\gamma &= \int \phi_g(y) \left[ \frac{1}{N_\gamma} \sum_{i=1}^{N_\gamma} \{K_{h_\gamma}(y - Y_{\gamma i}) - \mathbb{E}K_{h_\gamma}(y - Y)\} \right] dy \\ &= \frac{1}{N_\gamma} \sum_{i=1}^{N_\gamma} \left\{ \psi_{h_\gamma}^\#(Y_{\gamma i}) - \mathbb{E}\psi_{h_\gamma}^\#(Y) \right\}, \end{aligned}$$

where  $\psi_{h_\gamma}^\# := K_{h_\gamma}^\sim * \phi_g$  and  $K^\sim(t) := K(-t)$ . Hence  $\mathbb{E}[\mathbf{U}_\gamma] = 0$  and

$$\begin{aligned} \text{Var}(\mathbf{U}_\gamma) &= \frac{1}{N_\gamma} \mathbb{E}_g [\psi_{h_\gamma}^\#(Y)^2] = \frac{1}{N_\gamma} \|K_{h_\gamma} * \phi_g\|_{L^2(g)}^2 \\ &\leq \frac{\|\phi_g\|_{L^2(g)}^2}{N_\gamma} \int \frac{(K_{h_\gamma}^2 * g)(y)}{g(y)} dy, \end{aligned}$$

by Lemma B.3. Under the localized risk Assumption **A10**,  $\int (K_{h_\gamma}^2 * g)/g \leq C_R h_\gamma^{-d} + \tau_R(h_\gamma)$  for any fixed  $R$ , uniformly for small  $h_\gamma$ . Choosing  $R = R_\gamma \uparrow \infty$  with  $\tau_{R_\gamma}(h_\gamma) \rightarrow 0$  yields

$$\text{Var}(\mathbf{U}_\gamma) \lesssim \frac{1}{N_\gamma h_\gamma^d} \quad \Rightarrow \quad \sqrt{n_{\text{V-eff},\gamma}} \mathbf{U}_\gamma \xrightarrow{\mathbb{P}} 0$$

by Assumption **A9**.

**Bias term  $\mathbf{B}_\gamma$ .** By Cauchy-Schwarz and  $R_\gamma \in [0, 2]$ ,

$$|\mathbf{B}_\gamma| \leq 2\|\phi_g\|_{L^2(g)} \left\| \frac{g * K_{h_\gamma} - g}{\sqrt{g}} \right\|_{L^2}.$$

Under either kernel route in Assumption **A7** and the localized bias bound in Assumption **A10**,  $\|(g * K_{h_\gamma} - g)/\sqrt{g}\|_2 \lesssim h_\gamma^\beta$ , hence  $\sqrt{n_{V\text{-eff},\gamma}}\mathbf{B}_\gamma \rightarrow 0$  by Assumption **A9**.

**Nonlinear remainder  $R_{\gamma,2}$ .** Apply Lemma B.4 with  $u = T_\gamma - \bar{f}_{\gamma,h}$  and  $u = \bar{f}_{\gamma,h} - g * K_{h_\gamma}$ , and  $\varphi = \phi_g$ :

$$|\mathbf{R}_{\gamma,2}| \leq (MH(\hat{f}_{HT,\gamma}, g) + \sqrt{\epsilon(M)}) \left( \left\| \frac{T_\gamma - \bar{f}_{\gamma,h}}{\sqrt{g}} \right\|_2 + \left\| \frac{\bar{f}_{\gamma,h} - g * K_{h_\gamma}}{\sqrt{g}} \right\|_2 \right).$$

Under the localized risk bounds,

$$\begin{aligned} \mathbb{E} \left\| \frac{T_\gamma - \bar{f}_{\gamma,h}}{\sqrt{g}} \right\|_2^2 &\lesssim \frac{1}{n_{V\text{-eff},\gamma} h_\gamma^d} \quad (\text{Poisson-PPS; WOR has FPC } (1 - \alpha_\gamma)), \\ \mathbb{E} \left\| \frac{\bar{f}_{\gamma,h} - g * K_{h_\gamma}}{\sqrt{g}} \right\|_2^2 &\lesssim \frac{1}{N_\gamma h_\gamma^d}. \end{aligned}$$

Moreover  $H(\hat{f}_\gamma, g) = o_{\mathbb{P}}(n_{V\text{-eff},\gamma}^{-1/2})$  under Assumptions **A9** and **A6**. Choose  $M = M_\gamma \uparrow \infty$  so that  $\epsilon(M_\gamma) \rightarrow 0$ . Then  $\sqrt{n_{V\text{-eff},\gamma}}\mathbf{R}_{\gamma,2} \xrightarrow{\mathbb{P}} 0$ .

**Self-normalizer  $\mathbf{N}_\gamma^{(S)}$ .** Expand  $S_\gamma^{-1} = 1 - (S_\gamma - 1) + \rho_\gamma$ , where  $|\rho_\gamma| \leq 2(S_\gamma - 1)^2$  on  $\{|S_\gamma - 1| \leq \frac{1}{2}\}$ , which holds with high probability by the Bernstein/Freedman bound in Assumption **A6**. Then,

$$\mathbf{N}_\gamma^{(S)} = -(S_\gamma - 1)\mathbf{M}_\gamma + \rho_\gamma\mathbf{M}_\gamma, \quad \mathbf{M}_\gamma := \frac{1}{N_\gamma} \sum_{i \in \mathcal{U}_\gamma} \frac{\delta_{\gamma i}}{\pi_{\gamma i}} \tilde{\psi}_{h_\gamma,\gamma}(Y_{\gamma i}),$$

with  $\tilde{\psi}_{h_\gamma,\gamma}(y) := \int \phi_g(x) K_{h_\gamma}(x - y) R_\gamma(x) dx$ . Decompose

$$\mathbf{M}_\gamma = \underbrace{\frac{1}{N_\gamma} \sum_{i \in \mathcal{U}_\gamma} \tilde{\psi}_{h_\gamma,\gamma}(Y_{\gamma i})}_{\mathbf{m}_\gamma} + \underbrace{\frac{1}{N_\gamma} \sum_{i \in \mathcal{U}_\gamma} (\delta_{\gamma i}/\pi_{\gamma i} - 1) \tilde{\psi}_{h_\gamma,\gamma}(Y_{\gamma i})}_{\text{HT fluctuation}}.$$

By Lemma B.1 and  $R_\gamma \rightarrow 1$  in  $L^1(g)$ ,  $\mathbf{m}_\gamma - \mathbb{E}_g[\psi_{h_\gamma}(Y)] \rightarrow 0$  in probability, and  $\mathbb{E}_g[\psi_{h_\gamma}(Y)] = \int \phi_g g = 0$  (because  $\nabla \Gamma_G(\theta_0) = 0$ ). Hence  $\mathbf{m}_\gamma = o_{\mathbb{P}}(1)$  and the HT fluctuation is  $O_{\mathbb{P}}((n_{V\text{-eff},\gamma} h_\gamma^d)^{-1/2})$ . Using  $\sqrt{n_{V\text{-eff},\gamma}}(S_\gamma - 1) = O_{\mathbb{P}}(1)$  and  $\sqrt{n_{V\text{-eff},\gamma}}\rho_\gamma = O_{\mathbb{P}}(n_{V\text{-eff},\gamma}^{-1/2})$ , we get  $\sqrt{n_{V\text{-eff},\gamma}}\mathbf{N}_\gamma^{(S)} \rightarrow 0$  in probability.

**CLT conclusion via Taylor expansion.** Collecting all the previous steps,

$$\sqrt{n_{\text{V-eff},\gamma}} \nabla_{\theta} \Gamma_{\gamma}(\theta_0) = \sqrt{n_{\text{V-eff},\gamma}} \mathbf{\Xi}_{\gamma} + o_{\mathbb{P}}(1) \longrightarrow \mathcal{N}_p(0, \mathbf{\Sigma}).$$

(For fixed-size sampling, multiply by  $(1 - \alpha)$ .) A second-order Taylor expansion around  $\theta_0$  gives  $0 = \nabla_{\theta} \Gamma_{\gamma}(\hat{\theta}_{\gamma}) = \nabla_{\theta} \Gamma_{\gamma}(\theta_0) - \mathbf{A}_{\gamma}(\hat{\theta}_{\gamma} - \theta_0) + r_{\gamma}$ , with  $\mathbf{A}_{\gamma} := -\nabla_{\theta}^2 \Gamma_{\gamma}(\tilde{\theta}_{\gamma}) \xrightarrow{\mathbb{P}} \mathbf{A}$  (dominated differentiation, uniform LLN under the localized risk) and  $r_{\gamma} = o_{\mathbb{P}}(\|\hat{\theta}_{\gamma} - \theta_0\|)$ . Thus

$$\sqrt{n_{\text{V-eff},\gamma}}(\hat{\theta}_{\gamma} - \theta_0) = \mathbf{A}_{\gamma}^{-1} \sqrt{n_{\text{V-eff},\gamma}} \nabla_{\theta} \Gamma_{\gamma}(\theta_0) + o_{\mathbb{P}}(1) \longrightarrow \mathcal{N}_p(0, \mathbf{A}^{-1} \mathbf{\Sigma} \mathbf{A}^{-\top})$$

for Poisson-PPS, and with covariance  $\mathbf{A}^{-1}[(1 - \alpha)\mathbf{\Sigma}]\mathbf{A}^{-\top}$  for fixed-size SRS-WOR.

□

## B.1 Technical Lemmas

**Lemma B.1** (Approximate identity in  $L^2(f)$  via localization). *Under assumptions **A1** and **A9**, let  $A_R \uparrow \mathbb{R}^d$  be compacts with  $0 < c_R \leq g \leq C_R < \infty$  on  $A_R$ . If  $\phi_g \in L^2(g)$  then  $\|K_{h_{\gamma}} * \phi_g - \phi_g\|_{L^2(g)} \rightarrow 0$  as  $\gamma \rightarrow \infty$ .*

*Proof.* On  $A_R$ ,  $g$  is bounded above and below, hence  $\int_{A_R} |K_{h_{\gamma}} * \phi_g - \phi_g|^2 g \leq C_R \|K_{h_{\gamma}} * \phi_g - \phi_g\|_{L^2}^2 \rightarrow 0$  by the  $L^2$  approximate identity (Young + density). On  $A_R^c$ ,  $\int_{A_R^c} |K_{h_{\gamma}} * \phi_g - \phi_g|^2 g \leq 2 \int_{A_R^c} |K_{h_{\gamma}} * \phi_g|^2 g + 2 \int_{A_R^c} |\phi_g|^2 g$ . The second term goes to 0 from above as  $R \uparrow \infty$  since  $\phi_g \in L^2(g)$ . For the first term, fix  $R$  and split  $\phi_g = \phi_g \mathbf{1}_{A_R} + \phi_g \mathbf{1}_{A_R^c}$ , use  $K_{h_{\gamma}} \in L^1$ , Young, and the previous tail smallness of  $\phi_g \mathbf{1}_{A_R^c}$  to make it  $< \varepsilon$  uniformly for small  $h_{\gamma}$ . Now take  $\gamma \rightarrow \infty$ , then  $R \rightarrow \infty$ .

□

**Lemma B.2** (Weighted LLN for triangular weights). *Let  $Z_{\gamma i} \in \mathbb{R}^q$  be i.i.d. with  $\mathbb{E}\|Z_{\gamma 1}\| < \infty$ . Let  $w_{\gamma i} \geq 0$  with  $\sum_{i \in \mathcal{W}_{\gamma}} w_{\gamma i} = 1$  and  $\max_{i \in \mathcal{W}_{\gamma}} w_{\gamma i} \rightarrow 0$ . Then  $\sum_{i \in \mathcal{W}_{\gamma}} w_{\gamma i} Z_{\gamma i} \xrightarrow{\mathbb{P}} \mathbb{E} Z_{\gamma 1}$ .*

*Proof.* Write the scalar case and apply Chebyshev with  $\text{Var}(\sum_{i \in \mathcal{W}_{\gamma}} w_i Z_i) \leq (\max_i w_i) \sum w_i \mathbb{E}\|Z_i - \mathbb{E} Z\|^2$  when  $\mathbb{E}\|Z\|^2 < \infty$  (obtainable by truncation if only the first moment exists). Extend to vectors by component-wise application.

□

**Lemma B.3** (Weighted convolution inequality). *For any  $\varphi \in L^2(g)$ ,*

$$\|K_{h_{\gamma}} * \varphi\|_{L^2(g)}^2 \leq \|\varphi\|_{L^2(g)}^2 \int_{\mathbb{R}^d} \frac{(K_{h_{\gamma}}^2 * g)(y)}{g(y)} dy.$$

*Proof.* By Cauchy–Schwarz with weight  $g(x)$  inside the convolution:

$$\begin{aligned} \left| (K_{h_\gamma} * \varphi)(y) \right| &= \left| \int \varphi(x) K_{h_\gamma}(x - y) dx \right| = \left| \int \varphi(x) \sqrt{g(x)} \frac{K_{h_\gamma}(x - y)}{\sqrt{g(x)}} dx \right| \\ &\leq \|\varphi\|_{L^2(g)} \left( \int \frac{K_{h_\gamma}^2(x - y)}{g(x)} dx \right)^{1/2}. \end{aligned}$$

Taking the square and multiplying by  $g(y)$ , then integrating over  $y$  to obtain the claim after Fubini.  $\square$

**Lemma B.4** (Remainder via Hellinger and  $\chi^2(g)$ ). *For any  $\varphi \in L^2(g)$ , any square-integrable  $u$ , and densities  $h, g$ ,*

$$\left| \int \varphi u (R_h - 1) \right| \leq MH(h, g) \left\| \frac{u}{\sqrt{g}} \right\|_2 + \sqrt{\epsilon(M)} \left\| \frac{u}{\sqrt{g}} \right\|_2,$$

where  $R_h = 2\sqrt{g}/(\sqrt{h} + \sqrt{g})$ ,  $H(h, g) = \|\sqrt{h} - \sqrt{g}\|_2$ , and  $\epsilon(M) := \int_{\{|\varphi| > M\}} \varphi^2 g dy$ .

*Proof.* Split the domain into  $\{|\varphi| \leq M\}$  and its complement. On  $\{|\varphi| \leq M\}$ ,

$$\int |\varphi u| |R_h - 1| \leq M \int \frac{|u|}{\sqrt{g}} \frac{|\sqrt{h} - \sqrt{g}|}{\sqrt{h} + \sqrt{g}} \sqrt{g} \leq M \left\| \frac{u}{\sqrt{g}} \right\|_2 \left( \int g |R_h - 1|^2 \right)^{1/2}.$$

Since  $|R_h - 1| = |\sqrt{h} - \sqrt{g}|/(\sqrt{h} + \sqrt{g})$  and  $\sqrt{h} + \sqrt{g} \geq \sqrt{g}$ ,  $\int g |R_h - 1|^2 \leq \int (\sqrt{h} - \sqrt{g})^2 = H(h, g)^2$ . On  $\{|\varphi| > M\}$ , Cauchy–Schwarz gives  $\int |\varphi u| |R_h - 1| \leq \sqrt{\epsilon(M)} \|u/\sqrt{g}\|_2$  since  $|R_h - 1| \leq 1$ . Combine the two bounds.  $\square$

## Appendix C Robustness proof

**Lemma C.1** (Uniform Hellinger inequality). *For any distributions  $F_1, F_2$  with densities  $f_1, f_2$  and any  $\theta$ ,*

$$\left| \Gamma_{F_1}(\theta) - \Gamma_{F_2}(\theta) \right| = \left| \int \sqrt{f_\theta(y)} \left( \sqrt{f_1(y)} - \sqrt{f_2(y)} \right) dy \right| \leq H(f_1, f_2).$$

*Proof.* Cauchy–Schwarz and  $\|\sqrt{f_\theta}\|_2 = 1$ .  $\square$

**Proposition C.2** (Continuity of  $T$  in the Hellinger topology). *Let  $(G_n)$  be a sequence of distributions with densities  $g_n$  s.t.  $H(g_n, g) \rightarrow 0$ . Under Assumption A12, any sequence of maximizers  $T(G_n)$  satisfies  $T(G_n) \rightarrow \theta_0 = T(G)$ .*

*Proof.* By Lemma C.1,  $\sup_\theta |\Gamma_{g_n}(\theta) - \Gamma_g(\theta)| \leq H(g_n, g) \rightarrow 0$ . Fix  $\varepsilon > 0$ . By separation,

$\sup_{\|\theta - \theta_0\| \geq \varepsilon} \Gamma_g(\theta) \leq \Gamma_g(\theta_0) - \Delta(\varepsilon)$ . For  $n$  large with  $H(g_n, g) \leq \Delta(\varepsilon)/3$ ,

$$\sup_{\|\theta - \theta_0\| \geq \varepsilon} \Gamma_{g_n}(\theta) \leq \Gamma_g(\theta_0) - \frac{2}{3}\Delta(\varepsilon), \quad \Gamma_{g_n}(\theta_0) \geq \Gamma_g(\theta_0) - \frac{1}{3}\Delta(\varepsilon).$$

Hence any maximizer  $T(G_n)$  must lie in  $B(\theta_0, \varepsilon)$ . Since  $\varepsilon$  is arbitrary,  $T(G_n) \rightarrow \theta_0$ .  $\square$

*Proof of Theorem 3.16.* Let  $G_\varepsilon = (1 - \varepsilon)G + \varepsilon H$  and  $\theta_\varepsilon := T(G_\varepsilon)$ . By definition,  $S_{G_\varepsilon}(\theta_\varepsilon) = 0$  for small  $\varepsilon$ . Taylor expanding  $S_{G_\varepsilon}(\theta)$  at  $(\theta_0, G)$ ,

$$0 = S_G(\theta_0) + \mathbf{Q}(\theta_\varepsilon - \theta_0) + \dot{S}_G(\theta_0; H)\varepsilon + R_\varepsilon,$$

where  $\|R_\varepsilon\| = o(\|\theta_\varepsilon - \theta_0\|) + o(\varepsilon)$  by the dominated smoothness in Assumption **A12**(ii) and the definition of  $\dot{S}_G$ . Since  $S_G(\theta_0) = 0$  and  $\mathbf{Q}$  is nonsingular,

$$\theta_\varepsilon - \theta_0 = -\mathbf{Q}^{-1}\dot{S}_G(\theta_0; H)\varepsilon + o(\varepsilon).$$

Dividing by  $\varepsilon$  and letting  $\varepsilon \downarrow 0$  gives the stated derivative.  $\square$

## Appendix D Additional simulation results

Here we present additional results for the simulations in Section 4.1 of the main manuscript. Unless otherwise noted, the simulation settings are identical to the main manuscript.

### D.1 Gamma Model with Calibrated Weights

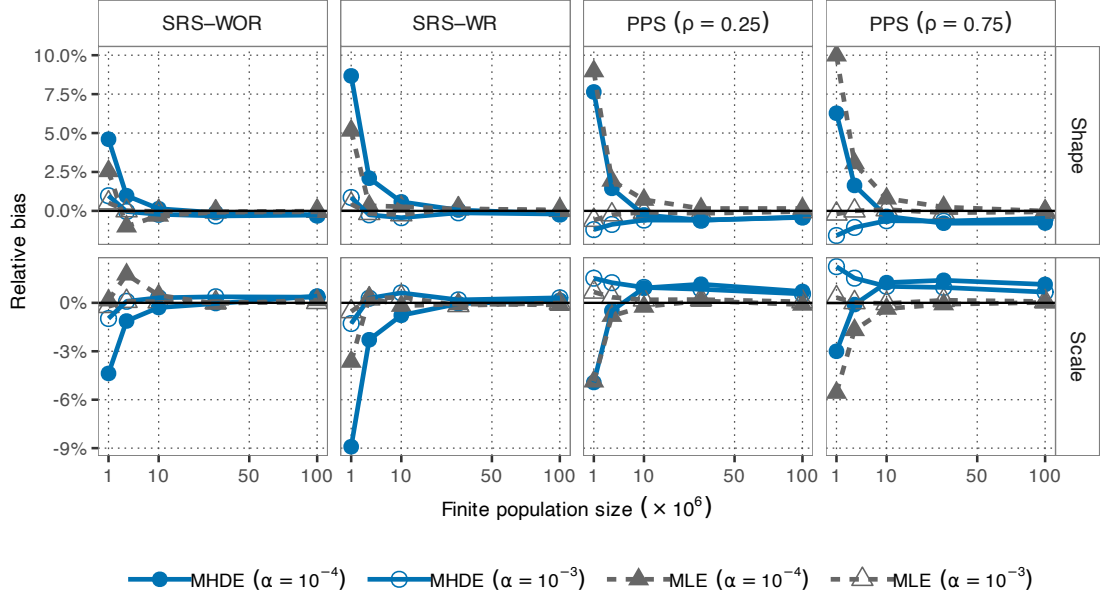
We now consider that each unit in the finite population  $\mathcal{U}$  belongs to one of five clusters. We then consider calibrated sampling weights based on an auxiliary variable  $X$ , with known cluster totals.

Specifically, for each  $i \in \mathcal{U}$  we know the cluster assignment via the membership function  $C: i \rightarrow \{1, \dots, 5\}$ . Moreover, we assume to know the cluster totals for the population,  $\bar{X}_c = \sum_{i \in \mathcal{U}} X_i \mathbf{1}_{\{C(i)=c\}}$ ,  $c \in \{1, \dots, 5\}$ . Given a sample  $\mathcal{S}$  and survey weights  $w_i$ , we determine the calibration adjustment factors  $\xi_c = \frac{1}{\bar{X}_c} \sum_{i \in \mathcal{S}} w_i X_i \mathbf{1}_{\{C(i)=c\}}$ . The calibrated weights are then given by  $\omega_i = w_i \xi_{C(i)}$  for  $i \in \mathcal{S}$ .

The results for calibrated weights are very similar to the results with the original survey weights. The relative bias in Figure 5 and the relative variance in Figure 6 still show that the MHDE is very close to the MLE across all scenarios.

#### D.1.1 Alternative Contamination Model

Here we consider a truncated  $t$  distribution as the source of contamination instead of the Normal distribution from the main manuscript. We replace  $\lfloor \varepsilon n \rfloor$  observations in the sample with



**Figure 5:** Relative bias of the MHDE (blue dots) and MLE (gray triangles) in a Gamma superpopulation model using calibrated weights and various sampling designs. The sample size in each simulation is determined by  $n = \alpha N$ , with  $\alpha \in \{10^{-3}, 10^{-4}\}$ .

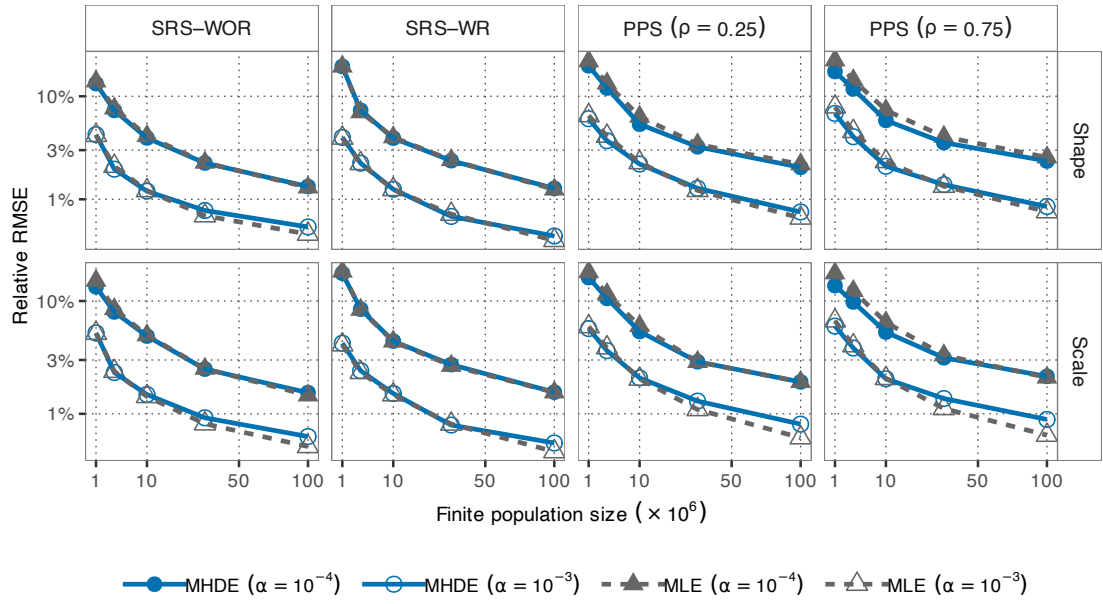
i.i.d. draws from a shifted and truncated (positive)  $t$  distribution with 3 degrees of freedom with mode at  $z > 0$ . We further scale the contamination to have the same variance as the true Gamma distribution from the superpopulation.

Figure 7 shows the influence function (top) for 10% contamination and varying

## D.2 Lognormal Model

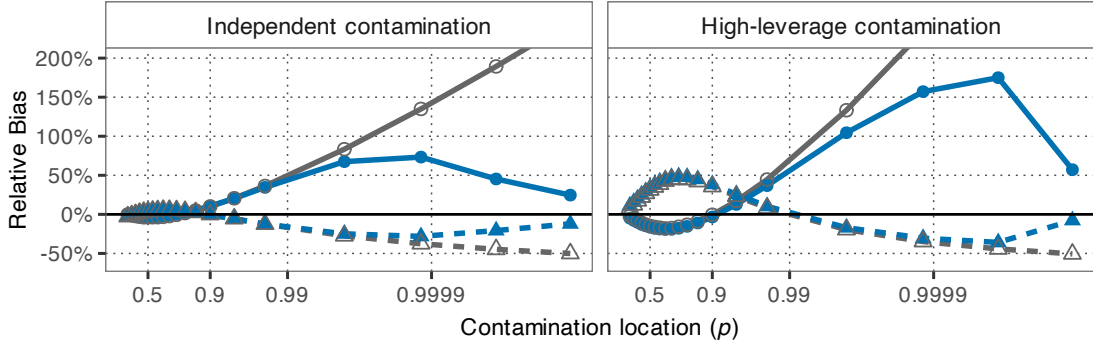
Instead of the Gamma model, we now consider a lognormal superpopulation model,  $Y \sim \text{LN}(\mu = 9, \sigma = 2)$ . The bias, shown in Figure 8, is close to 0 for both the mean and the SD parameters. Similarly, the variance goes to 0 quickly as the finite population size and the sample size increase, with practically no difference between the MHDE and the MLE. Due to the minimal bias, inference using the asymptotic distribution of the MHDE is also highly reliable. The empirical coverage probability of the 95% CI in Figure 10 is very close to the nominal level across all sampling schemes, even for small sample sizes.



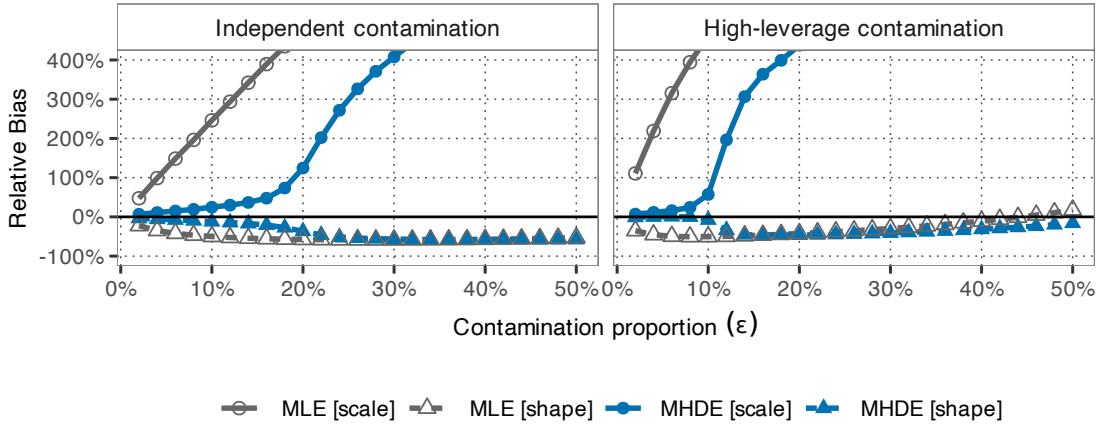


**Figure 6:** Relative RMSE of the MHDE (blue dots) and MLE (gray triangles) under a Gamma superpopulation model using calibrated weights and various sampling designs. The sample size in each simulation is determined by  $n = \alpha N$ , with  $\alpha \in \{10^{-3}, 10^{-4}\}$ .

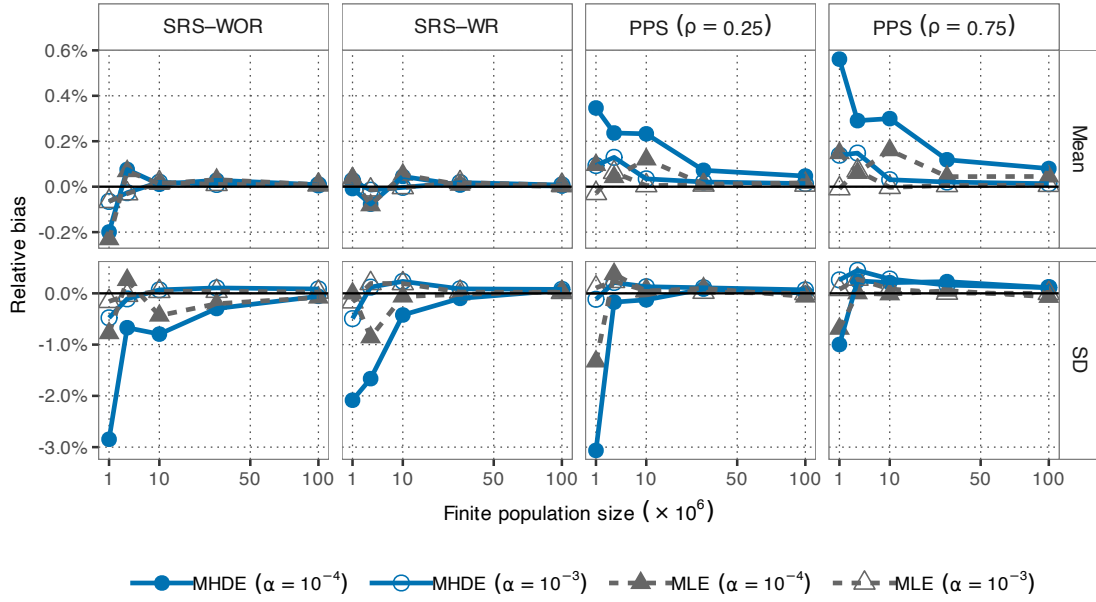
(a) Influence function



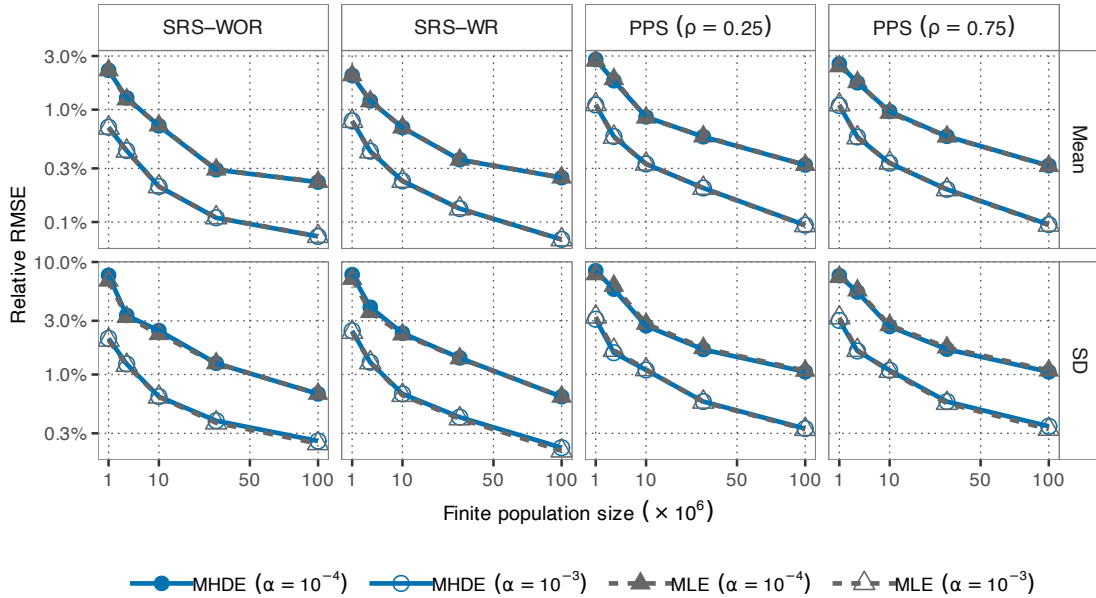
(b)  $\alpha$ -influence curve



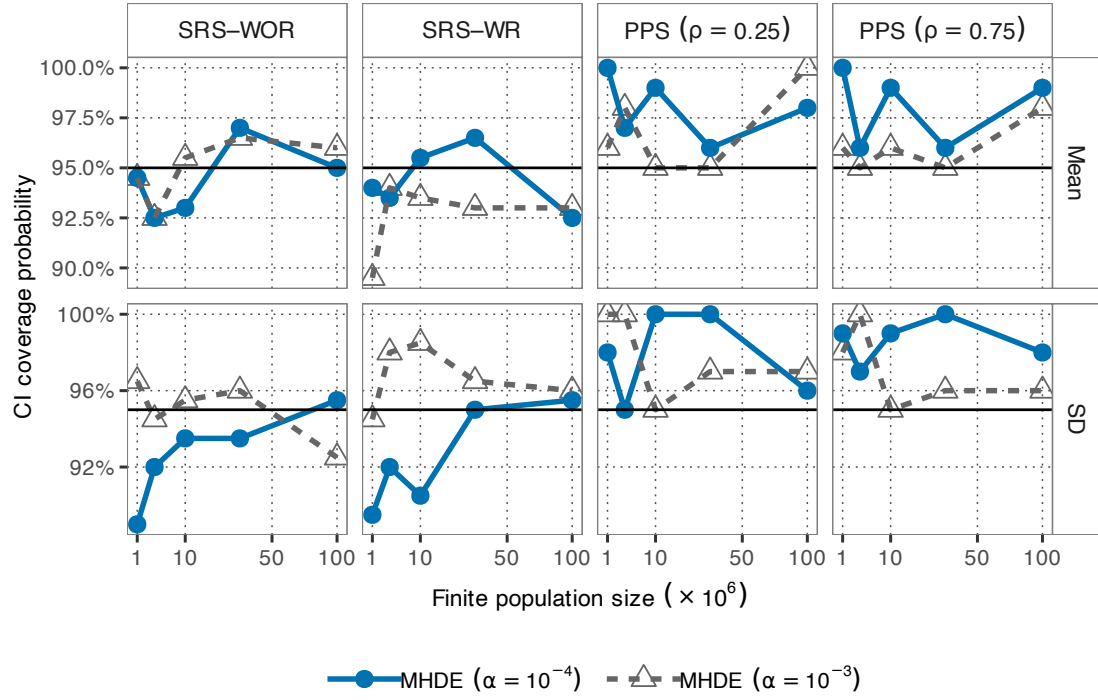
**Figure 7:** Influence functions (top) and alpha curves (bottom) for the MHDE and MLE of the scale ( $\circ$ ) and shape ( $\triangle$ ) parameters in the Gamma model. The contamination proportion in the influence function at the top is set to  $\varepsilon = 0.1$ . The horizontal axis shows the location of the mode of the truncated  $t$  distribution in terms of the quantile of the true superpopulation model, i.e.,  $z = G^{-1}(p)$ . For the alpha curve the mode of the  $t$  distribution is located at  $z = G^{-1}(p = 0.99) \approx 232.342$ .



**Figure 8:** Relative bias of the MHDE (blue dots) and MLE (gray triangles) in a Lognormal superpopulation model using various sampling designs. The sample size in each simulation is determined by  $n = \alpha N$ , with  $\alpha \in \{10^{-3}, 10^{-4}\}$ .



**Figure 9:** Relative RMSE of the MHDE (blue dots) and MLE (gray triangles) in a Lognormal superpopulation model using various sampling designs. The sample size in each simulation is determined by  $n = \alpha N$ , with  $\alpha \in \{10^{-3}, 10^{-4}\}$ .



**Figure 10:** Coverage probability of the 95% confidence interval for the MHDE in a Lognormal superpopulation model as the finite population size  $N$  increases using different sampling designs. The sample size in each simulation is determined by  $n = \alpha N$ , with  $\alpha \in \{10^{-3}, 10^{-4}\}$ .

## References

- [1] L.-C. Zhang, “Post-Stratification and Calibration—A Synthesis,” *The American Statistician*, vol. 54, no. 3, pp. 178–184, Aug. 2000.
- [2] R. Beran, “Minimum Hellinger Distance Estimates for Parametric Models,” *The Annals of Statistics*, vol. 5, no. 3, pp. 445–463, May 1977.
- [3] D. L. Donoho and R. C. Liu, “The ”Automatic” Robustness of Minimum Distance Functionals,” *The Annals of Statistics*, vol. 16, no. 2, Jun. 1988.
- [4] D. G. Simpson, “Hellinger Deviance Tests: Efficiency, Breakdown Points, and Examples,” *Journal of the American Statistical Association*, vol. 84, no. 405, pp. 107–113, Mar. 1989.
- [5] B. G. Lindsay, “Efficiency Versus Robustness: The Case for Minimum Hellinger Distance and Related Methods,” *The Annals of Statistics*, vol. 22, no. 2, Jun. 1994.
- [6] Z. Lu, Y. V. Hui, and A. H. Lee, “Minimum Hellinger Distance Estimation for Finite Mixtures of Poisson Regression Models and Its Applications,” *Biometrics*, vol. 59, no. 4, pp. 1016–1026, Dec. 2003.
- [7] E. Castilla, N. Martín, and L. Pardo, “Minimum phi-divergence estimators for multinomial logistic regression with complex sample design,” *AStA Advances in Statistical Analysis*, vol. 102, no. 3, pp. 381–411, Jul. 2018.
- [8] E. Castilla, A. Ghosh, N. Martin, and L. Pardo, “Robust semiparametric inference for polytomous logistic regression with complex survey design,” *Advances in Data Analysis and Classification*, vol. 15, no. 3, pp. 701–734, Sep. 2021.
- [9] C.-E. Särndal, B. Swensson, and J. Wretman, *Model Assisted Survey Sampling*, ser. Springer Series in Statistics. New York, NY: Springer New York, 1992.
- [10] CDC. Center for Disease Control and Prevention, “National Center for Health Statistics (NCHS). National Health and Nutrition Examination Survey Data,” 2025.
- [11] D. G. Horvitz and D. J. Thompson, “A Generalization of Sampling Without Replacement from a Finite Universe,” *Journal of the American Statistical Association*, vol. 47, no. 260, pp. 663–685, Dec. 1952.
- [12] J. W. Tukey, *A Survey of Sampling from Contaminated Distributions*. Princeton, New Jersey: Princeton University, 1959.
- [13] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 1st ed. Wiley, Oct. 2001.

- [14] J. A. Nelder and R. Mead, “A Simplex Method for Function Minimization,” *The Computer Journal*, vol. 7, no. 4, pp. 308–313, Jan. 1965.
- [15] A.-l. Cheng and A. N. Vidyashankar, “Minimum Hellinger distance estimation for randomized play the winner design,” *Journal of Statistical Planning and Inference*, vol. 136, no. 6, pp. 1875–1910, Jun. 2006.
- [16] F. R. Hampel, “The Influence Curve and its Role in Robust Estimation,” *Journal of the American Statistical Association*, vol. 69, no. 346, pp. 383–393, 1974.
- [17] H. J. Adrogué and N. E. Madias, “Hyponatremia,” *New England Journal of Medicine*, vol. 342, no. 21, pp. 1581–1589, May 2000.
- [18] L. Kish, “Weighting for unequal P<sub>i</sub>,” *Journal of Official Statistics*, vol. 8, no. 2, p. 183, 1992.
- [19] T. Lumley, P. Gao, and B. Schneider, “Survey: Analysis of Complex Survey Samples,” pp. 4.4–8, Jan. 2003.
- [20] N. Cressie and T. R. Read, “Multinomial Goodness-Of-Fit Tests,” *Journal of the Royal Statistical Society Series B: Statistical Methodology*, vol. 46, no. 3, pp. 440–464, Jul. 1984.
- [21] L. Pardo, *Statistical Inference Based on Divergence Measures*, ser. Statistics: Textbooks and Monographs. Boca Raton, Fla.: Chapman & Hall/CRC, 2006, no. 185.
- [22] D. A. Freedman, “On Tail Probabilities for Martingales,” *The Annals of Probability*, vol. 3, no. 1, Feb. 1975.