

Article

A Ratio Estimator for the Mean Using a Mixture Optional Enhance Trust (MOET) Randomized Response Model

Sat Gupta ¹ , Michael Parker ^{1,*} and Sadia Khalil ²¹ Department of Mathematics and Statistics, UNC Greensboro, Greensboro, NC 27413, USA; sngupta@uncg.edu² Department of Statistics, Lahore College for Women University, Lahore 54000, Pakistan; sadia.khalil@lcwu.edu.pk

* Correspondence: mlparke4@uncg.edu

Abstract: When researchers conduct surveys seeking sensitive, socially stigmatized information, respondents, on average, modify their answers to represent themselves favorably. To overcome this issue, researchers may use Randomized Response Technique (RRT) models. Recently, Parker et al. proposed a model that incorporates some of the most critical recent quantitative RRT advancements—mixture, optionality, and enhanced trust—into a single model, which they called a Mixture Optional Enhanced (MOET) model. We now improve upon the MOET model by incorporating auxiliary information into the analysis. Positively correlated auxiliary information can improve the mean response estimation through use of a ratio estimator. In this study, we propose just such an estimator for the MOET model. Further, we investigate the conditions under which the ratio estimator outperforms the basic MOET estimator proposed by Parker et al. in 2024. We also consider the possibility that the collection of auxiliary information may compromise privacy; and we study the impact of privacy reduction on the overall model performance as assessed by the unified measure (UM) proposed by Gupta et al. in 2018.

Keywords: Randomized Response Technique (RRT); respondent privacy; social desirability bias (SDB); unified measure (UM); ratio estimator; auxiliary variable

MSC: 62D05



Citation: Gupta, S.; Parker, M.; Khalil, S. A Ratio Estimator for the Mean Using a Mixture Optional Enhance Trust (MOET) Randomized Response Model. *Mathematics* **2024**, *12*, 3617. <https://doi.org/10.3390/math12223617>

Academic Editors: Davide Valenti and Manuel Alberto M. Ferreira

Received: 22 September 2024
Revised: 5 November 2024
Accepted: 15 November 2024
Published: 20 November 2024



Copyright: © 2024 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

Researchers in disciplines as far-flung as business, public health, and psychology need answers to questions that are hard to obtain, because they involve sensitive personal information. In such scenarios, “social desirability bias” comes into play. Latkin et al. (2017) describe SDB as the tendency of respondents to underreport socially undesirable attitudes or behaviors and to overreport more desirable attributes [1].

For example, income level is often perceived to be an indicator of success and status. Consequently, respondents may falsify their responses if queried directly about their income. For this reason, a businessperson doing market research, who needs accurate information about the income level of people in a target market, may be stymied. Similarly, a public health professional who makes decisions about how to allocate public funds may need reliable information about the prevalence of sexually transmitted diseases within a certain population but may have difficulty getting individuals to self-report their STD history honestly.

Warner (1965) and Greenberg (1969) first proposed a class of models, known as Randomized Response Technique (RRT) models, that were designed to encourage reliable responses to yes/no-type survey questions [2,3]. Later, Warner (1971) and Greenberg (1971) proposed new models designed to study questions which responses were numerical in nature; these became known as Quantitative RRT models [4,5]. Warner’s model relied on the “scrambling” of responses based on additive and multiplicative scrambling variables,

which would keep the respondents' true answers to a sensitive study question hidden, freeing them to answer truthfully. Greenberg's model was entirely different. Rather than perturbing the response, Greenberg's methodology involved hiding the question. Specifically, the question posed to each respondent would be one of two questions—the sensitive question under study or some unrelated question which response would appear similar to that of the sensitive question. Since the identity of the question would be unknown to the researcher, the respondent's anonymity would be maintained, and the respondent would be free to answer the question honestly.

Following Warner's and Greenberg's models, RRT was studied extensively, and many innovations and advancements were made. Several statisticians studied the scrambling that underlay Warner's quantitative model. Pollock and Beck (1976) studied the attributes of additive versus multiplicative scrambling [6]. Eichhorn and Hayre (1983) introduced a scrambling paradigm involving multiplicative scrambling, which Diana and Perri (2011) explored later in greater detail [7,8]. Singh et al. (2018) compared several of these models with an added emphasis on considering not only the efficiency of estimation but also the importance of enhanced privacy protection [9].

Gupta et al. (2002) improved the efficiency of RRT models by introducing the concept of "optionality", where survey respondents were instructed to answer the sensitive question truthfully (with no RRT) if they did not find the question to be sensitive to them [10]. This advancement also led to a means of assessing the sensitivity level of sensitive questions. Mehta and Aggarwal (2018) proposed a Bayesian approach to measuring the sensitivity of binary questions [11]. Sharma and Singh (2015) recognized the importance of nonresponse and measurement error in RRT scenarios (RRT questions, due to their sensitive nature, are likely to provoke both nonresponse and measurement error) [12]. They proposed means of accounting for these issues, thereby avoiding unrealistically favorable statistical inferences. Statisticians also considered the advantages that may be realized by "mixing" models. Perri (2008) first proposed a binary blank card RRT model that had interesting theoretical value in that it incorporated Warner-like features (direct and indirect questions) with Greenberg-like features (unrelated questions) into a single model [13]. Perri's model can therefore be thought of as an early exploration of "mixture" models (an important underpinning of the MOET model considered in this study). Following the considerable theoretical advancements over 50 years, Blair et al. (2015), among others, focused on practical design and technique issues that would facilitate the implementation of advanced models [14].

In Section 2 of this study, a brief review the Mixture Optional Enhanced Trust model (MOET) model proposed by Parker et al. (2024) is provided [15]. In Section 3, we propose a mean response ratio estimator for the MOET model, along with estimators that evaluate the MOET ratio estimator's bias and efficiency and the MOET model's privacy. In Section 4, we consider two characteristics—variance and correlation—that together assess the quality of auxiliary information, under the standard (but not always valid) assumption that auxiliary information does not impact model privacy. Under this assumption, we derive a relationship that assesses whether auxiliary information will improve the estimation of the mean response to the sensitive question. Then, in Section 5, we consider the possibility that auxiliary information does, in fact, compromise privacy, and we use the unified measure (UM) proposed by Gupta et al. (2018) to study the reduction in the ratio estimator's overall quality as privacy declines [16]. In Section 6, we provide tabular demonstrations and empirical model simulations that confirm the accuracy of the estimators developed in Section 3 of this study, and we consider the behavior of the ratio estimator across different values of key parameters.

2. MOET Model (2024) Review

As this paper will study the role of auxiliary information in improving the estimation of mean responses to sensitive questions when using the MOET model, we will begin with

a brief review of the MOET model itself. Figure 1 presents a diagram of the model, which appears in the study by Parker et al. (2024) [15].

$$Z = \begin{cases} Y + S & \text{with probability } W\alpha A \\ TY + S & \text{with probability } W(1 - A)(\alpha + p - \alpha p) \\ Y & \text{with probability } W(1 - \alpha)pA + (1 - W) \\ R & \text{with probability } W(1 - \alpha)(1 - p) \end{cases} \quad (1)$$

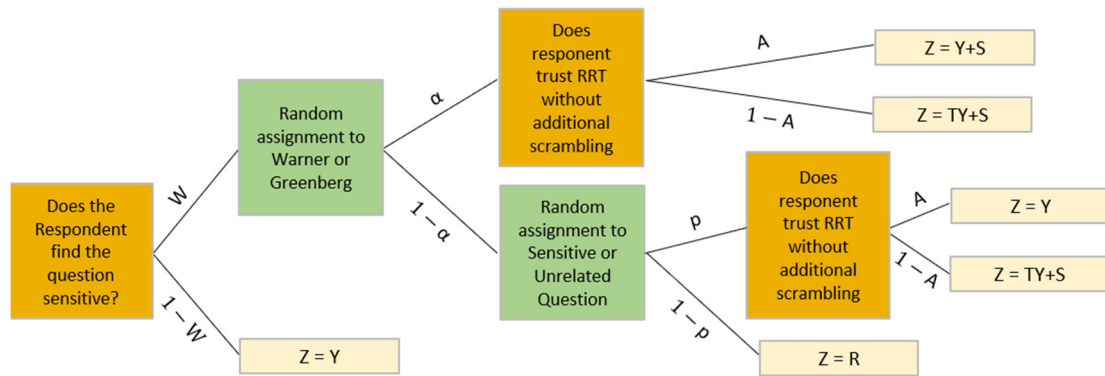


Figure 1. MOET model diagram.

Parker et al. (2024) further proposed a mean estimator and derived the expressions below, according to the MOET model, using a split sample approach [15].

Mean Estimator (Parker et al. 2024)

$$\hat{\mu}_Y = \frac{1 - p_1}{p_2 - p_1} \bar{Z}_2 - \frac{1 - p_2}{p_2 - p_1} \bar{Z}_1. \quad (2)$$

MSE of Mean Estimator (Parker et al. 2024)

$$\text{MSE}(\hat{\mu}_Y) = \left(\frac{1 - p_1}{p_2 - p_1} \right)^2 \text{Var}(\bar{Z}_2) + \left(\frac{1 - p_2}{p_2 - p_1} \right)^2 \text{Var}(\bar{Z}_1), \quad (3)$$

where

$$\text{var}(\bar{Z}_i) = \frac{2}{n} \left\{ W[1 - \lambda_i - A\phi_i] \sigma_S^2 + [W(1 - \lambda_i)((1 - A)\sigma_T^2 + 1) + 1 - W](\sigma_Y^2 + \mu_Y^2) + W\lambda_i(\sigma_R^2 + \mu_R^2) - [\mu_Y + (\mu_R - \mu_Y)W\lambda_i]^2 \right\},$$

$$\lambda_i = (1 - \alpha)(1 - p_i), \quad i = 1, 2,$$

$$\phi_i = p_i(1 - \alpha), \quad i = 1, 2,$$

$$p_1 \neq p_2.$$

In the expressions above, the following symbols are used.

- Y : The true response to the sensitive question. This random variable has mean μ_Y and variance σ_Y^2 .
- Z_i : The response collected from the respondent in the i^{th} sub-sample, $i = 1, 2$. This random variable has mean μ_{Z_i} and variance $\sigma_{Z_i}^2$.
- S : An additive scrambling variable with mean $\mu_S = 0$ and variance σ_S^2 .
- T : A multiplicative scrambling variable with mean $\mu_T = 1$ and variance σ_T^2 . T is independent of S and Y .
- R : Response to unrelated question with mean μ_R and variance σ_R^2 .
- n : The sample size. In split sampling, n is split into n_1 and n_2 , where $n_1 + n_2 = n$.
- p_i : The probability that an individual that has been assigned to the Greenberg sub-model within sub-sample i is assigned the sensitive question, as opposed to the unrelated question.

- A : The probability that a respondent will trust the RRT methodology without additional scrambling.
- W : The sensitivity level of the sensitive question. That is, a proportion $(1 - W)$ of the respondents do not consider the question sensitive and are willing to provide true responses without scrambling.

The Privacy Intrinsic to the MOET Model is given by:

$$V^a = \frac{1}{2} \left\{ 2\alpha A \sigma_S^2 + (2 - \lambda_1 - \lambda_2)(1 - A) \left[(\sigma_Y^2 + \mu_Y^2) \sigma_T^2 + \sigma_S^2 \right] + (\lambda_1 + \lambda_2) \left[\sigma_Y^2 + \sigma_R^2 (\mu_Y - \mu_R)^2 \right] \right\}. \tag{4}$$

The superscript a in the above expression reminds us that this measure has been adjusted to reflect the Gupta et al. (2018) assertion that optionality does not undermine privacy for the proportion of respondents $(1 - W)$ who do not consider the question sensitive [16].

The Unified Measure Intrinsic to the MOET Model is given by:

$$\delta^a = \frac{\text{MSE}(\hat{\mu}_Y)}{V^a}, \tag{5}$$

where $\text{MSE}(\hat{\mu}_Y)$ and V^a are defined by Equations (3) and (4). The unified measure (UM) proposed by Gupta et al. (2018) enables a quantification of the overall model quality, taking into account the competing factors of estimator efficiency and model privacy [16]. This measure highlights one of the key features of the MOET model—its mixture capability. The MOET model is fundamentally a mix between a Greenberg-type model and a Warner-type model, where the Greenberg model is generally more efficient and less private, but the Warner model is more private and less efficient. The model is, as detailed by Parker et al. (2024), sufficiently flexible in that it allows the researcher to choose the optimal balance between the two factors by strategically setting the “mixture parameter”, α , according to the researcher’s specific needs [15]. The researcher can also adjust the balance between privacy and efficiency, selecting a higher or lower scrambling variance. More scrambling (higher σ_S^2 and higher σ_T^2) will result in a model with more privacy but at a cost to efficiency.

The Estimator for the Sensitivity Parameter (W) is given by:

$$\hat{W} = \frac{\bar{Z}_1 - \bar{Z}_2}{\lambda_1(\mu_R - \bar{Z}_2) - \lambda_2(\mu_R - \bar{Z}_1)}, \tag{6}$$

$$\lambda_i = (1 - \alpha)(1 - p_i), \mu_Y \neq \mu_R, p_1 \neq p_2.$$

Parker et al. (2024) showed that the quantitative MOET model had important advantages over the OET model, which Gupta et al. (2022) showed was superior to the basic Warner model, especially with regard to estimator efficiency [15,17]. The mixture feature of the MOET model enabled blending, which was not a part of the OET model. Importantly, this feature made MOET superior to both a fully Warner-based or fully Greenberg-based model.

3. Ratio Estimator

Thompson (2012) [18] showed how auxiliary information that is strongly and positively correlated with the response to a sensitive question can be used as the basis of a ratio estimator. Some concrete examples of valuable auxiliary information might include “home value”, which may assist in the estimation of “personal wealth”, and “frequency of fast-food credit card purchases”, which may assist in the estimation of “body mass index”.

If auxiliary information is of adequate quality, the estimates achieved through use of a ratio estimator will be superior to those yielded by a “basic” estimator (in this study, we use the term “basic” to describe an estimator that does not incorporate auxiliary information). In this section, we propose just such a ratio estimator, and we derive expressions that

capture the ratio estimator’s MSE and Bias, the MOET model’s privacy, and the sensitivity level of the sensitive question.

3.1. Ratio Estimator Development

Ratio estimators are of the form

$$\hat{\mu}_{YR} = \hat{\mu}_Y \left(\frac{\mu_X}{\hat{\mu}_X} \right), \tag{7}$$

where X is defined as the auxiliary variable with mean μ_X and variance σ_X^2 . The covariance between the auxiliary information (X) and the response to the sensitive question (Y) is σ_{XY} , and the correlation is ρ_{XY} .

$\hat{\mu}_Y$ is an estimate of the mean of Y in the absence of an auxiliary variable, μ_X is the known mean of the auxiliary variable, and \bar{x} is the observed mean of the auxiliary responses sampled. To the extent that \bar{x} is smaller than the known mean of X , then the ratio $\frac{\mu_X}{\bar{x}}$ will be greater than 1, and the estimator’s estimate will be bumped up to reflect the information contributed by the auxiliary variable. The reverse will be true when \bar{x} is larger than μ_X .

The Parker et al. (2024) mean estimator given in Equation (2) can be transformed into a ratio estimator [15]:

$$\hat{\mu}_{YR} = \left[\frac{1 - p_1}{p_2 - p_1} \bar{Z}_2 - \frac{1 - p_2}{p_2 - p_1} \bar{Z}_1 \right] \left[\frac{1}{2} \left(\frac{\mu_X}{\bar{x}_1} + \frac{\mu_X}{\bar{x}_2} \right) \right]. \tag{8}$$

The term $\left[\frac{1}{2} \left(\frac{\mu_X}{\bar{x}_1} + \frac{\mu_X}{\bar{x}_2} \right) \right]$ is the average ratio adjustment from the two sub-samples. Assuming the auxiliary information used is of high quality, it is reasonable to expect that, while slightly biased, this ratio will yield a “better” (having superior efficiency) estimate of the mean response to the sensitive question.

3.2. Bias

Unlike the MOET basic estimator in Equation (2), the proposed ratio estimator in Equation (8) is biased. The following expression represents this bias:

$$\text{Bias}(\hat{\mu}_{YR}) = E(\hat{\mu}_{YR}) - \mu_Y. \tag{9}$$

Using a first order Taylor expansion, this expression can be approximated by

$$\text{Bias}(\hat{\mu}_{YR}) = \frac{\mu_Y}{2} \left(\frac{1}{n_1} + \frac{1}{n_2} \right) \left(\frac{\sigma_X}{\mu_X} \right)^2 + \frac{\sigma_{XY}}{2\mu_X} \left[\frac{1}{n_2} \left(\frac{1 - p_1}{p_1 - p_2} \right) - \frac{1}{n_1} \left(\frac{1 - p_2}{p_1 - p_2} \right) \right], \tag{10}$$

and under an equally split sample assumption, this expression further simplifies to

$$\text{Bias}(\hat{\mu}_{YR}) = \frac{2\mu_Y}{n} \left(\frac{\sigma_X}{\mu_X} \right)^2 - \frac{\sigma_{XY}}{n\mu_X}. \tag{11}$$

It is clear that $\hat{\mu}_{YR}$ is asymptotically unbiased. Furthermore, Equation (11) can be rewritten as

$$\text{Bias}(\hat{\mu}_{YR}) = \left(\frac{\sigma_X}{\mu_X} \right) \left(\frac{2\mu_Y\sigma_X}{n\mu_X} - \frac{\rho_{XY}\sigma_Y}{n} \right), \tag{12}$$

which reveals that the bias will be small when the variance of the auxiliary variable is small relative to its mean. Figure 2 reflects on the behavior of bias relative to n and σ_X/μ_X .

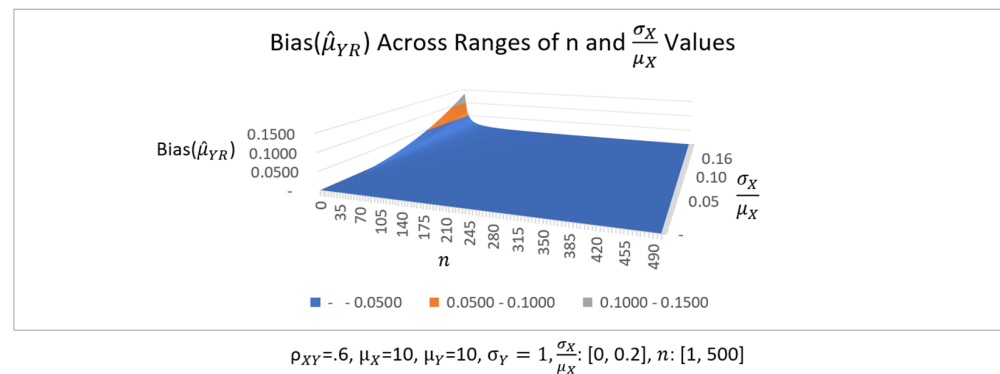


Figure 2. Bias of the ratio estimator.

The key conclusion of this analysis is that bias remains small, except in extreme scenarios. In the rear corner of the graph, the bias rises abruptly when, concurrently, n is small and σ_X/μ_X is large. But when either n is large or σ_X/μ_X is small, the bias remains small. For $n \geq 100$ and $\frac{\sigma_X}{\mu_X} \leq 0.2$, the bias is never greater than 0.007 (0.07% of μ_Y).

Bias is further studied numerically in Tables 1 and 2 of Section 6.

3.3. Efficiency of Ratio Estimator

For any mean estimator $\hat{\mu}$,

$$MSE(\hat{\mu}) = E[(\hat{\mu} - \mu)^2] \tag{13}$$

It follows that, for the proposed ratio estimator in Equation (8), MSE can be written

$$MSE(\hat{\mu}_{YR}) = E\left\{ \left[\frac{1}{2} \left(\frac{1-p_1}{p_2-p_1} \bar{Z}_2 - \frac{1-p_2}{p_2-p_1} \bar{Z}_1 \right) \left(\frac{\mu_X}{x_1} + \frac{\mu_X}{x_2} \right) - \mu_Y \right]^2 \right\}. \tag{14}$$

Assuming bivariate normality of (Y, X) and a first-order Taylor expansion, Equation (14) becomes

$$MSE(\hat{\mu}_{YR}) = \frac{\mu_Y^2}{4\mu_X^2} (\sigma_X^2) \left(\frac{1}{n_1} + \frac{1}{n_2} \right) + \left(\frac{1-p_2}{p_2-p_1} \right)^2 \left(\frac{\sigma_{Z_1}^2}{n_1} \right) + \left(\frac{1-p_1}{p_2-p_1} \right)^2 \left(\frac{\sigma_{Z_2}^2}{n_2} \right) + \frac{\mu_Y \sigma_{YX}}{\mu_X(p_2-p_1)} \left(\frac{1-p_2}{n_1} - \frac{1-p_1}{n_2} \right) + \frac{\mu_Y(1-p_1)(1-p_2)}{\mu_X(p_2-p_1)} W(1-\alpha)(\sigma_{YX} + \mu_R \mu_X) \left(\frac{1}{n_2} - \frac{1}{n_1} \right). \tag{15}$$

Under the equally split sample assumption, this expression simplifies to

$$MSE(\hat{\mu}_{YR}) = \left[\left(\frac{1-p_1}{p_2-p_1} \right)^2 \text{Var}(\bar{Z}_2) + \left(\frac{1-p_2}{p_2-p_1} \right)^2 \text{Var}(\bar{Z}_1) \right] + \frac{1}{n} \left(\frac{\mu_Y}{\mu_X} \right)^2 \sigma_X^2 - \frac{2}{n} \left(\frac{\mu_Y}{\mu_X} \right) \sigma_{XY}, \tag{16}$$

where

$$\text{Var}(\bar{Z}_i) = \frac{2}{n} \left\{ W[1-\lambda_i - A\phi_i] \sigma_S^2 + [W(1-\lambda_i)((1-A)\sigma_T^2 + 1) + 1-W](\sigma_Y^2 + \mu_Y^2) + W\lambda_i(\sigma_R^2 + \mu_R^2) - [\mu_Y + (\mu_R - \mu_Y)W\lambda_i]^2 \right\},$$

$$\lambda_i = (1-\alpha)(1-p_i), \quad i = 1, 2,$$

$$\phi_i = p_i(1-\alpha), \quad i = 1, 2,$$

$$p_1 \neq p_2.$$

Note that $\left[\left(\frac{1-p_1}{p_2-p_1} \right)^2 \text{Var}(\bar{Z}_2) + \left(\frac{1-p_2}{p_2-p_1} \right)^2 \text{Var}(\bar{Z}_1) \right]$, the first of three additive terms in Equation (16), exactly equals the MSE of the basic mean estimator proposed by Parker et al. (2024); see Equation (3) [15]. This portion of the MSE expression can be estimated as the weighted sum of the variances of observed \bar{Z}_i values, $i = 1, 2$. The remaining two terms in Equation (16) quantify the additional MSE attributable to the auxiliary

variable. μ_X is assumed to be known, but μ_Y and σ_{XY} are not. One can plug in parameter estimates for these parameters based on the collected data.

Alternatively, one can express the last two additive terms of Equation (16) using coefficients of variation (CV):

$$\begin{aligned} \frac{1}{n} \left(\frac{\mu_Y}{\mu_X} \right)^2 \sigma_X^2 - \frac{2}{n} \left(\frac{\mu_Y}{\mu_X} \right) \sigma_{XY} &= \frac{1}{n} \left(\frac{CV(X)}{CV(Y)} \right)^2 \sigma_Y^2 - \frac{2}{n} \left(\frac{CV(X)}{CV(Y)} \right) \rho_{YX} \sigma_Y^2 \\ &= \frac{1}{n} \left(\frac{CV(X)}{CV(Y)} \right) \sigma_Y^2 \left(\frac{CV(X)}{CV(Y)} - 2\rho_{YX} \right). \end{aligned} \tag{17}$$

This representation makes it clear that, under the standard assumption that $CV(X) = CV(Y)$, the last two terms of Equation (16) will sum to a negative number if $\rho_{YX} > 1/2$. The circumstance that $\rho_{YX} > 1/2$ is likely, in that auxiliary data are only used in cases where high correlation exists between auxiliary information and the question under study. It follows that the first portion of Equation (16) alone can be used to approximate the MSE of μ_{YR} and that this will always be a conservative estimate when $\rho_{YX} > 1/2$ and $CV(X) = CV(Y)$. Formally,

$$MSE(\widehat{\mu}_{YR}) = MSE(\widehat{\mu}_Y) = \frac{2}{n} \left[\left(\frac{1 - p_1}{p_2 - p_1} \right)^2 \text{Var}(Z_2) + \left(\frac{1 - p_2}{p_2 - p_1} \right)^2 \text{Var}(Z_1) \right], \tag{18}$$

where $\text{Var}(Z_i)$ can be estimated by the variance of the observed Z_i values. Following Table 2 in Section 6, additional discussion of this simplification is provided based on numerical demonstration.

3.4. Privacy of Ratio Estimator Model (∇)

The concept of privacy is critical to all RRT models, because it is the privacy offered by a RRT model that removes SDB influences (embarrassment, shame, even illegality) and allows the respondent to answer truthfully. However, beyond this, there are many ethical considerations related to privacy. Not only must respondents be educated to understand that their information is made private by RRT, but practitioners must closely follow procedures to avoid collecting private information in a way that could link an individual respondent to their response, especially as it relates to the proportion of respondents $(1 - W)$ who elect to answer the sensitive question without any means of RRT perturbation. No records that could be used to identify individual respondents, including even the scheduled time of their interview, should be documented along with the response. Institutional review board approval should be sought before launching any RRT-based study.

Yan et al. (2008) proposed the following measure, which is commonly used to evaluate the privacy provided by quantitative models [19]:

$$\nabla = E[(Z - Y)^2], \tag{19}$$

where Y represents a respondent’s true response to a sensitive question, and Z represents a respondent’s reported response (which may be scrambled). Parker et al. (2024) showed, based on Yan’s definition of privacy, that MOET model privacy, under an equally split subsample assumption, can be calculated by Equation (4) [15]. Under the standard assumption that auxiliary information does not impair model privacy, this expression for privacy will remain valid for the MOET model even when auxiliary information is collected.

3.5. Sensitivity Estimator (\widehat{W}_X)

Gupta et al. (2002) showed that, while the primary value of optionality is that it improves the efficiency of RRT estimates, it also has an important secondary value [10]. Specifically, $W \in [0, 1]$ can be used to quantify the sensitivity of a question. W values close

to one imply that a question is highly sensitive, while W values close to zero imply low sensitivity. To find an estimator for W that takes advantage of auxiliary information, we let

$$\mu_{YR} = \left[\frac{\mu_X}{\bar{x}} \right] \left[\frac{E[Z] - W(1 - \alpha)(1 - p)\mu_R}{1 - W(1 - \alpha)(1 - p)} \right]. \tag{20}$$

This formula can be applied separately to the two splits in a split sample as follows:

$$[1 - W(1 - \alpha)(1 - p_i)]\bar{x}\mu_{YR} = \mu_X[E[Z_i] - W(1 - \alpha)(1 - p_i)\mu_R], \quad i = 1, 2. \tag{21}$$

Estimating $E[Z_i]$ by \bar{Z}_i and μ_{YR} by $\hat{\mu}_{YR}$ in these expressions yields

$$[1 - \hat{W}(1 - \alpha)(1 - p_i)]\bar{x}\hat{\mu}_{YR} = \mu_X[\bar{Z}_i - \hat{W}(1 - \alpha)(1 - p_i)\mu_R], \quad i = 1, 2. \tag{22}$$

Finally, solving the two split sample equations simultaneously and calling this estimator \hat{W}_X , to signify that it takes auxiliary information into account, leads to

$$\hat{W}_X = \frac{\bar{x}\hat{\mu}_{YR}(\lambda_2 + \lambda_1) - \mu_X(\bar{Z}_1\lambda_2 + \bar{Z}_2\lambda_1)}{2\lambda_1\lambda_2(\bar{x}\hat{\mu}_{YR} - \mu_R\mu_X)}, \tag{23}$$

$$\begin{aligned} \lambda_1 &= (1 - \alpha)(1 - p_1), \\ \lambda_2 &= (1 - \alpha)(1 - p_2), \\ \mu_R &\neq \mu_Y, \\ p_1 &\neq p_2. \end{aligned}$$

Table 4 of Section 6 demonstrates the superiority of this sensitivity estimator to the basic sensitivity estimator that does not leverage auxiliary information in Equation (6).

4. The Impact of Auxiliary Information Quality on Estimation

Not all auxiliary information is perfect. In fact, frequently, auxiliary information will be so imperfect that it should not be used. This could be a result of issues like nonresponse and missing data, which typically plague online and mail-in surveys, but these issues are rare in RRT surveys, because RRT surveys are based on face-to-face interviews. More commonly, candidate auxiliary information may be inadequate because it does not correlate well enough with the sensitive variable and therefore will not assist in the estimation. In this section, we consider the impact of auxiliary information (with different levels of quality) on mean estimation by comparing the MOET mean ratio estimator that makes use of auxiliary information to the basic mean estimator, which does not. Throughout this section, we make the standard assumption that auxiliary information does not reduce model privacy. Since the ratio estimator will therefore have the same privacy level as the basic estimator, the two estimators can be compared based on efficiency alone.

It is intuitively clear that auxiliary information well correlated with the response to a sensitive question could be used to improve the estimation of the mean response to that question. It is furthermore reasonable to expect that, the more correlated X and Y are, the more accurate the estimate will be. That is, higher values of ρ_{XY} will always be preferred.

It is also reasonable to expect that an auxiliary variable with very little variability (σ_X near zero) will offer little prognostic value, as such auxiliary information will lead to a ratio estimator which ratio value will be close to 1, yielding no benefit. On the other hand, it seems logical that wildly variable auxiliary information (high σ_X) would be chaotic and would be of little value in estimating μ_Y . Based on these informal ruminations, we suspect that the ideal σ_X value will be neither extremely large nor extremely small. To validate these expectations mathematically, we rewrite Equation (16) above in the format below:

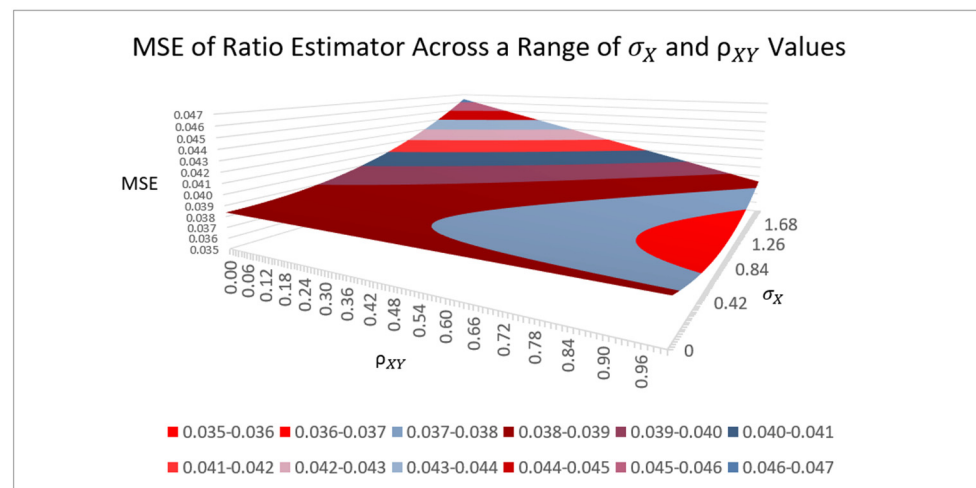
$$MSE(\hat{\mu}_{YR}) = MSE(\hat{\mu}_Y) + \frac{1}{n} \left(\frac{\mu_Y}{\mu_X} \right)^2 \sigma_X^2 - \frac{2}{n} \left(\frac{\mu_Y}{\mu_X} \right) \rho_{XY} \sigma_X \sigma_Y. \tag{24}$$

Note that, within the scope of this study, ρ_{XY} will always be a positive factor (the case of negative correlation can be handled using a conceptually analogous product estimator). As ρ_{XY} appears in the final, negative term of Equation (24), we can therefore confirm that large correlation values will be preferred in that they will reduce MSE (except in the unlikely scenario when this term is more than double the size of $MSE(\hat{\mu}_Y) + \frac{1}{n} \left(\frac{\mu_Y}{\mu_X}\right)^2 \sigma_X^2$). It is also clear that, if Equation (24) is considered as a function of σ_X with fixed ρ_{XY} , this function will be an upward-facing parabola, thereby confirming our expectation that the ideal σ_X will neither be extremely large nor extremely small. Basic calculus convinces us that maximum efficiency will be achieved when

$$\sigma_{X(\text{ideal})} = \left(\frac{\mu_X}{\mu_Y}\right) \rho_{XY} \sigma_Y. \tag{25}$$

This expression shows that the ideal auxiliary information will have a standard deviation σ_X close to the standard deviation of the response to the sensitive question σ_Y after adjustment for the relative sizes of the means of X and Y $\left(\frac{\mu_X}{\mu_Y}\right)$ and adjustment for correlation strength ρ_{XY} . Note that, in the unique circumstance where $\mu_X = \mu_Y$ and where $\rho_{XY} = 1$, the expression reduces to $\sigma_{X(\text{ideal})} = \sigma_Y$.

The graphical representation shown in Figure 3 demonstrates a concrete scenario and shows that, indeed, high correlation and “favorable σ_X ” (σ_X close to ideal) lead to low MSE .



$n=500, \mu_Y=10, \mu_X=10, \mu_S=0, \mu_T=1, \mu_R=10, \sigma_Y=1, \sigma_S=1, \sigma_T=1, \sigma_R=1, \sigma_X: [0, 2], \rho_{XY}: [0, 1]$

Figure 3. Ratio estimator efficiency.

The MSE when $\sigma_X = 0$ (the MSE value along the ρ_{XY} axis of the graph) is exactly the MSE of the basic estimator, because an auxiliary variable with zero variance will lead to $\frac{\mu_X}{x} = 1$, and so, the ratio estimator will “become” the basic estimator. Other points on the graph’s surface that rise above the basic estimator’s MSE value represent cases where the ratio estimator has higher MSE than the basic estimator. This happens most significantly when the correlation is small and σ_X is large, circumstances represented by the rear corner of the graph where ρ_{XY} is close to 0 and σ_X is close to 2. This result should not be surprising; it is obvious that the use of uncorrelated, wildly variable data as the basis of a ratio estimator would be ineffective. But for a set of data where the correlation is high and σ_X is favorable (near the value implied by Equation (25)), the ratio estimator will outperform the basic estimator. This occurrence is represented by points on the graph’s surface that fall below the basic estimator’s MSE value. The dark burgundy region represents correlation–variance combinations that lead to a ratio estimator MSE approximately equal to that of the basic estimator, implying that the two estimators will have equal efficiency. The light blue and

peach regions represent circumstances where the ratio estimator is superior to the basic estimator. For the scenario underlying the graph above, the maximum ratio estimator efficiency for highly correlated data will be achieved near $\sigma_X = \sigma_Y = 1$.

We embody the concepts above in the following relationship:

If $\rho_{XY} > \frac{1}{2} \left(\frac{\mu_Y}{\mu_X} \right)^2 \left(\frac{\sigma_X}{\sigma_Y} \right)$, then the ratio estimator will be more efficient than the basic estimator.

Clearly, high-quality auxiliary information may be used to improve the estimator accuracy. Alternatively, a researcher may use auxiliary information as a means of reducing the sample size. To calculate this reduction in sample size (Δ_n), we write the MSE of the ratio estimator as a function of $n - \Delta_n$:

$$MSE_{\hat{\mu}_{YR}}(n - \Delta_n) = \frac{2}{n - \Delta_n} \left[\left(\frac{1 - p_1}{p_2 - p_1} \right)^2 \sigma_{Z2}^2 + \left(\frac{1 - p_2}{p_2 - p_1} \right)^2 \sigma_{Z1}^2 \right] + \frac{1}{n - \Delta_n} \left(\frac{\mu_Y}{\mu_X} \right)^2 \sigma_X^2 - \frac{2}{n - \Delta_n} \left(\frac{\mu_Y}{\mu_X} \right) \rho_{XY} \sigma_X \sigma_Y \tag{26}$$

The MSE of the basic estimator can be written

$$MSE_{\hat{\mu}_Y}(n) = \frac{2}{n} \left[\left(\frac{1 - p_1}{p_2 - p_1} \right)^2 \sigma_{Z2}^2 + \left(\frac{1 - p_2}{p_2 - p_1} \right)^2 \sigma_{Z1}^2 \right]. \tag{27}$$

Equating the right-hand sides of Equations (26) and (27) and solving for Δ_n , we identify the reduction in sample size that can be achieved through use of an auxiliary variable:

$$\Delta_n = \left\lfloor \frac{2\mu_X\mu_Y\sigma_{XY} - \mu_Y^2\sigma_X^2}{\mu_X^2MSE(\hat{\mu}_Y)} \right\rfloor. \tag{28}$$

5. The Impact of Auxiliary Information on Privacy

Lanke (1976) recognized the critical importance of privacy to RRT [20], saying that, formerly, statisticians were “considering matters from the statistician’s point of view only” and that, from the general public’s point of view, a much more central question would be “to what extent do the different [RRT] methods protect the privacy of interviewees”. In his study, Lanke proposed a means of quantifying the level of privacy provided by binary RRT models. Yan et al. (2008) later proposed a means of measuring quantitative model RRT privacy in Equation (19) [19].

In Section 4, we relied on the standard assumption that the collection of auxiliary information does not reduce privacy, but this may not be true. Considering that auxiliary information is chosen because it correlates with the response to a sensitive question, auxiliary information may, in fact, compromise privacy. For example, in a case where X and Y were perfectly correlated, the knowledge of X would lead directly to knowledge of Y . This reality further complicates the ethical privacy considerations alluded to in Section 3.4. To the extent that auxiliary information, along with survey information, can help to identify individual respondents and therefore expose their answer to the sensitive question, such auxiliary information should not be collected.

The quantitative MOET RRT model proposed by Parker et al. (2024) [15], which does not involve the collection of auxiliary information, is represented by the exhibit displayed at the beginning of Section 2 in Figure 1. The output of that model is the scrambled response Z , and the model’s privacy is a function of the true response to the sensitive question (Y), so ∇ for this model may be represented $\nabla(Y)$.

If the collection of auxiliary information is part of the RRT model, then the model can instead be represented as in Figure 4.

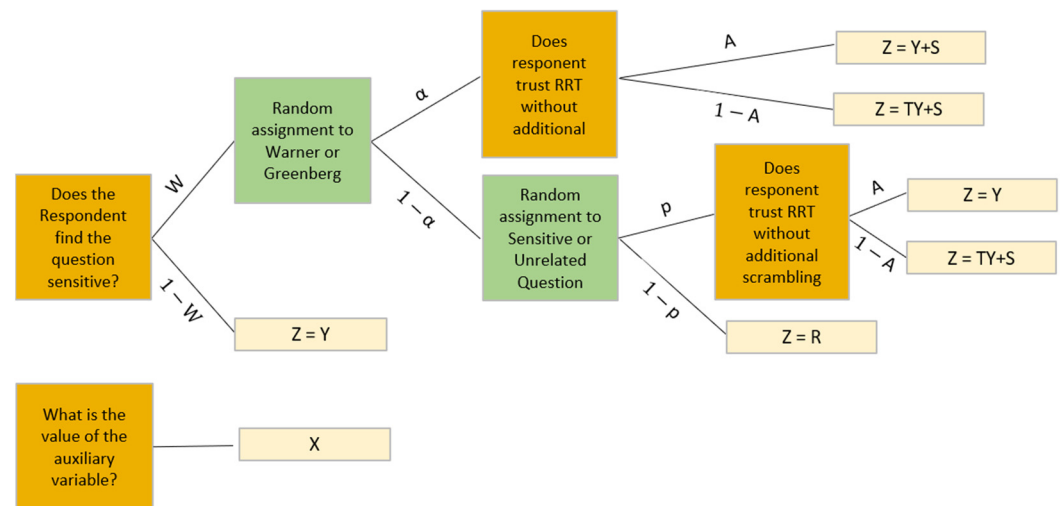


Figure 4. Diagram of the MOET model with auxiliary information.

The bivariate output of this model is (X, Y) , and the privacy of this model may be represented $\nabla(X, Y)$. The exact value of $\nabla(X, Y)$ is not obvious, but we do know that the privacy of this model can be no greater than the privacy level of the basic “no auxiliary information” model, because collecting additional information about a respondent could not possibly increase the respondent’s privacy. Moreover, privacy is always floored at zero. Therefore, we have that

$$0 \leq \nabla(X, Y) \leq \nabla(Y). \tag{29}$$

In order to study the role of privacy reduction due to auxiliary data, we define ϕ as the percentage reduction in privacy that occurs when auxiliary data are collected.

It follows that

$$\nabla(X, Y) = (1 - \phi)\nabla(Y). \tag{30}$$

At this point, we will make use of the unified measure (UM) proposed by Gupta et al. (2018) to assess the overall model value, taking into account the opposing considerations of efficiency and privacy [16]. We represent the UM metric by the symbol δ :

$$\delta(\hat{\mu}_Y) = \frac{\text{MSE}(\hat{\mu}_Y)}{\nabla(Y)}. \tag{31}$$

Similarly, when auxiliary information is collected for ratio estimation:

$$\delta(\hat{\mu}_{YR}) = \frac{\text{MSE}(\hat{\mu}_{YR})}{\nabla(X, Y)}. \tag{32}$$

As this metric is made small by either small MSE values or large privacy values, smaller values of UM are preferred. In the illustrative figure labeled Figure 5, we show the UM calculated across a range of Φ values for both the basic estimator and the ratio estimator.

The parameters underlying Figure 5 were chosen, because they illustrate the concept of UM “crossover” well. A more extensive numerical analysis across a broad range of values is provided in Table 3 of Section 6. Figure 5 shows, as expected, that the UM value for the ratio estimator is best (smallest) when Φ is small (that is, when the collection of auxiliary information does not reduce privacy a lot). The UM gets worse (higher) from left to right, as the collection of auxiliary data undermines privacy more and more. The UM of the “basic” estimator (which does not depend on auxiliary information) is, of course, unaffected by Φ . In the scenario underlying the figure, the basic estimator will be superior to the ratio

estimator overall when the collection of auxiliary information reduces respondent privacy by no more than $\Phi_{\text{crossover}} \approx 11\%$. The identity of this crossover point may be conceived as

$$\Phi_{\text{crossover}} = \frac{2\mu_Y\mu_X\sigma_X\sigma_Y\rho_{XY} - \mu_Y^2\sigma_X^2}{n\mu_X^2\text{MSE}(\hat{\mu}_Y)}. \tag{33}$$

This expression cannot be exactly calculated, as it relies on μ_Y , σ_Y , and $\text{MSE}(\hat{\mu}_Y)$, which are unknown. However, if reasonable a priori estimates for these values are available and ρ_{XY} is assumed to be less than 1, it follows that

$$\Phi_{\text{crossover}(\text{approx})} \leq \frac{2\mu_{Y_{ap}}\mu_X\sigma_X\sigma_{Y_{ap}} - \mu_{Y_{ap}}^2\sigma_X^2}{n\mu_X^2(\text{MSE}(\hat{\mu}_Y)_{ap})}, \tag{34}$$

where $\mu_{Y_{ap}}$, $\sigma_{Y_{ap}}$, and $\text{MSE}(\hat{\mu}_Y)_{ap}$ are a priori estimates.

Note that $\Phi_{\text{crossover}}$ will be small when n is large. This means that even small losses of privacy will undermine the overall estimator quality (UM) if the sample size is large. This makes sense, as estimates without the help of auxiliary information will generally already be accurate in these circumstances.

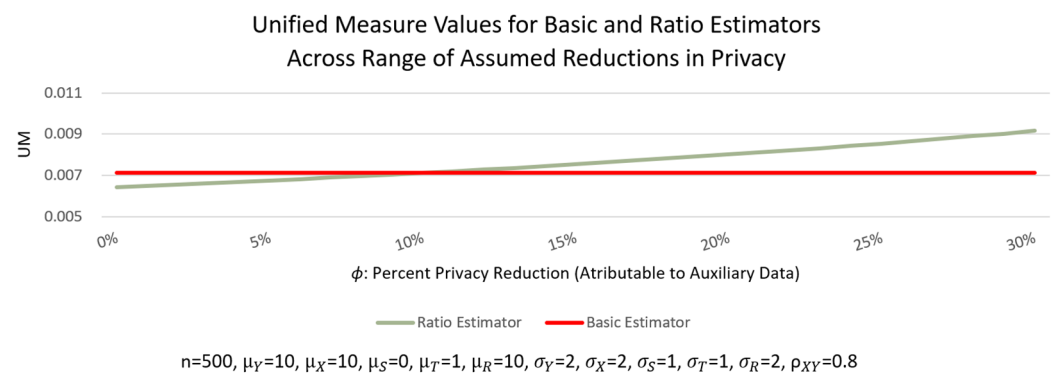


Figure 5. Overall performance of estimators when auxiliary information reduces privacy.

6. Simulations

In this section, we provide four tables that present theoretical and simulated values associated with the MOET model when auxiliary information is used. Each row of each table represents the output for a particular set of parameter values. Except when the concepts underlying the tables call for unique values (the reasoning for such selections will be specified along with each table), all four tables are based on a common set of parameter values and modeling specifications. The value $\mu_Y = 10$ is chosen arbitrarily, as, in isolation, its value is not important to the results or conclusions. The value of $\sigma_Y = 5$ is set relative to μ_Y and was chosen to result in a similar σ_Y/μ_Y value to other recent related papers such as Gupta et al. (2022) and Parker et al. (2024) [15,17]. We set μ_X equal to μ_Y , because, in isolation, the size of μ_X is unimportant, and setting it equal to μ_Y facilitates easy comparisons. It is the relationship between σ_Y/μ_Y and σ_X/μ_X that is important, so σ_X is tested at values below, equal to, and above σ_Y (3, 5, and 7). ρ_{XY} is a key quality of auxiliary information, so it too is calculated across low, medium, and high values (0.50, 0.75, and 0.95). The values of μ_R and σ_R are set equal to the values of μ_Y and σ_Y , because answers to the Greenberg unrelated question should mimic answers to the sensitive question. $\mu_S = 0, \mu_T = 1, \sigma_S = 1, \text{ and } \sigma_T = 1$ are standard choices for the additive and multiplicative scrambling values that facilitate simplified expectation calculations, and the choices of $\alpha = 0.15, p_1 = 0.85, \text{ and } p_2 = 0.15$ are chosen because they lead to high model efficiency, per Parker et al. (2024) [15]. When not otherwise specified, W and A are set to the moderate values of 0.90 and 0.95 respectively.

Table 1 is provided for two reasons: to demonstrate the fact that the values of the estimators derived in this study fall close to theoretical values across a wide variety of simulated scenarios and to show that bias remains small in realistic scenarios. Table 1A assumes a high sensitivity ($W = 0.9$) and low trust ($A = 0.85$) scenario, while Table 1B reflects low sensitivity ($W = 0.5$) and high trust ($A = 0.98$). These values were chosen because it makes sense that increased question sensitivity would undermine respondent trust. In addition to running scenarios across ranges of σ_X and ρ_{XY} values (the two factors that together characterize auxiliary information quality), a range of Y values are considered (500, 250, and 100) to enable the study of bias across sample sizes. The scenarios in Table 1A,B are run 100,000 times each to show that the estimation procedure is robust even for small quantities like bias.

In Table 1A,B, $\hat{\mu}_{YR}$ represents the estimated mean response to the sensitive question based on the ratio estimator. The letters MSE stand for mean squared error (a measurement of estimator efficiency), the symbol ∇ represents privacy, and the symbol δ represents the unified measure. The subscripts T and E indicate theoretical and empirical results, respectively. In the boxed row of Table 1B above, the theoretical and empirical values of MSE are 0.1470 and 0.1473, the privacy values are 22.85 and 22.91, the UM values are 0.0064 and 0.0064, and the bias values are 0.0021 and 0.0020. Across the board, Table 1A,B show that the empirical values match the theoretical values closely.

Table 1A,B also show that the bias of the ratio estimator is generally small. This is made most clear by the final column of the table ($\text{Bias}(\hat{\mu}_{YR})_T / \mu_Y$), where bias is shown to be a small percentage of the mean response to the sensitive question. It is only when, concurrently, n is small and σ_X is large that bias becomes more significant. But large σ_X , which leads to large $CV(X)$ values that are dissimilar from $CV(Y)$, represent a circumstance when the superiority of the ratio estimator over the basic mean estimator is not guaranteed.

Table 2 is shown to demonstrate the performance of the MOET ratio estimator relative to the performance of the MOET basic estimator across a variety of auxiliary information scenarios. Importantly, σ_X and ρ_{XY} together define the auxiliary information quality, so all permutations of low, medium, and high values of ρ_{XY} were considered, and a broad range of σ_X values (1, 3, 5, 7 and 9) were considered. The scenarios in Tables 2–4 are run 10,000 times each.

Table 1. Ratio estimator using MOET—theoretical vs. empirical values. (A) ($A = 0.85, W = 0.90, \alpha = 0.15, p_1 = 0.85, p_2 = 0.15, \mu_Y = 10, \mu_X = 10, \mu_R = 10, \mu_S = 0, \mu_T = 1, \sigma_Y = 5, \sigma_R = 5, \sigma_S = 1, \text{ and } \sigma_T = 1$). (B) ($A = 0.98, W = 0.50, \alpha = 0.15, p_1 = 0.85, p_2 = 0.15, \mu_Y = 10, \mu_X = 10, \mu_R = 10, \mu_S = 0, \mu_T = 1, \sigma_Y = 5, \sigma_R = 5, \sigma_S = 1, \text{ and } \sigma_T = 1$).

(A)												
n	ρ_{XY}	σ_X	$\hat{\mu}_{YR}$	$\text{MSE}(\hat{\mu}_{YR})$	$\text{MSE}(\hat{\mu}_{YR})_E$	$\nabla (\hat{\mu}_{YR})_T$	$\nabla (\hat{\mu}_{YR})_E$	$\delta(\hat{\mu}_{YR})_T$	$\delta(\hat{\mu}_{YR})_E$	$\text{Bias}(\hat{\mu}_{YR})_T$	$\text{Bias}(\hat{\mu}_{YR})_E$	$\frac{\text{Bias}(\hat{\mu}_{YR})_T}{\mu_Y}$
500	0.95	3	10.0013	0.2021	0.2027	32.25	32.37	0.0063	0.0063	0.0008	0.0007	0.01%
500	0.95	5	10.0038	0.1961	0.1973	32.25	32.37	0.0061	0.0061	0.0053	0.0054	0.05%
500	0.95	7	10.0128	0.2061	0.2079	32.25	32.37	0.0064	0.0064	0.0130	0.0133	0.13%
500	0.75	3	10.0023	0.2141	0.2134	32.25	32.37	0.0066	0.0066	0.0014	0.0005	0.01%
500	0.75	5	10.0056	0.2161	0.2176	32.25	32.37	0.0067	0.0067	0.0063	0.0072	0.06%
500	0.75	7	10.0156	0.2341	0.2355	32.25	32.37	0.0073	0.0073	0.0144	0.0149	0.14%
500	0.50	3	10.0025	0.2291	0.2309	32.25	32.37	0.0071	0.0071	0.0021	0.0020	0.02%
500	0.50	5	10.0073	0.2411	0.2421	32.25	32.37	0.0075	0.0075	0.0075	0.0086	0.08%
500	0.50	7	10.0181	0.2691	0.2709	32.25	32.37	0.0083	0.0084	0.0161	0.0161	0.16%
250	0.95	3	10.0021	0.4043	0.4045	32.25	32.49	0.0125	0.0125	0.0015	0.0024	0.02%
250	0.95	5	10.0130	0.3923	0.3961	32.25	32.49	0.0122	0.0122	0.0105	0.0096	0.11%
250	0.95	7	10.0254	0.4123	0.4204	32.25	32.50	0.0128	0.0129	0.0259	0.0251	0.26%
250	0.75	3	10.0056	0.4283	0.4275	32.25	32.50	0.0133	0.0132	0.0027	0.0026	0.03%
250	0.75	5	10.0092	0.4323	0.4364	32.25	32.49	0.0134	0.0134	0.0125	0.0126	0.13%
250	0.75	7	10.0289	0.4683	0.4779	32.25	32.49	0.0145	0.0147	0.0287	0.0281	0.29%
250	0.50	3	10.0048	0.4583	0.4580	32.25	32.49	0.0142	0.0141	0.0042	0.0042	0.04%
250	0.50	5	10.0124	0.4823	0.4819	32.25	32.49	0.0150	0.0148	0.0150	0.0129	0.15%
250	0.50	7	10.0360	0.5383	0.5499	32.25	32.50	0.0167	0.0169	0.0322	0.0320	0.32%

Table 1. Cont.

50	0.95	3	10.0112	2.0214	2.0241	32.25	33.48	0.0627	0.0605	0.0075	0.0077	0.08%
50	0.95	5	10.0552	1.9614	2.0400	32.25	33.49	0.0608	0.0609	0.0525	0.0519	0.53%
50	0.95	7	10.1353	2.0614	2.2505	32.25	33.49	0.0639	0.0672	0.1295	0.1363	1.30%
50	0.75	3	10.0085	2.1414	2.1604	32.25	33.48	0.0664	0.0645	0.0135	0.0145	0.14%
50	0.75	5	10.0664	2.1614	2.2312	32.25	33.49	0.0670	0.0666	0.0625	0.0613	0.63%
50	0.75	7	10.1432	2.3414	2.5277	32.25	33.48	0.0726	0.0755	0.1435	0.1446	1.44%
50	0.50	3	10.0154	2.2914	2.3177	32.25	33.45	0.0711	0.0693	0.0210	0.0215	0.21%
50	0.50	5	10.0781	2.4114	2.4964	32.25	33.49	0.0748	0.0745	0.0750	0.0758	0.75%
50	0.50	7	10.1580	2.6914	2.9396	32.25	33.46	0.0835	0.0879	0.1610	0.1657	1.61%
(B)												
n	ρ_{XY}	σ_X	$\hat{\mu}_{YR}$	$MSE(\hat{\mu}_{YR})$	$MSE(\hat{\mu}_{YR})_E$	$\nabla(\hat{\mu}_{YR})_T$	$\nabla(\hat{\mu}_{YR})_E$	$\delta(\hat{\mu}_{YR})_T$	$\delta(\hat{\mu}_{YR})_E$	$Bias(\hat{\mu}_{YR})_T$	$Bias(\hat{\mu}_{YR})_E$	$\frac{Bias(\hat{\mu}_{YR})_T}{\mu_Y}$
500	0.95	3	10.0018	0.1200	0.1205	22.85	22.92	0.0053	0.0053	0.0008	0.0007	0.01%
500	0.95	5	10.0036	0.1140	0.1146	22.85	22.91	0.0050	0.0050	0.0053	0.0054	0.05%
500	0.95	7	10.0132	0.1240	0.1255	22.85	22.92	0.0054	0.0055	0.0130	0.0133	0.13%
500	0.75	3	10.0017	0.1320	0.1320	22.85	22.92	0.0058	0.0058	0.0014	0.0005	0.01%
500	0.75	5	10.0048	0.1340	0.1349	22.85	22.92	0.0059	0.0059	0.0063	0.0072	0.06%
500	0.75	7	10.0156	0.1520	0.1523	22.85	22.91	0.0067	0.0066	0.0144	0.0149	0.14%
500	0.50	3	10.0021	0.1470	0.1473	22.85	22.91	0.0064	0.0064	0.0021	0.0020	0.02%
500	0.50	5	10.0078	0.1590	0.1596	22.85	22.92	0.0070	0.0070	0.0075	0.0087	0.08%
500	0.50	7	10.0183	0.1870	0.1903	22.85	22.91	0.0082	0.0083	0.0161	0.0161	0.16%
250	0.95	3	10.0009	0.2401	0.2416	22.85	22.98	0.0105	0.0105	0.0015	0.0024	0.02%
250	0.95	5	10.0135	0.2281	0.2297	22.85	22.98	0.0100	0.0100	0.0105	0.0096	0.11%
250	0.95	7	10.0219	0.2481	0.2534	22.85	22.99	0.0109	0.0110	0.0259	0.0251	0.26%
250	0.75	3	10.0035	0.2641	0.2647	22.85	22.99	0.0116	0.0115	0.0027	0.0025	0.03%
250	0.75	5	10.0082	0.2681	0.2715	22.85	22.99	0.0117	0.0118	0.0125	0.0125	0.13%
250	0.75	7	10.0314	0.3041	0.3140	22.85	22.99	0.0133	0.0137	0.0287	0.0282	0.29%
250	0.50	3	10.0044	0.2941	0.2967	22.85	22.99	0.0129	0.0129	0.0042	0.0042	0.04%
250	0.50	5	10.0137	0.3181	0.3195	22.85	22.98	0.0139	0.0139	0.0150	0.0129	0.15%
250	0.50	7	10.0343	0.3741	0.3837	22.85	22.99	0.0164	0.0167	0.0322	0.0319	0.32%
50	0.95	3	10.0083	1.2004	1.2145	22.85	23.54	0.0525	0.0516	0.0075	0.0078	0.08%
50	0.95	5	10.0506	1.1404	1.2113	22.85	23.55	0.0499	0.0514	0.0525	0.0521	0.53%
50	0.95	7	10.1361	1.2404	1.4213	22.85	23.54	0.0543	0.0604	0.1295	0.1373	1.30%
50	0.75	3	10.0133	1.3204	1.3399	22.85	23.55	0.0578	0.0569	0.0135	0.0146	0.14%
50	0.75	5	10.0675	1.3404	1.4083	22.85	23.55	0.0587	0.0598	0.0625	0.0615	0.63%
50	0.75	7	10.1468	1.5204	1.7068	22.85	23.55	0.0665	0.0725	0.1435	0.1451	1.44%
50	0.50	3	10.0192	1.4704	1.4861	22.85	23.53	0.0644	0.0632	0.0210	0.0213	0.21%
50	0.50	5	10.0813	1.5904	1.6685	22.85	23.54	0.0696	0.0709	0.0750	0.0764	0.75%
50	0.50	7	10.1611	1.8704	2.0808	22.85	23.54	0.0819	0.0884	0.1610	0.1650	1.61%

Red border box referenced in following text to demonstrate match between theoretical and empirical values.

Table 2. Comparison between the ratio estimator and basic estimator using MOET. ($A = 0.95$, $W = 0.90$, $\alpha = 0.15$, $p_1 = 0.85$, $p_2 = 0.15$, $\mu_Y = 10$, $\mu_X = 10$, $\mu_R = 10$, $\mu_S = 0$, $\mu_T = 1$, $\sigma_Y = 5$, $\sigma_R = 5$, $\sigma_S = 1$, $\sigma_T = 1$, and $n = 500$).

Scenario			Basic Estimator			Ratio Estimator			$\frac{1}{2} \left(\frac{\mu_Y}{\mu_X} \right)^2 \left(\frac{\sigma_X}{\sigma_Y} \right)$	$\rho_{XY} > \frac{1}{2} \left(\frac{\mu_Y}{\mu_X} \right)^2 \left(\frac{\sigma_X}{\sigma_Y} \right)$	$\sigma_X(\text{ideal})$	$\frac{MSE(\hat{\mu}_{YR})}{MSE(\hat{\mu}_{YR})_T}$	Δn
ρ_{XY}	σ_X	$CV(X)$	$\hat{\mu}_Y$	$MSE(\hat{\mu}_Y)_T$	$MSE(\hat{\mu}_{YR})$	$\hat{\mu}_Y$	$MSE(\hat{\mu}_{YR})_T$	$MSE(\hat{\mu}_{YR})$					
0.95	1	0.1	10.002	0.1823	0.1823	10.002	0.1653	0.1823	0.10	yes	4.50	91%	46
0.95	3	0.3	10.002	0.1823	0.1822	10.002	0.1433	0.1822	0.30	yes	4.50	79%	106
0.95	5	0.5	9.995	0.1823	0.1819	10.004	0.1373	0.1819	0.50	yes	4.50	75%	123
0.95	7	0.7	10.002	0.1823	0.1819	10.014	0.1473	0.1819	0.70	yes	4.50	81%	96
0.95	9	0.9	10.006	0.1823	0.1823	10.026	0.1733	0.1823	0.90	yes	4.50	95%	24
0.75	1	0.1	9.999	0.1823	0.1821	9.998	0.1693	0.1821	0.10	yes	3.75	93%	35
0.75	3	0.3	10.003	0.1823	0.1827	10.003	0.1553	0.1827	0.30	yes	3.75	85%	74
0.75	5	0.5	10.000	0.1823	0.1818	10.007	0.1573	0.1818	0.50	yes	3.75	86%	68
0.75	7	0.7	10.004	0.1823	0.1816	10.012	0.1753	0.1816	0.70	yes	3.75	96%	19
0.75	9	0.9	10.006	0.1823	0.1820	10.031	0.2093	0.1820	0.90	no	3.75	115%	0
0.50	1	0.1	10.002	0.1823	0.1820	10.002	0.1743	0.1820	0.10	yes	2.50	96%	21
0.50	3	0.3	10.000	0.1823	0.1820	10.002	0.1703	0.1820	0.30	yes	2.50	93%	32
0.50	5	0.5	9.994	0.1823	0.1820	10.002	0.1823	0.1820	0.50	no	2.50	100%	0
0.50	7	0.7	9.996	0.1823	0.1815	10.015	0.2103	0.1815	0.70	no	2.50	115%	0
0.50	9	0.9	9.995	0.1823	0.1820	10.020	0.2543	0.1820	0.90	no	2.50	139%	0

Table 2 compares the ratio estimator performance to the basic mean estimator performance. We can see that, unlike the basic estimator, the ratio estimator is slightly biased. We also see that the ratio estimator should be used only when auxiliary information is of sufficiently high quality to support its use (see Section 4 for the discussion of high-quality auxiliary data). The eleven rows of the table where the ratio estimator is superior to the basic estimator are identified by the word “yes”. These are the scenarios that pass the superiority condition derived in Section 4: $\rho_{XY} > \frac{1}{2} \left(\frac{\mu_Y}{\mu_X} \right)^2 \left(\frac{\sigma_X}{\sigma_Y} \right)$.

Table 2 further shows that the ratio estimator tends to be most effective when the correlation is high and σ_X is “favorable” (close to “ $\sigma_{X(\text{ideal})}$ ” per Equation (25)). Note, for example, that when $\rho_{XY} = 0.95$ and $\sigma_X = 5$, the ratio estimator’s MSE is 25% lower than that of the basic estimator. But when rho is 0.5 and $\sigma_X = 9$, the ratio estimator’s MSE is 139% of that of the basic estimator. The “ $\sigma_{X(\text{ideal})}$ ” column represents the σ_X that yields the lowest MSE for the scenario. The $\text{MSE}(\widehat{\mu}_{YR})_T / \text{MSE}(\widehat{\mu}_Y)_T$ column shows the ratio estimator MSE as a percentage of the basic estimator MSE. The Δn column shows the number of units by which a researcher could reduce his sample size while retaining a UM equal to that of this basic estimator.

Lastly, based on Table 2, we can consider the performance of the MSE approximation suggested in Section 3.3. The estimator of the ratio estimator’s MSE was represented by $\text{MSE}(\widehat{\mu}_{YR})$, per Equation (18), and we can observe the conservative nature of this estimator by comparing $\text{MSE}(\widehat{\mu}_{YR})$ to $\text{MSE}(\widehat{\mu}_{YR})_T$ (that is, comparing the simplified estimation of MSE to its true value). As expected, it is specifically in the circumstances where $\rho_{YX} > 1/2$ and $\text{CV}(X)$ is close to $\text{CV}(Y) = 0.50$ that $\text{MSE}(\widehat{\mu}_{YR})$ provides a reliably conservative estimate of the ratio estimator’s MSE. The estimator becomes less precise and more conservative for higher values of ρ_{YX} . When the correlation is low and $\text{CV}(X)$ is not close to $\text{CV}(Y)$, the simplified estimator is no longer reliable. But these are circumstances when ratio estimation should not be used in the first place.

Table 3. Overall (UM) performance of mean estimators when the collection of auxiliary information reduces privacy. ($A = 0.95, W = 0.90, \alpha = 0.15, p_1 = 0.85, p_2 = 0.15, \mu_Y = 10, \mu_X = 10, \mu_R = 10, \mu_S = 0, \mu_T = 1, \sigma_Y = 5, \sigma_R = 5, \sigma_S = 1, \sigma_T = 1$, and $n = 500$).

Scenario			Basic Estimator						Ratio Estimator				
ρ_{XY}	σ_X	ϕ	$\widehat{\mu}_Y$	$\text{MSE}(\widehat{\mu}_Y)_T$	$\text{MSE}(\widehat{\mu}_Y)$	$\nabla(\widehat{\mu}_Y)_T$	$\delta(\widehat{\mu}_Y)_T$	$\widehat{\mu}_{YR}$	$\text{MSE}(\widehat{\mu}_{YR})_T$	$\text{MSE}(\widehat{\mu}_{YR})$	$\nabla(\widehat{\mu}_{YR})_T$	$\delta(\widehat{\mu}_{YR})_T$	$\delta(\widehat{\mu}_{YR}) < \delta(\widehat{\mu}_Y)$
0.95	7	0%	10.001	0.1823	0.1819	25.0	0.0073	10.013	0.1473	0.1819	25.0	0.0059	yes
0.95	7	5%	9.996	0.1823	0.1818	25.0	0.0073	10.011	0.1473	0.1818	23.8	0.0062	yes
0.95	7	10%	9.994	0.1823	0.1818	25.0	0.0073	10.009	0.1473	0.1818	22.5	0.0065	yes
0.95	7	15%	10.002	0.1823	0.1820	25.0	0.0073	10.013	0.1473	0.1820	21.3	0.0069	yes
0.95	7	20%	10.001	0.1823	0.1820	25.0	0.0073	10.018	0.1473	0.1820	20.0	0.0074	no
0.75	5	0%	9.998	0.1823	0.1818	25.0	0.0073	10.003	0.1573	0.1818	25.0	0.0063	yes
0.75	5	5%	10.000	0.1823	0.1820	25.0	0.0073	10.004	0.1573	0.1820	23.8	0.0066	yes
0.75	5	10%	9.993	0.1823	0.1818	25.0	0.0073	10.001	0.1573	0.1818	22.5	0.0070	yes
0.75	5	15%	9.994	0.1823	0.1819	25.0	0.0073	9.999	0.1573	0.1819	21.3	0.0074	no
0.75	5	20%	10.002	0.1823	0.1820	25.0	0.0073	10.008	0.1573	0.1820	20.0	0.0079	no
0.50	3	0%	9.995	0.1823	0.1814	25.0	0.0073	9.997	0.1703	0.1814	25.0	0.0068	yes
0.50	3	5%	10.004	0.1823	0.1822	25.0	0.0073	10.006	0.1703	0.1822	23.8	0.0072	yes
0.50	3	10%	9.995	0.1823	0.1813	25.0	0.0073	9.997	0.1703	0.1813	22.5	0.0076	no
0.50	3	15%	10.001	0.1823	0.1817	25.0	0.0073	10.003	0.1703	0.1817	21.3	0.0080	no
0.50	3	20%	9.998	0.1823	0.1825	25.0	0.0073	10.002	0.1703	0.1825	20.0	0.0085	no

Table 3 studies the behavior of the basic and ratio estimators in the possible real-life scenario where the collection of auxiliary information does reduce privacy (see Section 5 of this study). As usual, we consider a range of ρ_{XY} and σ_X values as a means of studying auxiliary data with various levels of quality. We also consider a range of percentage reductions in privacy (Φ), as privacy reduction is the focus of this table, and we calculate UM for all scenarios to quantify the overall model quality. While the basic estimator is indifferent to the collection of auxiliary information (UM is constant across all scenarios), the privacy associated with the ratio estimator declines as Φ rises, causing the UM (δ) to become ever less favorable (higher). Whenever the reduction in privacy is less than the $\Phi_{\text{crossover}}$ value per Equation (33), the UM for the ratio estimator is greater than the UM for the basic estimator, and we can claim that the ratio estimator is superior overall.

For Table 4, $\mu_R = 8$ rather than the usual assumption $\mu_R = \mu_Y = 10$. This was done because the estimators \widehat{W} in Equation (6) and \widehat{W}_X in Equation (23) become unstable when $\mu_R = \mu_Y$, as can be seen by observing their formulas. The sensitivity estimator proposed in this paper that incorporates auxiliary information outperforms the estimator proposed by Parker et al. (2024) given in every scenario (relative error is smaller) [15]. Note that both \widehat{W} and \widehat{W}_X are most accurate when the sensitivity levels are high; these are the circumstances when RRT is most valuable.

Table 4. Sensitivity estimation—compare the basic estimator to the estimator that incorporates auxiliary information. ($A = 0.95$, $\alpha = 0.15$, $p_1 = 0.85$, $p_2 = 0.15$, $\mu_Y = 10$, $\mu_X = 10$, $\mu_R = 8$, $\mu_S = 0$, $\mu_T = 1$, $\sigma_Y = 5$, $\sigma_R = 5$, $\sigma_S = 1$, and $\sigma_T = 1$).

Scenario		Basic Sensitivity Estimator			Sensitivity Estimator that Incorporates Auxiliary Information		
n	W	\hat{W}	Absolute Error	Relative Error	\hat{W}_X	Absolute Error	Relative Error
500	0.90	0.8630	−0.037	4.1%	0.8732	−0.027	3.0%
500	0.70	0.6554	−0.045	6.4%	0.6667	−0.033	4.8%
500	0.30	0.2374	−0.063	20.9%	0.2499	−0.050	16.7%
250	0.90	0.8172	−0.083	9.2%	0.8357	−0.064	7.1%
250	0.70	0.5746	−0.125	17.9%	0.6165	−0.083	11.9%
250	0.30	0.1483	−0.152	50.6%	0.2473	−0.053	17.6%
50	0.90	0.8305	−0.070	7.7%	1.4627	0.563	62.5%
50	0.70	36.6184	35.918	5131.2%	2.1976	1.498	213.9%
50	0.30	8.3543	8.054	2684.8%	1.8016	1.502	500.5%

7. Conclusions

The use of auxiliary information, through the implementation of a ratio estimator, will improve the efficiency of mean estimation, provided the auxiliary information is of sufficiently high quality. In some scenarios shown in this study, efficiency improvements of 25% were observed.

In Section 3.3, we showed that, in the expected circumstances, $\rho_{YX} > 1/2$ and $CV(X) = CV(Y)$, the efficiency of the ratio mean estimator can be conservatively estimated using the MSE estimator for the basic mean estimator in Equation (3), because the incremental MSE associated with auxiliary data will be negative.

In Section 5, we used the unified measure (which takes both efficiency and privacy into account) to show that the overall quality of the ratio mean estimator is greater than that of the basic mean estimator, provided the auxiliary information used does not impair respondent privacy. But we also noted that, in some cases, auxiliary information does undermine privacy. In such cases, the use of auxiliary information should be considered carefully from an ethical standpoint before being collected. Even when no ethical breaches are identified, it is possible that privacy loss can undermine UM to the point that the overall quality of the ratio estimator is less than that of the basic estimator.

Author Contributions: Initial Conceptualization, S.G.; methodology, M.P., S.G. and S.K.; software, S.K. and M.P.; validation, S.K. and M.P.; formal analysis, M.P., S.K. and S.G.; writing—original draft preparation, M.P.; writing—review and editing, S.G., S.K. and M.P.; supervision, S.G. and S.K. All authors have read and agreed to the published version of the manuscript.

Funding: This research received no external funding.

Data Availability Statement: This study is based on a simulation and does not involve a real dataset. The R code used to run the simulations will be provided upon request.

Conflicts of Interest: The authors declare no conflicts of interest.

References

- Latkin, C.; Edwards, C.; Davey-Rothwell, M.; Tobin, K. The relationship between social desirability bias and self-reports of health, substance use, and social network factors among urban substance users in Baltimore, Maryland. *Addict. Behav.* **2017**, *73*, 133–136. [\[CrossRef\]](#) [\[PubMed\]](#)
- Warner, S. Randomized response: A survey technique for eliminating evasive answer bias. *J. Am. Stat. Assoc.* **1965**, *60*, 63–69. [\[CrossRef\]](#) [\[PubMed\]](#)
- Greenberg, B.G.; Abul-Ela, A.; Simmons, W.R.; Horvitz, D.G. The unrelated question randomized response model: Theoretical framework. *J. Am. Stat. Assoc.* **1969**, *64*, 520–539. [\[CrossRef\]](#)
- Warner, S. The Linear Randomized Response Model. *J. Am. Stat. Assoc.* **1971**, *66*, 884–888. [\[CrossRef\]](#)
- Greenberg, B.G.; Kuebler, R.; Abernathy, J.R.; Horvitz, D.G. Application of the Randomized Response Technique in Obtaining Quantitative Data. *J. Am. Stat. Assoc.* **1971**, *66*, 243–250. [\[CrossRef\]](#)
- Pollock, K.H.; Bek, Y. A comparison of three randomized response models for quantitative data. *J. Am. Stat. Assoc.* **1976**, *71*, 884–886. [\[CrossRef\]](#)
- Eichhorn, B.H.; Hayre, L.S. Scrambled randomized response methods for obtaining sensitive quantitative data. *J. Stat. Plan. Inference* **1983**, *7*, 307–316. [\[CrossRef\]](#)

8. Diana, G.; Perri, P.F. A class of estimators for quantitative sensitive data. *Stat. Pap.* **2011**, *52*, 633–650. [[CrossRef](#)]
9. Singh, G.N.; Kumar, A.; Vishwakarma, G.K. Some alternative additive randomized response models for estimation of population mean of quantitative sensitive variable in the presence of scramble variable. *Commun. Stat.-Simul. Comput.* **2018**, *49*, 2785–2807. [[CrossRef](#)]
10. Gupta, S.; Gupta, B.; Singh, S. Estimation of the sensitivity level of personal interview survey questions. *J. Stat. Plan. Inference* **2002**, *100*, 239–247. [[CrossRef](#)]
11. Mehta, S.; Aggarwal, P. Bayesian estimation of sensitivity level and population proportion of a sensitive characteristic in a binary optional unrelated question RRT model. *Commun. Stat.-Theory Methods* **2018**, *47*, 4021–4028. [[CrossRef](#)]
12. Sharma, P.; Singh, R. Method of Estimation in the Presence of Nonresponse and Measurement Errors Simultaneously. *J. Mod. Appl. Stat. Methods* **2015**, *14*, 107–121.
13. Perri, P.F. Modified Randomized Devices for Simmons' model. *Model Assist. Stat. Appl.* **2008**, *3*, 233–239. [[CrossRef](#)]
14. Blair, G.; Imai, K.; Zhou, Y.Y. Design and Analysis of the Randomized Response Technique. *J. Am. Stat. Assoc.* **2015**, *110*, 1304–1319. [[CrossRef](#)]
15. Parker, M.; Gupta, S.; Khalil, S. A Mixture Quantitative Randomized Response Model That Improves Trust in RRT Methodology. *Axioms* **2024**, *13*, 11. [[CrossRef](#)]
16. Gupta, S.; Mehta, S.; Shabbir, J.; Khalil, S. A Unified Measure of respondent privacy and model efficiency in quantitative RRT models. *J. Stat. Theory Pract.* **2018**, *12*, 506–511. [[CrossRef](#)]
17. Gupta, S.; Zhang, J.; Khalil, S.; Sapra, P. Mitigating Lack of Trust in Quantitative Randomized Response Techniques Models. *Commun. Stat. Part B Simul. Comput.* **2022**, *53*, 2624–2632. [[CrossRef](#)]
18. Thompson, S.K. *Sampling*, 3rd ed.; John Wiley & Sons Inc.: Hoboken, NJ, USA, 2012; pp. 93–124, ISBN 978-0-470-40231-3.
19. Yan, Z.; Wang, J.; Lai, J. An Efficiency and Protection Degree-Based Comparison Among the Quantitative Randomized Response Strategies. *Theory Methods* **2008**, *38*, 400–408. [[CrossRef](#)]
20. Lanke, J. On the degree of protection in randomized interviews. *Int. Stat. Rev.* **1976**, *44*, 33. [[CrossRef](#)]

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.