Scientific Research Publishing

# Estimating a Finite Population Mean under Random Non-Response in Two Stage Cluster Sampling with Replacement

## Nelson Kiprono Bii[1], Christopher Ouma Onyango[2], John Odhiambo[1]

[1]Institute of Mathematical Sciences, Strathmore University, Nairobi, Kenya
[2]Department of Statistics, Kenyatta University, Nairobi, Kenya
Email: nkiprono@strathmore.edu, Chrisouma2004@yahoo.co.uk, jodhiambo@strathmore.edu

## Abstract

Non-response is a regular occurrence in Sample Surveys. Developing estimators when non-response exists may result in large biases when estimating population parameters. In this paper, a finite population mean is estimated when non-response exists randomly under two stage cluster sampling with replacement. It is assumed that non-response arises in the survey variable in the second stage of cluster sampling. Weighting method of compensating for non-response is applied. Asymptotic properties of the proposed estimator of the population mean are derived. Under mild assumptions, the estimator is shown to be asymptotically consistent.

## Keywords

Non-Response, Nadaraya-Watson Estimation, Two Stage Cluster Sampling

## 1. Introduction

In survey sampling, non-response is one source of errors in data analysis. Non-response introduces bias into the estimation of population characteristics. It also causes samples to fail to follow the distributions determined by the original sampling design. This paper seeks to reduce the non-response bias in the estimation of a finite population mean in two stage cluster sampling.

Use of regression models is recognized as one of the procedures for reducing bias due to non-response using auxiliary information. In practice, information on the variables of interest is not available for non-respondents but information on auxiliary variables may be available for non-respondents. It is therefore desirable to model the response behavior and incorporate the auxiliary data into

the estimation so that the bias arising from non-response can be reduced. If the auxiliary variables are correlated with the response behavior, then the regression estimators would be more precise in estimation of population parameters, given the auxiliary information is known.

Many authors have developed estimators of population mean where non-response exists in the study and auxiliary variables. But there exist cases that do not exhibit non-response in the auxiliary variables, such as: number of people in a family, duration one takes to go through education. Imputation techniques have been used to account for non-response in the study variable. For instance, [1] applied compromised method of imputation to estimate a finite population mean under two stage cluster sampling, this method however produced a large bias. In this study, the Nadaraya-Watson regression technique is applied in deriving the estimator for the finite population mean. Kernel weights are used to compensate for non-response.

### Reweighting Method

Non-response causes loss of observations and therefore reweighting means that the weights are increased for all or almost all of the elements that fail to respond in a survey. The population mean, $\bar{Y}$, is estimated by selecting a sample of size $n$ at random with replacement. If responding units to item $y$ are independent so that the probability of unit $j$ responding in cluster $i$ is $p_{ij}$ $(i = 1, 2, \cdots, n; j = 1, 2, \cdots, m)$ then an imputed estimator, $\bar{y}_I$, for $\bar{Y}$, is given by

$$\bar{y}_I = \frac{1}{\sum_{i,j \in s} w_{ij}} \left[ \sum_{i,j \in s_r} w_{ij} y_{ij} + \sum_{i,j \in s_m} w_{ij} y_{ij}^* \right] \tag{1.0}$$

where $w_{ij} = \dfrac{1}{\pi_{ij}}$ gives sample survey weight tied to unit $j$ in cluster $i$ and

$\pi_{ij} = p[i, j \in s]$ is its second order probability of inclusion, $s_r$, is the set of $r$ units responding to item $y$ and $s_m$ is the set of $m$ units that failed to respond to item $y$ so that $r + m = n$ and $y_{ij}^*$ is the imputed value generated so that the missing value $y_{ij}$ is compensated for, [2].

## 2. The Proposed Estimator of Finite Population Mean

Consider a finite population of size $M$ consisting of $N$ clusters with $N_i$ elements in the $i^{\text{th}}$ cluster. A sample of $n$ clusters is selected so that $n_1$ units respond and $n_2$ units fail to respond. Let $y_{ij}$ denote the value of the survey variable $Y$ for unit $j$ in cluster $i$, for $i = 1, 2, \cdots, N$, $j = 1, 2, \cdots, N_i$ and let population mean be given by

$$\bar{\bar{Y}} = \frac{1}{MN_i} \sum_{i=1}^{N} \sum_{j=1}^{M_i} Y_{ij} \tag{2.1}$$

Let an estimator of the finite population mean be defined by $\hat{\bar{\bar{Y}}}$ as follows:

$$\hat{\bar{\bar{Y}}} = \frac{1}{M} \left\{ \frac{1}{n_1} \sum_{i \in s} \sum_{j \in s} \frac{Y_{ij}}{\pi_{ij}} \delta_{ij} + \frac{1}{n_2} \sum_{i \in s} \sum_{j \notin s} \left( 1 - \frac{1}{\pi_{ij}} \right) \hat{Y}_{ij} \delta_{ij} \right\} \tag{2.2}$$

where $\delta_{ij}$ is an indicator variable defined by

$$\delta_{ij} = \begin{cases} 1, & \text{if } j^{\text{th}} \text{ unit in the } i^{\text{th}} \text{ cluster responds} \\ 0, & \text{elsewhere} \end{cases}$$

and $n_1$ and $n_2$ are the number of units that respond and those that fail to respond respectively.

$\pi_{ij}$ is the probability of selecting the $j^{\text{th}}$ unit in the $i^{\text{th}}$ cluster into the sample.

Let $w(x_{ij}) = \dfrac{1}{\pi_{ij}}$ to be the inverse of the second order inclusion probabilities

and that $x_{ij}$ is the $i^{\text{th}}$ auxiliary random variable from the $j^{\text{th}}$ cluster. It follows that Equation (2.2) becomes

$$\hat{\bar{\bar{Y}}} = \frac{1}{M}\left\{ \frac{1}{n_1}\sum_{i \in s}\sum_{j \in s} w(x_{ij})Y_{ij}\delta_{ij} + \frac{1}{n_2}\sum_{i \in s}\sum_{j \notin s}\left(1 - w(x_{ij})\right)\hat{Y}_{ij}\delta_{ij} \right\} \tag{2.3}$$

Suppose $\delta_{ij}$ is known to be Bernoulli random variables with probability of success $\delta_{ij}^*$, then, $E(\delta_{ij}) = p_r(\delta_{ij} = 1) = \delta_{ij}^*$ and $(\delta_{ij}) = \delta_{ij}^*(1 - \delta_{ij}^*)$, [3]. Thus, the expected value of the estimator of population mean is given by

$$E\left(\hat{\bar{\bar{Y}}}\right) = \frac{1}{M}\left\{ \frac{1}{n_1}\sum_{i \in s}\sum_{j \in s} E\left(w(x_{ij})Y_{ij}\right)\delta_{ij} + \frac{1}{n_2}\sum_{i \in s}\sum_{j \notin s} E\left(\left(1 - w(x_{ij})\right)\hat{Y}_{ij}\right)\delta_{ij}^* \right\} \tag{2.4}$$

Assuming non-response in the second stage of sampling, the problem is therefore to estimate the values of $\hat{Y}_{ij}$. To do this, a linear regression model applied by [4] and [5] given below is used;

$$\hat{Y}_{ij} = m(\hat{x}_{ij}) + \hat{e}_{ij} \tag{2.5}$$

where $m(.)$ is a smooth function of the auxiliary variables and $\hat{e}_{ij}$ is the residual term with mean zero and variance which is strictly positive, Substituting Equation (2.5) in Equation (2.4) the following result is obtained:

$$\begin{aligned} E\left(\hat{\bar{\bar{Y}}}\right) = \frac{1}{M}\Bigg\{ &\frac{1}{n_1}\sum_{i \in s}\sum_{j \in s} E\left(\left(m(\hat{x}_{ij}) + \hat{e}_{ij}\right)w(x_{ij})\right)\delta_{ij} \\ &+ \frac{1}{n_2}\sum_{i \in s}\sum_{j \notin s} E\left(1 - w(x_{ij})\right)\left(m(\hat{x}_{ij}) + \hat{e}_{ij}\right)\delta_{ij}^* \Bigg\} \end{aligned} \tag{2.6}$$

Assuming that $n_1 = n_2 = n$, and simplifying Equation (2.6) we obtain the following

$$\begin{aligned} E\left(\hat{\bar{\bar{Y}}}\right) = \frac{1}{Mn}\Bigg\{ &\sum_{i \in s}\sum_{j \in s} E\left(\left(m(\hat{x}_{ij}) + \hat{e}_{ij}\right)w(x_{ij})\right)\delta_{ij} \\ &+ \sum_{i \in s}\sum_{j \notin s} E\left(1 - w(x_{ij})\right)\left(m(\hat{x}_{ij}) + \hat{e}_{ij}\right)\delta_{ij}^* \Bigg\} \end{aligned} \tag{2.7}$$

A detailed work done by [5] proved that $E(\hat{e}_{ij}) = 0$. Therefore Equation (2.7) reduces to

$$\begin{aligned} E\left(\hat{\bar{\bar{Y}}}\right) = \frac{1}{Mn}\Bigg\{ &\sum_{i \in s}\sum_{j \in s} E\left(m(\hat{x}_{ij})\right)E\left(w(x_{ij})\right)\delta_{ij} \\ &+ \sum_{i \in s}\sum_{j \notin s} E\left(1 - w(x_{ij})\right)E\left(m(\hat{x}_{ij}) + \hat{e}_{ij}\right)\delta_{ij}^* \Bigg\} \end{aligned} \tag{2.8}$$

The second term in Equation (2.8) is simplified as follows:

$$\frac{1}{Mn}\left\{\sum_{i\notin s}\sum_{j\notin s}E\left(1-w\left(x_{ij}\right)\right)E\left(m\left(\hat{x}_{ij}\right)+\hat{e}_{ij}\right)\delta_{ij}^{*}\right\}$$

$$=\frac{1}{Mn}\left\{\sum_{i\notin s}\sum_{j\notin s}E\left(1-w\left(x_{ij}\right)\right)m\left(\hat{x}_{ij}\right)\delta_{ij}\right\} \tag{2.9}$$

$$+\frac{1}{Mn}\left\{\sum_{i\notin s}\sum_{j\notin s}E\left(1-w\left(x_{ij}\right)\right)e_{ij}\delta_{ij}\right\}$$

But $E\left(m\left(x_{ij}\right)\right)=m\left(\hat{x}_{ij}\right)=m\left(x_{ij}\right)$, [6]. Thus we get the following:

$$\frac{1}{Mn}\left\{\sum_{i\notin s}\sum_{j\notin s}E\left(1-w\left(x_{ij}\right)\right)E\left(m\left(\hat{x}_{ij}\right)+\hat{e}_{ij}\right)\delta_{ij}^{*}\right\}$$

$$=\frac{1}{Mn}\left\{\sum_{i=m+1}^{M}\sum_{j=n+1}^{N}\delta_{ij}m\left(x_{ij}\right)-w\left(x_{ij}\right)\delta_{ij}m\left(x_{ij}\right)\right\} \tag{2.10}$$

$$+\frac{1}{Mn}\left\{\sum_{i=m+1}^{M}\sum_{j=n+1}^{N}E\left(e_{ij}\delta_{ij}\right)-E\left(w\left(x_{ij}\right)\left(e_{ij}\delta_{ij}\right)\right)\right\}$$

$$\frac{1}{Mn}\left\{\sum_{i\notin s}\sum_{j\notin s}E\left(1-w\left(x_{ij}\right)\right)E\left(m\left(\hat{x}_{ij}\right)+\hat{e}_{ij}\right)\delta_{ij}^{*}\right\}$$

$$=\frac{1}{Mn}\left\{\left(M-(m+1)\right)\left(N-(n+1)\right)\left[\left(\delta_{ij}\right)m\left(x_{ij}\right)-w\left(x_{ij}\right)\delta_{ij}m\left(x_{ij}\right)\right]\right\} \tag{2.11}$$

$$+\frac{1}{Mn}\left\{\left(M-(m+1)\right)\left(N-(n+1)\right)\left[\delta_{ij}E\left(e_{ij}\right)-E\left(e_{ij}\right)\delta_{ij}w\left(x_{ij}\right)\right]\right\}$$

But $E\left(e_{ij}\right)=0$, for details see [5].

On simplification, Equation (2.11) reduces to

$$\frac{1}{Mn}\left\{\sum_{i\notin s}\sum_{j\notin s}E\left(1-w\left(x_{ij}\right)\right)E\left(m\left(\hat{x}_{ij}\right)+\hat{e}_{ij}\right)\delta_{ij}^{*}\right\}$$

$$=\frac{\left(M-(m+1)\right)\left(N-(n+1)\right)}{Mn}\left\{\delta_{ij}m\left(x_{ij}\right)\left(1-w\left(x_{ij}\right)\right)\right\} \tag{2.12}$$

Recall $w\left(x_{ij}\right)=\dfrac{1}{\pi_{ij}}$

so that Equation (2.12) may be re-written as follows:

$$\frac{1}{Mn}\left\{\sum_{i\notin s}\sum_{j\notin s}E\left(1-w\left(x_{ij}\right)\right)E\left(m\left(\hat{x}_{ij}\right)+\hat{e}_{ij}\right)\delta_{ij}^{*}\right\}$$

$$=\frac{\left(M-(m+1)\right)\left(N-(n+1)\right)}{Mn}\left\{\delta_{ij}m\left(x_{ij}\right)\left(\frac{\pi_{ij}-1}{\pi_{ij}}\right)\right\} \tag{2.13}$$

Assume the sample sizes are large *i.e.* as $n \to N$ and $m \to M$, Equation (2.13) simplifies to

$$\frac{1}{Mn}\left\{\sum_{i\notin s}\sum_{j\notin s}E\left(1-w\left(x_{ij}\right)\right)E\left(m\left(\hat{x}_{ij}\right)+\hat{e}_{ij}\right)\delta_{ij}^{*}\right\}$$

$$=\frac{1}{Mn}\left\{\delta_{ij}m\left(x_{ij}\right)\left(\frac{\pi_{ij}-1}{\pi_{ij}}\right)\right\} \tag{2.14}$$

Combining Equation (2.14) with the first term in Equation (2.08) becomes;

$$E\left(\hat{\bar{\bar{Y}}}\right) = \frac{1}{Mn}\left\{\sum_{i\in s}\sum_{j\in s}E\left(m\left(x_{ij}\right)\right)E\left(\frac{\delta_{ij}}{\pi_{ij}}\right) + \sum_{i\in s}\sum_{j\notin s}\delta_{ij}\left(m\left(\hat{x}_{ij}\right)\right)\left(\frac{\pi_{ij}-1}{\pi_{ij}}\right)\right\} \quad (2.15)$$

Since the first term represents the response units, their values are all known. The problem is to estimate the non-response units in the second term. Let the indicator variable $\delta_{ij}=1$, the problem now reduces to that of estimating the function $m\left(\hat{x}_{ij}\right)$, which is a function of the auxiliary variables, $x_{ij}$. Hence the expected value of the estimator of the finite population mean under non-response is given as;

$$E\left(\hat{\bar{\bar{Y}}}\right) = \frac{1}{Mn}\left\{\sum_{i\in s}\sum_{j\in s}Y_{ij} + \sum_{i\in s}\sum_{j\notin s}\delta_{ij}\left(m\left(\hat{x}_{ij}\right)\right)\left(\frac{\pi_{ij}-1}{\pi_{ij}}\right)\right\} \quad (2.16)$$

In order to derive the asymptotic properties of the expected value of the proposed estimator in 2.16, first a review of Nadaraya-Watson estimator is given below.

## 3. Review of Nadaraya-Watson Estimator

Given a random sample of bivariate data $\left(x_i, y_i\right), \cdots, \left(x_n, y_n\right)$ having a joint pdf $g\left(x, y\right)$ with the regression model given by

$Y_{ij} = m\left(x_{ij}\right) + e_{ij}$ as in Equation (2.5), where $m\left(.\right)$ is unknown. Let the error term satisfy the following conditions:

$$E\left(e_{ij}\right)=0, \ Var\left(e_{ij}\right)=\sigma_{ij}^2, \ cov\left(e_i, e_j\right)=0 \ \text{for } i \neq j \quad (3.0)$$

Furthermore, let $K\left(.\right)$ denote a symmetric kernel density function which is twice continuously differentiable with:

$$\left.\begin{array}{c}\int_{-\infty}^{\infty}k\left(w\right)\mathrm{d}w=1 \\ \int_{-\infty}^{\infty}wk\left(w\right)\mathrm{d}w=0 \\ \int_{-\infty}^{\infty}k^2\left(w\right)\mathrm{d}w<\infty \\ \int_{-\infty}^{\infty}w^2k\left(w\right)\mathrm{d}w=d_k \\ k\left(w\right)=k\left(-w\right)\end{array}\right\} \quad (3.1)$$

In addition, let the smoothing weights be defined by

$$w\left(x_{ij}\right) = \frac{K\left(\dfrac{x-X_{ij}}{b}\right)}{\sum_{i\in s}\sum_{i\in s}K\left(\dfrac{x-X_{ij}}{b}\right)}, \ i=1,2,\cdots,n; j=1,2,\cdots,m \quad (3.2)$$

where $b$ is a smoothing parameter, normally referred to as the bandwidth such that, $\sum_i\sum_j w\left(x_{ij}\right)=1$.

Using Equation (3.2), the Nadaraya-Watson estimator of $m\left(x_{ij}\right)$ is given by:

$$m\left(\hat{x}_{ij}\right)=\sum_{i\in s}\sum_{j\in s}w\left(x_{ij}\right)Y_{ij}=\frac{\sum_{i\in s}\sum_{j\in s}K\left(\dfrac{x-X_{ij}}{b}\right)Y_{ij}}{\sum_{i\in s}\sum_{j\in s}K\left(\dfrac{x-X_{ij}}{b}\right)},\ i=1,2,\cdots,n;\ j=1,2,\cdots,m \quad (3.3)$$

Given the model $\hat{Y}_{ij}=m\left(\hat{x}_{ij}\right)+\hat{e}_{ij}$ and the conditions of the error term as explained in 3.0 above, the expression for the survey variable $Y_{ij}$ relative to the auxiliary variable $X_{ij}$ can be given as a joint pdf of $g\left(x_{ij},y_{ij}\right)$ as follows:

$$m\left(x_{ij}\right)=E\left(Y_{ij}/X_{ij}=x_{ij}\right)=\int yg\left[y/x\right]\mathrm{d}y=\frac{\int yg\left(x,y\right)\mathrm{d}y}{\int g\left(x,y\right)\mathrm{d}y} \quad (3.4)$$

where $\int g\left(x,y\right)\mathrm{d}y$ is the marginal density of $X_{ij}$. The numerator and the denominator of Equation (3.4) can be estimated separately using kernel functions as follows:

$g\left(x,y\right)$ is estimated by;

$$\hat{g}\left(x,y\right)=\frac{1}{mn}\sum_{i}\sum_{j}\left(\frac{1}{b}K\left(\frac{x-X_{ij}}{b}\right)\frac{1}{b}K\left(\frac{y-Y_{ij}}{b}\right)\right) \quad (3.5)$$

and

$$\int y\hat{g}\left(x,y\right)\mathrm{d}y=\frac{1}{mn}\sum_{i}\sum_{j}\int\left(\frac{1}{b}K\left(\frac{x-X_{ij}}{b}\right)\frac{1}{b}K\left(\frac{y-Y_{ij}}{b}\right)\right)y\mathrm{d}y \quad (3.6)$$

Using change of variables technique; let

$$\left.\begin{aligned}w&=\frac{y-Y_{ij}}{b}\\y&=wb+Y_{ij}\\\mathrm{d}y&=b\mathrm{d}w\end{aligned}\right\} \quad (3.7)$$

So that

$$\int y\hat{g}\left(x,y\right)\mathrm{d}y=\frac{1}{mn}\sum_{i}\sum_{j}\int\frac{1}{b}K\left(\frac{x-X_{ij}}{b}\right)\frac{1}{b}\left(bw+Y_{ij}\right)K\left(w\right)b\mathrm{d}w \quad (3.8)$$

$$=\frac{1}{mnb}\sum_{i}\sum_{j}K\left(\frac{x-X_{ij}}{b}\right)\left[\int wK\left(w\right)b\mathrm{d}w+\frac{1}{b}Y_{ij}\int K\left(w\right)b\mathrm{d}w\right] \quad (3.9)$$

From the conditions specified in Equation (3.1), the following (3.9) simplifies to

$$\int y\hat{g}\left(x,y\right)\mathrm{d}y=\frac{1}{mnb}\sum_{i}\sum_{j}K\left(\frac{x-X_{ij}}{b}\right)\left[0+Y_{ij}\right] \quad (3.10)$$

which reduces to:

$$\int y\hat{g}\left(x,y\right)\mathrm{d}y=\frac{1}{mnb}\sum_{i}\sum_{j}K\left(\frac{x-X_{ij}}{b}\right)Y_{ij} \quad (3.11)$$

Following the same procedure, the denominator can be obtained as follows:

$$\int \hat{g}(x, y) \mathrm{d}y = \frac{1}{mn} \sum_i \sum_j \int \left( \frac{1}{b} K\left( \frac{x - X_{ij}}{b} \right) \frac{1}{b} K\left( \frac{y - Y_{ij}}{b} \right) \right) \mathrm{d}y$$

$$= \frac{1}{mnb} \sum_{i=1}^{n} \sum_{j=1}^{m} K\left( \frac{x - X_{ij}}{b} \right) \int \frac{1}{b} K\left( \frac{y - Y_{ij}}{b} \right) \mathrm{d}y \qquad (3.12)$$

Using change of variable technique as in Equation (3.7), Equation (3.12) can be re-written as follows:

$$\int \hat{g}(x, y) \mathrm{d}y = \frac{1}{mnb} \sum_{i=1}^{n} \sum_{j=1}^{m} K\left( \frac{x - X_{ij}}{b} \right) \int \frac{1}{b} K(w) b \mathrm{d}w \qquad (3.13)$$

which yields

$$\int \hat{g}(x, y) \mathrm{d}y = \frac{1}{mnb} \sum_{i=1}^{n} \sum_{j=1}^{m} K\left( \frac{x - X_{ij}}{b} \right) \qquad (3.14)$$

Since $\int \frac{1}{b} K(w) b \mathrm{d}w$ is a pdf and therefore integrates to 1.

It follows from Equations ((3.11) and (3.14)) that the estimator $m(\hat{x}_{ij})$ is as given in Equation (3.3). Thus the estimator of $m(x_{ij})$ is a linear smoother since it is a linear function of the observations, $Y_{ij}$. Given a sample and a specified kernel function, then for a given auxiliary value $x_{ij}$, the corresponding y-estimate is obtained by the estimator outlined in Equation (3.3), which can be written as:

$$\hat{y}_{ij} = m_{NW}\left( \hat{x}_{ij} \right) = \sum_i \sum_j W_{ij}\left( x_{ij} \right) Y_{ij} \qquad (3.15)$$

where $m_{NW}\left( \hat{x}_{ij} \right)$ is the Nadaraya-Watson estimator for estimating the unknown function $m(.)$, for details see [7] [8].

This provides a way of estimating for instance the non-response values of the survey variable $Y_{ij}$, given the auxiliary values $x_{ij}$, for a specified kernel function.

## 4. Asymptotic Bias of the Mean Estimator $\hat{\bar{\bar{Y}}}$

Equation (2.16) may be written as

$$E\left( \hat{\bar{\bar{Y}}} \right) = \frac{1}{MN} \left\{ \sum_{i=1}^{n} \sum_{j=1}^{m} Y_{ij} + \sum_{i=n+1}^{N} \sum_{j=m+1}^{M} m_{NW}\left( \hat{y}_{ij} \right) \right\} \qquad (4.1)$$

Replacing $x$ by $x_{ij}$ and re-writing Equation (3.15) using the property of symmetry associated with Nadaraya-Watson estimator, then

$$m_{NW}\left( \hat{x}_{ij} \right) = \frac{\sum_{i \in s} \sum_{j \in s} K\left( \frac{X_{ij} - x_{ij}}{b} \right) Y_{ij}}{\sum_{i \in s} \sum_{j \in s} K\left( \frac{X_{ij} - x_{ij}}{b} \right)}, i = 1, 2, \cdots, n; j = 1, 2, \cdots, m \qquad (4.2)$$

$$= \frac{1}{\hat{g}(x_{ij})} \left[ \frac{1}{mnb} \sum_i \sum_j K\left( \frac{X_{ij} - x_{ij}}{b} \right) Y_{ij} \right] \qquad (4.3)$$

where $\hat{g}(x_{ij})$ is the estimated marginal density of auxiliary variables $X_{ij}$.

But for a finite population mean, the expected value of the estimator is given in Equation (4.1). The bias is given by

$$\text{Bias}\left(\hat{\bar{\bar{Y}}}\right) = E\left(\hat{\bar{\bar{Y}}} - \bar{\bar{Y}}\right) \tag{4.4}$$

$$\text{Bias}\left(\hat{\bar{\bar{Y}}}\right) = E\left\{\frac{1}{MN}\left[\sum_{i=1}^{n}\sum_{j=1}^{m}Y_{ij} + \sum_{i=n+1}^{N}\sum_{j=m+1}^{M}m(\hat{x}_{ij})\right] \right.$$
$$\left. - \frac{1}{MN}\left[\sum_{i=1}^{n}\sum_{j=1}^{m}Y_{ij} + \sum_{i=n+1}^{N}\sum_{j=m+1}^{M}Y_{ij}\right]\right\} \tag{4.5}$$

Which reduces to

$$\text{Bias}\left(\hat{\bar{\bar{Y}}}\right) = \frac{1}{MN}\left\{\sum_{i=n+1}^{N}\sum_{j=m+1}^{M}m(\hat{x}_{ij}) - \sum_{i=n+1}^{N}\sum_{j=m+1}^{M}Y_{ij}\right\} \tag{4.6}$$

$$= \frac{1}{MN}\left\{\sum_{i=n+1}^{N}\sum_{j=m+1}^{M}m(\hat{x}_{ij}) - \sum_{i=n+1}^{N}\sum_{j=m+1}^{M}m(x_{ij})\right\} \tag{4.7}$$

Re-writing the regression model given by $Y_{ij} = m(X_{ij}) + e_{ij}$ as

$$Y_{ij} = m(x_{ij}) + \left[m(X_{ij}) - m(x_{ij})\right] + e_{ij} \tag{4.8}$$

So that from Equation (4.3) the first term in Equation (4.7) before taking the expectation is given as:

$$\frac{1}{MN}\left\{\frac{\frac{1}{mnb}\sum_{i=n+1}^{N}\sum_{j=m+1}^{M}K\left(\frac{X_{ij}-x_{ij}}{b}\right)Y_{ij}}{\hat{g}(x_{ij})}\right\}$$

$$= \frac{1}{MN}\left\{\frac{1}{\hat{g}(x_{ij})}\left\{\sum_{i=n+1}^{N}\sum_{j=m+1}^{M}K\left(\frac{X_{ij}-x_{ij}}{b}\right)m(x_{ij})\right.\right.$$
$$+ \frac{1}{mnb}\sum_{i=n+1}^{N}\sum_{j=m+1}^{M}K\left(\frac{X_{ij}-x_{ij}}{b}\right)\left[m(X_{ij}) - m(x_{ij})\right]$$
$$\left.\left. + \frac{1}{mnb}\sum_{i=n+1}^{N}\sum_{j=m+1}^{M}K\left(\frac{X_{ij}-x_{ij}}{b}\right)e_{ij}\right\}\right\} \tag{4.9}$$

Simplifying Equation (4.9) the following is thus obtained:

$$\frac{1}{MN}\left\{\frac{1}{mnb\hat{g}(x_{ij})}\sum_{i=n+1}^{N}\sum_{j=m+1}^{M}K\left(\frac{X_{ij}-x_{ij}}{b}\right)Y_{ij}\right\}$$

$$= \frac{1}{MN}\left\{\frac{\sum_{i=n+1}^{N}\sum_{j=m+1}^{M}\left[\hat{g}(x_{ij})m(x_{ij}) + \hat{m}_1(x_{ij}) + \hat{m}_2(x_{ij})\right]}{mnb\hat{g}(x_{ij})}\right\} \tag{4.10}$$

where

$$\hat{m}_1(x_{ij}) = \sum_{i=n+1}^{N}\sum_{j=m+1}^{M}K\left(\frac{X_{ij}-x_{ij}}{b}\right)\left[m(X_{ij}) - m(x_{ij})\right]$$

$$\hat{m}_2\left(x_{ij}\right) = \sum_{i=n+1}^{N} \sum_{j=m+1}^{M} K\left(\frac{X_{ij} - x_{ij}}{b}\right) e_{ij}$$

Taking conditional expectation of Equation (4.10) we get

$$E\left[\frac{\sum_{i=n+1}^{N}\sum_{j=m+1}^{M} M\left(\hat{x}_{ij}\right)}{x_{ij}}\right]$$

$$= \frac{1}{MN} E\left[\frac{1}{mnb} \sum_{i=n+1}^{N} \sum_{j=m+1}^{M}\left[m\left(x_{ij}\right) + \frac{\hat{m}_1\left(x_{ij}\right)}{\hat{g}\left(x_{ij}\right)} + \frac{\hat{m}_2\left(x_{ij}\right)}{\hat{g}\left(x_{ij}\right)}\right]\right] \tag{4.11}$$

To obtain the relationship between the conditional mean and the selected bandwidth, the following theorem due to [6] is applied;

### Theorem: (Dorfman, 1992)

Let $k(w)$ be a symmetric density function with $\int wk(w)\mathrm{d}w = 0$ and $\int w^2 k(w)\mathrm{d}w = k_2$. Assume $n$ and $N$ increase together such that $\frac{n}{N} \to \pi$ with $0 < \pi < 1$. Besides, assume the sampled and non-sampled values of $x$ are in the interval $[c,d]$ and are generated by densities $d_s$ and $d_{p-s}$ respectively both bounded away from zero on $[c,d]$ and assumed to have continuous second derivatives. If for any variable $\mathcal{Z}$, $E(\mathcal{Z}/U = u) = A(u) + O(B)$ and $Var(\mathcal{Z}/U = u) = O(C)$, then $\mathcal{Z} = A(u) + O_p\left(B + C^{\frac{1}{2}}\right)$.

Applying this theorem, we have

$$MSE\left(\frac{\hat{\bar{Y}}}{x_{ij}}\right) = \frac{1}{\left(MN\right)^2}\left\{ \frac{\left(MN - mn\right)^2 \int k\left(w^2\right)\mathrm{d}w}{mnbg\left(x_{ij}\right)} \right.$$

$$+ \frac{\left(MN - mn\right)^2}{4m^2 n^2} b^4 k_2^2(k)\left[m''\left(x_{ij}\right) + \frac{2g''\left(x_{ij}\right)m'\left(x_{ij}\right)}{g\left(x_{ij}\right)}\right]^2 \tag{4.12}$$

$$\left. + O\left(b^4\right) + O\left[\frac{\left(MN - mn\right)^2}{mnb} + \frac{1}{mnb}\right]\right\}$$

This theorem is stated without proof. To prove it, we partition it into the bias and variance terms and separately prove them as follows:

From Equation (3.0) it follows that $E\left(e_{ij}/X_{ij}\right) = 0$. Therefore, $E\left[\hat{m}_2\left(x_{ij}\right)\right] = 0$. Thus $E\left[\hat{m}_1\left(x_{ij}\right)\right]$ can be obtained as follows:

$$E \sum_{i=n+1}^{N} \sum_{j=m+1}^{M} \left[\hat{m}_1\left(x_{ij}\right)\right]$$

$$= \frac{1}{MN}\left\{\frac{1}{mnb} E\left\{\sum_{i=n+1}^{N} \sum_{j=m+1}^{M} K\left(\frac{X_{ij} - x_{ij}}{b}\right)\left[m\left(X_{ij}\right) - m\left(x_{ij}\right)\right]\right\}\right\} \tag{4.13}$$

Using substitution and change of variable technique below

$$w = \frac{V - x_{ij}}{b} \text{ so that } V = x_{ij} + bw \text{ and } \mathrm{d}V = b\mathrm{d}w \tag{4.14}$$

Equation (4.13) can simplify to:

$$E \sum_{i=n+1}^{N} \sum_{j=m+1}^{M} \left[ \hat{m}_1 \left( x_{ij} \right) \right]$$

$$= \frac{1}{MN} \left\{ \frac{MN-mn}{mnb} \int k(w) \left[ m \left( x_{ij} + bw \right) - m \left( x_{ij} \right) \right] \int g \left( x_{ij} + bw \right) b \mathrm{d}w \right\} \tag{4.15}$$

$$= \frac{1}{MN} \left\{ \frac{MN-mn}{mn} \int k(w) \left[ m \left( x_{ij} + bw \right) - m \left( x_{ij} \right) \right] g \left( x_{ij} + bw \right) \mathrm{d}w \right\} \tag{4.16}$$

Using the Taylor's series expansion about the point $x_{ij}$, the $k^{\text{th}}$ order kernel can be derived as follows:

$$g \left( x_{ij} + bw \right) = g \left( x_{ij} \right) + g' \left( x_{ij} \right) bw + \frac{1}{2} g'' \left( x_{ij} \right) b^2 w^2 + \cdots + \frac{1}{k!} g^k \left( x_{ij} \right) b^k w^k + O \left( b^2 \right) \tag{4.17}$$

Similarly,

$$m \left( x_{ij} + bw \right) = m \left( x_{ij} \right) + m' \left( x_{ij} \right) bw + \frac{1}{2} m'' \left( x_{ij} \right) b^2 w^2 + \cdots + \frac{1}{k!} m^k \left( x_{ij} \right) b^k w^k + O \left( b^2 \right) \tag{4.18}$$

Expanding up to the 3$^{\text{rd}}$ order kernels, Equation (4.18) becomes

$$\left[ m \left( x_{ij} + bw \right) - m(x_{ij}) \right] = m' \left( x_{ij} \right) bw + \frac{1}{2} m'' \left( x_{ij} \right) b^2 w^2 + \frac{1}{3!} m''' \left( x_{ij} \right) b^3 w^3 \tag{4.19}$$

In a similar manner, the expansion of Equation (4.16) up to order $O \left( b^2 \right)$ is given by:

$$E \sum_{i=n+1}^{N} \sum_{j=m+1}^{M} \left[ \hat{m}_1 \left( x_{ij} \right) \right]$$

$$= \frac{1}{MN} \left\{ \frac{MN-mn}{mn} \int k(w) \left( m' \left( x_{ij} \right) bw + \frac{1}{2} m'' \left( x_{ij} \right) b^2 w^2 \right) \left( g \left( x_{ij} \right) + g' \left( x_{ij} \right) bw \right) \mathrm{d}w \right\} \tag{4.20}$$

Simplifying Equation (4.20) gives;

$$E \sum_{i=n+1}^{N} \sum_{j=m+1}^{M} \left[ \hat{m}_1 \left( x_{ij} \right) \right] = \frac{1}{MN} \left\{ \left( \frac{MN-mn}{mn} \right) g \left( x_{ij} \right) m' \left( x_{ij} \right) b \int w k(w) \mathrm{d}w \right.$$

$$+ \left( \frac{MN-mn}{mn} \right) g' \left( x_{ij} \right) m' \left( x_{ij} \right) b^2 \int w^2 k(w) \mathrm{d}w \tag{4.21}$$

$$\left. + \left( \frac{MN-mn}{mn} \right) \frac{1}{2} g \left( x_{ij} \right) m'' \left( x_{ij} \right) b^2 \int w^2 k(w) \mathrm{d}w + O \left( b^2 \right) \right\}$$

Using the conditions stated in Equation (3.1), the derivation in (4.21) can further be simplified to obtain:

$$E \sum_{i=n+1}^{N} \sum_{j=m+1}^{M} \left[ \hat{m}_1 \left( x_{ij} \right) \right]$$

$$= \frac{1}{MN} \left\{ \left( \frac{MN-mn}{mn} \right) \left[ g' \left( x_{ij} \right) m' \left( x_{ij} \right) + \frac{1}{2} g \left( x_{ij} \right) m'' \left( x_{ij} \right) \right] b^2 d_k + O \left( b^2 \right) \right\} \tag{4.22}$$

Hence the expected value of the second term in Equation (4.11) then becomes:

$$E \sum_{i=n+1}^{N} \sum_{j=m+1}^{M} \left[ \hat{m}_1 \left( x_{ij} \right) \right]$$

$$= \frac{1}{MN} \left\{ \left( \frac{MN-mn}{mn} \right) \left[ \frac{1}{2} m'' \left( x_{ij} \right) + \frac{g' \left( x_{ij} \right) m' \left( x_{ij} \right)}{g \left( x_{ij} \right)} \right] b^2 d_k + O \left( b^2 \right) \right\} \tag{4.23}$$

$$= \frac{1}{MN}\left\{\left(\frac{MN-mn}{mn}\right)\left[\frac{m''(x_{ij})}{2}+\left[g(x_{ij})\right]^{-1}g'(x_{ij})m'(x_{ij})\right]b^2 d_k + O(b^2)\right\} \quad (4.24)$$

$$= \frac{1}{MN}\left\{\left(\frac{MN-mn}{mn}\right)b^2 d_k C(x)+O(b^2)\right\} \quad (4.25)$$

where

$$C(x)=\frac{1}{2}m''(x_{ij})+\left[g(x_{ij})\right]^{-1}g'(x_{ij})m'(x_{ij}) \quad (4.26)$$

and $d_k$ is as stated in Equation (3.1)

Using equation of the bias given in (4.4) and the conditional expectation in Equation (4.11), we obtain the following equation for the bias of the estimator:

$$\begin{aligned}\text{Bias}\left(\hat{\bar{\bar{Y}}}\right) &= \frac{1}{MN}\left\{\left(\frac{MN-mn}{mn}\right)b^2 d_k C(x)+O(b^2)\right\} \\ &= \frac{1}{MN}\left\{\left(\frac{MN-mn}{mn}\right)b^2 d_k C(x)+O(b^2)\right\}\end{aligned} \quad (4.27)$$

## 5. Asymptotic Variance of the Estimator, $\hat{\bar{\bar{Y}}}$

From Equations ((4.9) and (4.11)),

$$m'_2(x_{ij})=\frac{1}{mnb}\sum_{i=1}^{n}\sum_{j=1}^{m}K\left(\frac{X_{ij}-x_{ij}}{b}\right)e_{ij} \quad (5.0)$$

Hence

$$\text{Var}\sum_{i=n+1}^{N}\sum_{j=m+1}^{M}\left[\hat{m}_2(x_{ij})\right]=\frac{1}{(MN)^2}\left(\frac{MN-mn}{mnb}\right)^2\sum_{i=1}^{n}\sum_{j=1}^{m}\text{Var}(D_x) \quad (5.1)$$

where

$$D_x=K\left(\frac{X_{ij}-x_{ij}}{b}\right)e_{ij}$$

Expressing Equation (5.1) in terms of expectation we obtain:

$$\text{Var}\sum_{i=n+1}^{N}\sum_{j=m+1}^{M}\left[\hat{m}_2(x_{ij})\right]=\frac{1}{(MN)^2}\left[\frac{(MN-mn)^2}{mnb^2}\right]\left\{E[D_x]^2-\left[E(D_x)\right]^2\right\} \quad (5.2)$$

Using the fact that the conditional expectation $E(e_{ij}/X_{ij})=0$, the second term in Equation (4.13) reduces to zero. Therefore,

$$\text{Var}\sum_{i=n+1}^{N}\sum_{j=m+1}^{M}\left[\hat{m}_2(x_{ij})\right]=\frac{1}{(MN)^2}\left[\frac{(MN-mn)^2}{mnb^2}\right]\sigma^2_{(x_{ij})} \quad (5.3)$$

where

$$E\left(e_{ij}/X_{ij}\right)^2=\sigma^2_{(x_{ij})}$$

Let $X=X_{ij}$, and $x=x_{ij}$, and making the following substitutions

$$w = \frac{X - x}{b} \\ X - x = bw \\ \mathrm{d}X = b\mathrm{d}w$$ (5.4)

$$\mathrm{Var} \sum_{i=n+1}^{N} \sum_{j=m+1}^{M} \left[ \hat{m}_2 \left( x_{ij} \right) \right] = \frac{(MN - mn)^2}{mnb^2 (MN)^2} \int K \left( \frac{X - x}{b} \right)^2 \sigma_{(x_{ij})}^2 g(X) \mathrm{d}X$$ (5.5)

$$= \frac{(MN - mn)^2}{mnb^2 (MN)^2} \int K(w)^2 \sigma_{(x_{ij})}^2 g(x + bw) b\mathrm{d}w$$ (5.6)

which can be simplified to get:

$$\mathrm{Var} \sum_{i=n+1}^{N} \sum_{j=m+1}^{M} \left[ \hat{m}_2 \left( x_{ij} \right) \right] = \frac{(MN - mn)^2}{mnb (MN)^2} \int K(w)^2 g(x) \sigma_{(x_{ij})}^2 \mathrm{d}w + O \left( \frac{1}{mnb} \right)$$ (5.7)

$$\mathrm{Var} \sum_{i=n+1}^{N} \sum_{j=m+1}^{M} \left[ \hat{m}_1 \left( x_{ij} \right) \right]$$

$$= \frac{1}{(MN)^2} \mathrm{Var} \sum_{i=n+1}^{N} \sum_{j=m+1}^{M} \left[ \frac{1}{mnb} \sum_{i=1}^{n} \sum_{j=1}^{m} K \left( \frac{X_{ij} - x_{ij}}{b} \right) \right] \left[ M \left( X_{ij} \right) - m \left( x_{ij} \right) \right]$$ (5.8)

$$\mathrm{Var} \sum_{i=n+1}^{N} \sum_{j=m+1}^{M} \left[ \hat{m}_1 \left( x_{ij} \right) \right] = \frac{(MN - mn)^2}{mnb^2 (MN)^2} \mathrm{Var} \, K \left( \frac{X_{ij} - x_{ij}}{b} \right) \left[ M \left( X_{ij} \right) - m \left( x_{ij} \right) \right]$$ (5.9)

Hence

$$\mathrm{Var} \sum_{i=n+1}^{N} \sum_{j=m+1}^{M} \left[ \hat{m}_1 \left( x_{ij} \right) \right]$$

$$= \frac{(MN - mn)^2}{mnb^2 (MN)^2} E \left[ \int K \left( \frac{X - x}{b} \right)^2 \left[ M(X) - m(x) \right]^2 \right] g(X) \mathrm{d}X$$ (5.10)

where $X = bw + x$ so that $\mathrm{d}X = b\mathrm{d}w$.

Changing variables and applying Taylor's series expansion then

$$\mathrm{Var} \sum_{i=n+1}^{N} \sum_{j=m+1}^{M} \left[ \hat{m}_1 \left( x_{ij} \right) \right]$$

$$= \frac{(MN - mn)^2}{mnb^2 (MN)^2} \int K(w)^2 \left[ m(x + bw) - m(x) \right]^2 g(x + bw) \mathrm{d}w$$ (5.11)

$$= \frac{(MN - mn)^2}{mnb^2 (MN)^2} \int K(w)^2 \left[ m(x) + m'(x) bw + \cdots - m(x) \right]^2 \left( g(x) + g'(x) bw \right) \mathrm{d}w$$ (5.12)

which simplifies to

$$\mathrm{Var} \sum_{i=n+1}^{N} \sum_{j=m+1}^{M} \left[ \hat{m}_1 \left( x_{ij} \right) \right] = O \left[ \frac{(MN - mn)^2 b^2}{mnb} \right]$$ (5.13)

For large samples, as $n \to N$, $m \to M$ and for $b \to 0$, then $mnb \to \infty$. Hence the variance in Equation (5.12) asymptotically tends to zero, that is,

$$\mathrm{Var} \sum_{i=n+1}^{N} \sum_{j=m+1}^{M} \left[ \hat{m}_1 \left( x_{ij} \right) \right] \to 0$$

$$\text{Var}\left(\hat{\bar{\bar{Y}}}\right) = \frac{(MN-mn)^2}{mnb(MN)^2} \sum_{i=n+1}^{N} \sum_{j=m+1}^{M} \text{Var}\left[m(x_{ij}) + \frac{m_1'(x_{ij}) + m_2''(x_{ij})}{\hat{g}(x_{ij})}\right] \qquad (5.14)$$

On simplification,

$$\text{Var}\left(\hat{\bar{\bar{Y}}}\right) = \frac{(MN-mn)^2}{mnb(MN)^2 \left[\hat{g}(x_{ij})\right]^2} \text{Var}\left\{\sum_{i=n+1}^{N} \sum_{j=m+1}^{M} \left[\hat{m}_2(x_{ij})\right]\right\} \qquad (5.15)$$

Substituting Equations ((5.7) into (5.15)) yields the following:

$$\text{Var}\left(\hat{\bar{\bar{Y}}}\right) = \frac{1}{(MN)^2} \left\{ \frac{(MN-mn)^2 \int K(w)^2 \sigma_{(x_{ij})}^2 dw}{mnb(g(x_{ij}))} + O\left[\frac{(MN-mn)^2}{mnb} + \frac{1}{mnb}\right]\right\} \qquad (5.16)$$

$$= \frac{1}{(MN)^2} \left\{ \frac{(MN-mn)^2 H(w)\sigma_{(x_{ij})}^2}{mnb(g(x_{ij}))} + O\left[\frac{(MN-mn)^2}{mnb} + \frac{1}{mnb}\right]\right\} \qquad (5.17)$$

where, $H(w) = \int K(w)^2 dw$

It is notable that the variance term still depends on the marginal density function, $g(x_{ij})$ of the auxiliary variables $X_{ij}$. It can also be observed that the variance is inversely related to the smoothing parameter $b$. This implies that an increase in $b$ results in a smaller variance. However, increasing the bandwidth would give a larger bias. Therefore there is a trade-off between the bias and variance of the estimated population mean. A bandwidth that provides a compromise between the two measures would therefore be desirable.

## 6. Mean Squared Error (MSE) of the Finite Population Mean Estimator $\hat{\bar{\bar{Y}}}$

The MSE of $\hat{\bar{\bar{Y}}}$ combines the bias and the variance terms of this estimator that is,

$$MSE\left(\hat{\bar{\bar{Y}}}\right) = E\left(\hat{\bar{\bar{Y}}} - \bar{\bar{Y}}\right)^2 \qquad (6.0)$$

$$MSE\left(\hat{\bar{\bar{Y}}}\right) = E\left(\hat{\bar{\bar{Y}}} - E\left[\hat{\bar{\bar{Y}}}\right] + E\left[\hat{\bar{\bar{Y}}}\right] - \bar{\bar{Y}}\right)^2 \qquad (6.1)$$

Expanding Equation (6.1) gives:

$$MSE\left(\hat{\bar{\bar{Y}}}\right) = E\left(\hat{\bar{\bar{Y}}} - E\left[\bar{\bar{Y}}\right]\right)^2 + E\left(E\left[\hat{\bar{\bar{Y}}}\right] - \left[\bar{\bar{Y}}\right]\right)^2$$
$$+ 2E\left(\hat{\bar{\bar{Y}}} - E\left[\hat{\bar{\bar{Y}}}\right]\right)\left(E\left[\hat{\bar{\bar{Y}}}\right] - \bar{\bar{Y}}\right) \qquad (6.2)$$

$$= \text{Var}\left(\hat{\bar{\bar{Y}}}\right) + Bias^2 + 0 \qquad (6.3)$$

Combining the bias in Equation (4.27) and the variance in Equation (5.17) and conditioning on the auxiliary values $x_{ij}$ of the auxiliary variables $X_{ij}$ then

$$MSE\left(\hat{\bar{\bar{Y}}}\bigg/X_{ij}=x_{ij}\right)$$

$$=\frac{1}{\left(MN\right)^2}\left\{\frac{\left(MN-mn\right)^2 H\left(w\right)\sigma^2_{\left(x_{ij}\right)}}{mnb\left(g\left(x_{ij}\right)\right)}+O\left(\frac{1}{MN}\left\{\frac{\left(MN-mn\right)^2}{mnb}+\frac{1}{mnb}\right\}\right)\right\}\quad(6.4)$$

$$+\frac{1}{MN}\left\{\left(\frac{MN-mn}{mn}\right)b^2 d_k C\left(x\right)+O\left(b^2\right)\right\}$$

$$MSE\left(\hat{\bar{\bar{Y}}}\bigg/X_{ij}=x_{ij}\right)$$

$$=\frac{1}{\left(MN\right)^2}\left\{\frac{\left(MN-mn\right)^2 H\left(w\right)\sigma^2_{\left(x_{ij}\right)}}{mnb\left(g\left(x_{ij}\right)\right)}\right.$$

$$+\left[\frac{\left(MN-mn\right)^2}{4\left(mn\right)^2\left(MN\right)^2}b^4 d_k^2\left[m''\left(x_{ij}\right)+\frac{2g'\left(x_{ij}\right)m'\left(x_{ij}\right)}{g\left(x_{ij}\right)}\right]^2\right.\quad(6.5)$$

$$\left.\left.+O\left(b^4\right)+\frac{1}{MN}\left(O\left\{\left(\frac{MN-mn}{mnb}\right)+\frac{1}{mnb}\right\}\right)\right]\right\}$$

where $H\left(w\right)=\int K\left(w\right)^2\,\mathrm{d}w$, $d_k=\int w^2 K\left(w\right)\mathrm{d}w$,

$C\left(x\right)=\frac{1}{2}m''\left(x_{ij}\right)+\left[g\left(x_{ij}\right)\right]^{-1}g'\left(x_{ij}\right)m'\left(x_{ij}\right)$ as used earlier in the rest of the derivations.

## 7. Conclusion

If the sample size is large enough, that is as $n\to N$ and $m\to M$ the *MSE* of $\left(\hat{\bar{\bar{Y}}}\right)$ in Equation (6.5) due to the kernel tends to zero for sufficiently a small bandwidth *b*. The estimator $\left(\hat{\bar{\bar{Y}}}\right)$ is therefore asymptotically consistent since its MSE converges to zero.

## References

[1] Singh, S. and Horn, S. (2000) Compromised Imputation in Survey Sampling. *Metrika*, **51**, 267-276. https://doi.org/10.1007/s001840000054

[2] Lee, H., Rancourt, E. and Särndal, C. (2002) Variance Estimation from Survey Data under Single Imputation. Survey Nonresponse, 315-328.

[3] Bethlehem, J.G. (2012) Using Response Probabilities for Assessing Representativity. Statistics Netherlands, *International Statistical Review*, **80**, 382-399.

[4] Ouma, C. and Wafula, C. (2005) Bootstrap Confidence Intervals for Model-Based Surveys. *East African Journal of Statistics*, **1**, 84-90.

[5] Onyango, C.O., Otieno, R.O. and Orwa, G.O. (2010) Generalized Model Based Confidence Intervals in Two Stage Cluster Sampling. *Pakistan Journal of Statistics and Operation Research*, **6**. https://doi.org/10.18187/pjsor.v6i2.128

[6] Dorfman, A.H. (1992) Nonparametric Regression for Estimating Totals in Finite Populations. In: *Proceedings of the Section on Survey Research Methods*, American

Statistical Association Alexandria, VA, 622-625.

[7]  Nadaraya, E.A. (1964) On Estimating Regression. *Theory of Probability & Its Applications*, **9**, 141-142. https://doi.org/10.1137/1109020

[8]  Watson, G.S. (1964) Smooth Regression Analysis. *Sankhya: The Indian Journal of Statistics, Series A*, 359-372.