



Rita Cristina Pinto de Sousa

Mestre em Estatística

Parameter Estimation in the Presence of Auxiliary Information

Dissertação para obtenção do Grau de Doutora em
Estatística e Gestão de Risco, Especialidade em Estatística

Orientador : Sat Gupta, Professor of Statistics,
University of North Carolina at Greensboro, USA

Co-orientador : Pedro Corte Real, Professor Auxiliar,
Universidade Nova de Lisboa, Portugal

Júri:

Presidente: Prof. Doutora Maria Luísa Dias de Carvalho de Sousa Leonardo

Arguentes: Prof. Doutor Dinis Duarte Ferreira Pestana
Prof. Doutora Maria Teresa Themido da Silva Pereira

Vogais: Prof. Doutor João Tiago Praça Nunes Mexia
Prof. Doutora Célia Maria Pinto Nunes
Prof. Doutor Sat Gupta
Prof. Doutor Pedro Alexandre da Rosa Corte Real



Rita Cristina Pinto de Sousa

Mestre em Estatística

Parameter Estimation in the Presence of Auxiliary Information

Dissertação para obtenção do Grau de Doutora em
Estatística e Gestão de Risco, Especialidade em Estatística

Orientador : Sat Gupta, Professor of Statistics,
University of North Carolina at Greensboro, USA

Co-orientador : Pedro Corte Real, Professor Auxiliar,
Universidade Nova de Lisboa, Portugal

Júri:

Presidente: Prof. Doutora Maria Luísa Dias de Carvalho de Sousa Leonardo

Arguentes: Prof. Doutor Dinis Duarte Ferreira Pestana
Prof. Doutora Maria Teresa Themido da Silva Pereira

Vogais: Prof. Doutor João Tiago Praça Nunes Mexia
Prof. Doutora Célia Maria Pinto Nunes
Prof. Doutor Sat Gupta
Prof. Doutor Pedro Alexandre da Rosa Corte Real

Parameter Estimation in the Presence of Auxiliary Information

Copyright © Rita Cristina Pinto de Sousa, Faculdade de Ciências e Tecnologia, Universidade Nova de Lisboa

A Faculdade de Ciências e Tecnologia e a Universidade Nova de Lisboa têm o direito, perpétuo e sem limites geográficos, de arquivar e publicar esta dissertação através de exemplares impressos reproduzidos em papel ou de forma digital, ou por qualquer outro meio conhecido ou que venha a ser inventado, e de a divulgar através de repositórios científicos e de admitir a sua cópia e distribuição com objectivos educacionais ou de investigação, não comerciais, desde que seja dado crédito ao autor e editor. O copyright dos capítulos 2, 3, 4, 5 e 6 foram transferidos dos autores para editoras e são reproduzidos sob permissão dos editores originais e sujeitos às restrições de cópia impostos pelos mesmos.

Aos meus Pais e Irmã

Acknowledgements

The last four years were a great challenge for me. I have to thank many people who made this interesting journey possible. Their scientific or emotional support was crucial to reach the end of my thesis.

- First of all I would like to thank my Mentors, the advisor Sat Gupta and the co-adviser Pedro Corte Real for sharing their great knowledge and experience with me. I have to point out their tremendous support and their orientation ability that made me easily forget and overcome the physical distance.

I thank Professor Pedro Corte Real especially for having encouraged me to invest in my school graduation which contributed to my growth as a researcher. His friendly receptivity, the sharing of his knowledge and his great ideas were fundamental for my motivation and for a better outcome. His experience and friendship are for me an important and true reference.

I thank Professor Sat Gupta especially for sharing with me his pleasure, knowledge and motivation in working on this theme. The chance to work with him allowed me to meet many people, to learn a lot and take a broad view of estimation with auxiliary information. He proved to be a very friendly person. I consider him an excellent reference as person, as a mentor and as a researcher.

I am deeply grateful to them. Their support greatly contributed for my development as research and as person.

- I would like to thank the Professors and people from the university staff who support the post-graduate students. A special thanks to the coordinator Manuel Esquivel for his support and his incentive to carry on, always understanding my constraints as working student.
- I would like to thank all my Co-Authors who helped me write the chapters of this thesis. Their collaboration improved and expanded my knowledge on several topics.

A special thanks to Javid for sharing with me his great theoretical knowledge and

experience and to Nursel for being a good reference as a young and active researcher.

- I also would like to thank the anonymous referees for making very constructive suggestions which resulted in a significant improvement over the original version of the papers to be found in the chapters of this thesis.
- Thanks to my work colleagues for their support and encouragement.
A special reference to Pedro and São, my classmates at Statistics Portugal, for their advice and for having always a friendly word encouraging me and giving me strength to carry on.
A special message of affection to my colleagues from afrolatINE for having always supported me and for having contributed to my well being and happiness. Who dances is happier.
- Thanks to my parents and my sister for their affection and unconditional support.
From the beginning my parents always supported me, encouraging me in my decision and doing everything they could to allow me to complete successfully this journey. They are my true inspiration.
My sister is a true example of a blood sister and especially a sister by heart. Thank you for being as you are.
- I would also like to thank my friends for all their support, for being always there and for their comprehension for my very often absence. They know who they are.
I can not forget to mention Graça, Joana, Magda, Marinha, Marta, Mónica, Rosário and Silvia. To them I owe my deep friendship and a lot of my well being as a person.
A special word for Rui who accompanied me in the final phase of this thesis. Thanks for your encouragement, for your patience, for your fellowship and affection.

Abstract

In survey research, there are many situations when the primary variable of interest is sensitive. The sensitivity of some queries can give rise to a refusal to answer or to false answers given intentionally. Survey can be conducted in a variety of settings, in part dictated by the mode of data collection, and these settings can differ in how much privacy they offer the respondent. The estimates obtained from a direct survey on sensitive questions would be subject to high bias. A variety of techniques have been used to improve reporting by increasing the privacy of the respondents.

The Randomized Response Technique (RRT), introduced by Warner in 1965, develops a random relation between the individual's response and the question. This technique provides confidentiality to respondents and still allows the interviewers to estimate the characteristic of interest at an aggregate level.

In this thesis we propose some estimators to improve the mean estimation of a sensitive variable based on a RRT by making use of available non-sensitive auxiliary information. In the first part of this thesis we present the ratio and the regression estimators as well as some generalizations in order to study the gain in the estimation over the ordinary RRT mean estimator. In chapters 4 and 5 we study the performance of some exponential type estimators, also based on a RRT. The final part of the thesis illustrates an approach to mean estimation in stratified sampling. This study confirms some previous results for a different sample design. An extensive simulation study and an application to a real dataset are done for all the study estimators to evaluate their performance. In the last chapter we present a general discussion referring to the main results and conclusions as well as showing an application to a real dataset which compares the performance of study estimators.

Keywords: Auxiliary variable; Exponential estimator; Randomized response technique; Ratio estimator; Regression estimator; Sensitive variable.

Resumo

Em estudos de pesquisa por inquérito existem muitas situações em que a variável de interesse é sensível. A sensibilidade de algumas questões pode dar origem a recusas na resposta ou a falsas respostas dadas de forma intencional. Os inquéritos podem assumir diversas configurações, em parte relacionadas com o método de recolha e com o grau de privacidade que é oferecido aos respondentes. As estimativas obtidas por inquérito direto em questões sensíveis estariam sujeitas a erros elevados. Muitas técnicas têm sido utilizadas para melhorar as respostas através do aumento de privacidade dos inquiridos. A Técnica de Resposta Aleatorizada, introduzida por Warner em 1965, desenvolve uma relação aleatória entre as respostas individuais e a questão. Esta técnica providencia confidencialidade aos respondentes e ainda permite aos entrevistadores estimar a característica de interesse num nível mais agregado.

Nesta tese propõem-se alguns estimadores para melhorar a estimação da média de uma variável sensível baseada numa técnica de resposta aleatorizada com recurso a informação auxiliar disponível não sensível. Na primeira parte da tese apresentam-se os estimadores da razão e da regressão bem como algumas generalizações para estudar o ganho na estimação face ao estimador ordinário da média. Nos capítulos 4 e 5 estuda-se a performance de alguns estimadores do tipo exponencial, também baseados numa técnica de resposta aleatorizada. A parte final da tese ilustra uma aproximação à estimação da média com amostragem estratificada. Este estudo vem confirmar resultados anteriores com um novo desenho amostral. Um extenso estudo de simulação e uma aplicação a dados reais são feitos para avaliar a performance de todos os estimadores. No último capítulo apresenta-se uma discussão geral, bem como uma aplicação a dados reais onde se compara a performance dos estimadores em estudo.

Palavras-chave: Estimador da razão; Estimador da regressão; Estimador exponencial; Técnica de resposta aleatorizada; Variável auxiliar; Variável sensível.

Contents

Contents	xiii
List of Figures	xv
List of Tables	xvii
Listings	xix
Abbreviations	xxi
1 General Introduction	1
References	3
2 Ratio Estimation of the Mean of a Sensitive Variable in the Presence of Auxiliary Information	5
Abstract	5
2.1 Introduction	6
2.2 Terminology	6
2.3 The Proposed Estimator	7
2.4 A Simulation Study	9
2.5 Numerical Example	11
2.6 Transformed Ratio Estimators	13
2.7 Conclusions	18
References	18
Appendix A - R Routines	21

3	Estimation of the Mean of a Sensitive Variable in the Presence of Auxiliary Information	31
	Abstract	31
3.1	Introduction	32
3.2	Terminology	32
3.3	The Ratio Estimator	33
3.4	Ordinary Regression Estimator	34
3.5	Generalized Regression-cum-ratio Estimator	35
3.6	The Simulation Study	37
3.7	Numerical Example	39
3.8	Conclusions	41
	References	41
	Appendix B - R Routines	43
4	Exponential Type Estimators of the Mean of a Sensitive Variable in the Presence of Non Sensitive Auxiliary Information	49
	Abstract	49
4.1	Introduction	50
4.2	Terminology	50
4.3	Estimators Review	51
4.4	Proposed Exponential Type Estimators	52
4.5	Comparison with Gupta et al. (2012) Estimators	55
4.6	Simulation Study	56
4.7	Numerical Example	58
4.8	Conclusions	60
	References	60
	Appendix C - R Routines	63
5	Improved Exponential Type Estimators of the Mean of a Sensitive Variable in the Presence of Non-Sensitive Auxiliary Information	71
	Abstract	71
5.1	Introduction	72
5.2	Terminology	72

5.3	Difference-cum-exponential Estimator (Koyuncu et al., 2013)	73
5.4	Proposed estimator	73
5.5	Simulation Study	75
5.6	Numerical Example	77
5.7	Conclusions	79
	References	79
	Appendix D - R Routines	81
6	Improved Mean Estimation of a Sensitive Variable Using Auxiliary Information in Stratified Sampling	89
	Abstract	89
6.1	Introduction	90
6.2	Terminology	90
6.3	Estimators Review	91
6.4	Proposed combined ratio estimator	92
6.5	Proposed combined regression estimator	94
6.6	A Simulation Study	95
6.7	Numerical Example	98
6.8	Conclusions	100
	References	100
	Appendix E - R Routines	103
7	General Discussion	113
7.1	Summary	113
7.2	Comparison of the main study estimators	113
7.3	Final Remarks	117
	References	118
	Appendix F - R Routines	120

List of Figures

7.1	Distribution of empirical <i>Bias</i>	115
7.2	Distribution of empirical <i>MSE</i>	117

List of Tables

2.1	Empirical <i>ARB</i> for RRT mean estimator and ratio estimator (bold).	10
2.2	Theoretical <i>ARB</i> for ratio estimator based on 1^{st} and 2^{nd} order (bold) approximation. . .	10
2.3	<i>MSE</i> correct up to 1^{st} and 2^{nd} order approximations and <i>PRE</i> for the ratio estimator relative to the RRT mean estimator.	11
2.4	Empirical <i>ARB</i> for the RRT mean estimator and the ratio estimator (bold).	12
2.5	Theoretical <i>ARB</i> for the RRT mean estimator and the ratio estimator.	13
2.6	<i>MSE</i> , corrected to 1^{st} order approximation, and <i>PRE</i> for the ratio estimator related to the RRT mean estimator.	13
2.7	Empirical <i>ARB</i> for the RRT mean estimator, the ratio estimator and for the transformed ratio estimators.	15
2.8	Theoretical <i>ARB</i> to 1^{st} order approximation for the RRT mean estimator, the ratio estimator and for the transformed ratio estimators.	16
2.9	Empirical <i>MSE</i> and theoretical (bold) <i>MSE</i> to 1^{st} order of approximation for the RRT mean estimator, the ratio estimator and for the transformed ratio estimators.	17
2.10	<i>PRE</i> for the transformed ratio estimator related to the ratio estimator based on 1^{st} order of approximation.	18
2.11	Calculations for the expression in (2.16).	18
3.1	<i>MSE</i> correct up to 1^{st} order approximation and <i>PRE</i> for the ratio estimator ($\hat{\mu}_R$), the regression estimator ($\hat{\mu}_{Reg}$) and the generalized regression-cum-ratio estimator ($\hat{\mu}_{GRR}$) relative to the RRT mean estimator.	38
3.2	<i>MSE</i> correct up to 1^{st} order approximation and <i>PRE</i> for the ratio estimator ($\hat{\mu}_R$), the regression estimator ($\hat{\mu}_{Reg}$) and the generalized regression-cum-ratio estimator ($\hat{\mu}_{GRR}$) relative to the RRT mean estimator.	40

4.1	Empirical <i>MSE</i> , theoretical <i>MSE</i> correct up to 1 st order approximation and <i>PRE</i> of all estimators.	57
4.2	Table 4.1 Continued.	58
4.3	<i>MSE</i> and <i>PRE</i> for the ratio estimator ($\hat{\mu}_R$), the regression estimator ($\hat{\mu}_{Reg}$), the generalized regression-cum-ratio estimator ($\hat{\mu}_{GRR}$) and the exponential estimator ($\hat{\mu}_{exp1}$) relative to the RRT mean estimator.	60
5.1	Empirical <i>ARB</i> for the difference-cum-exponential estimator ($\hat{\mu}_{DE}$) and for the improved exponential estimator ($\hat{\mu}_{IE}$).	76
5.2	Theoretical <i>ARB</i> for the difference-cum-exponential estimator ($\hat{\mu}_{DE}$) and for the improved exponential estimator ($\hat{\mu}_{IE}$).	76
5.3	Empirical and theoretical <i>MSE</i> for the difference-cum-exponential estimator ($\hat{\mu}_{DE}$) and for the improved exponential estimator ($\hat{\mu}_{IE}$).	77
5.4	Empirical <i>ARB</i> for the difference-cum-exponential estimator ($\hat{\mu}_{DE}$) and for the improved exponential estimator ($\hat{\mu}_{IE}$).	78
5.5	Theoretical <i>ARB</i> for the difference-cum-exponential estimator ($\hat{\mu}_{DE}$) and for the improved exponential estimator ($\hat{\mu}_{IE}$).	78
5.6	Empirical and theoretical <i>MSE</i> for the difference-cum-exponential estimator ($\hat{\mu}_{DE}$) and for the improved exponential estimator ($\hat{\mu}_{IE}$).	79
6.1	Empirical and Theoretical <i>MSE</i> , for the RRT mean estimator, ratio estimator (underlined) and regression estimator (bold); and corresponding <i>PRE</i> relative to the RRT mean estimator.	97
6.2	Empirical, theoretical <i>MSE</i> , <i>PRE</i> for the ratio estimator (underlined) and for the regression estimator (bold) relative to the RRT mean estimator and <i>PRE</i> for the simple random sample (<i>SRS</i>) relative to the stratified sample (<i>Str</i>).	99
7.1	Theoretical <i>ARB</i> for the estimators in comparison.	115
7.2	Empirical <i>MSE</i> , theoretical <i>MSE</i> correct up to 1 st order approximation and <i>PRE</i> for all the estimators in comparison relative to the RRT mean estimator.	116

Listings

2.1	R Code for Simulation Study of Proposed Estimator in Chapter 2	21
2.2	R Code for Simulation Study of Transformed Ratio Estimators in Chapter 2	24
2.3	R Code for Numerical Example of Proposed Estimator in Chapter 2	28
3.1	R Code for Simulation Study of Proposed Estimator in Chapter 3	43
3.2	R Code for Numerical Example of Proposed Estimator in Chapter 3	46
4.1	R Code for Simulation Study of Proposed Estimator in Chapter 4	63
4.2	R Code for Numerical Example of Proposed Estimator in Chapter 4	67
5.1	R Code for Simulation Study of Proposed Estimator in Chapter 5	81
5.2	R Code for Numerical Example of Proposed Estimator in Chapter 5	85
6.1	R Code for Simulation Study of Proposed Estimator in Chapter 6	103
6.2	R Code for Numerical Example of Proposed Estimator in Chapter 6	108
7.1	R Code for Numerical Example of Proposed Estimator in Chapter 7	120

List of Abbreviations

ARB - Absolute Relative Bias

Deff - Design Effect

ICT - Information and Communication Technologies

MES - Monthly Economic Survey

MSE - Mean Square Error

NACE - Statistical Classification of Economic Activities in the European Community

PRE - Percent Relative Efficiency

RRT - Randomized Response Technique

SRS - Simple Random Sample

SRSWOR - Simple Random Sampling Without Replacement

Str - Stratified Sample



General Introduction

One of the major problems in survey research involving sensitive questions is the social desirability response bias (Edwards, 1957). For various reasons individuals in a sample survey may prefer not to confide to the interviewer the correct answers to certain questions. In such cases the individuals may elect not to reply at all or to reply with incorrect answers. The resulting evasive answer bias is ordinarily difficult to assess. That bias is potentially removable through allowing the interviewer to maintain privacy using a randomization device (Warner, 1965).

Randomized response is a research method used in structured survey interview. It was first proposed by Warner in 1965 and later modified by Greenberg et al. in 1969. This technique allows respondents to respond to sensitive issues while maintaining confidentiality. It provides confidentiality to respondents through a random relation between the individual's response and the question. It still allows the interviewers to estimate the characteristic of interest at an aggregate level.

Gupta and Thornton (2002) showed that Randomized Response Technique (RRT) is effective in circumventing the social desirability response bias, and is more friendly and portable than other methods such as the method which uses a bogus pipeline (Jones and Sigall, 1971).

RRT models may be classified as Full RRT model, Partial RRT model or Optional RRT model depending on the level of scrambling. In the Full RRT model (Eichhorn and Hayre, 1983) all the respondents are asked to provide a scrambled response. When a predetermined proportion of randomly selected respondents are asked to provide a true response we have a Partial RRT model (Mangat and Singh, 1990). Gupta et al. (2002) proposed an Optional RRT model where the respondents are allowed to report a true response or a

scrambled response depending on whether the respondents find the question sensitive or not.

In RRT work, generally the focus is on the estimation of the mean of a sensitive variable or the prevalence of a sensitive characteristic in the population. The mean can be estimated by using one of many RRT but we propose some estimators which improve the mean estimation considerably by using non-sensitive auxiliary information. In such cases, one will be able to observe an auxiliary variable directly but will have to rely on some RRT to collect information on the variable of interest, resulting from a sensitive issue. Given that our main aim is to evaluate the performance of the mean estimator in the presence of auxiliary information, we opt for using an additive Full RRT method to scramble the sensitive variable.

The main goal of this thesis is to improve the parameter estimation of a sensitive variable in the presence of auxiliary information. For that purpose we introduce some estimators for the population mean based on the additive Full RRT technique. Expressions are derived for the *Bias* and Mean Square Error (*MSE*) for all the proposed estimators. Furthermore, an extensive simulation study and an application to a real dataset are done for all the study estimators. All the applications are developed using the statistical software R [1].

This thesis is based on five papers to be found in chapters 2–6. Each chapter presents, at least, a new estimator and evaluates its performance comparing it to the other estimators previously proposed. The contents of this thesis are as follows:

- In **Chapter 2** we propose a ratio estimator for the mean of a sensitive variable using information from a non-sensitive auxiliary variable. We generalize the proposed estimator to the case of transformed ratio estimators. We show that there is hardly any difference in the first order and second order approximations for *MSE* even for small sample sizes. We also show that the proposed estimator does better than the ordinary RRT mean estimator which does not use the auxiliary information (Sousa et al., 2010).
- In **Chapter 3** we introduce a regression estimator which performs better than the ratio estimator even for modest correlation between the primary and the auxiliary variables. We consider a generalized regression-cum-ratio estimator that has even smaller *MSE*. It is shown that the proposed regression estimator performs better than the ratio estimator and the ordinary RRT mean estimator that does not utilize the auxiliary information (Gupta et al., 2012).
- In **Chapter 4** we propose exponential type estimators using one and two auxiliary variables to improve the efficiency of mean estimator based on a RRT. It is shown

the proposed exponential type estimators are more efficient than the existing estimators described in Sousa et al. (2010) and Gupta et al. (2012)(Koyuncu et al., 2013).

- In **Chapter 5** we propose an improved exponential type estimator which is more efficient than the Koyuncu et al. (2013) estimator, which in turn was shown to be more efficient than the usual mean estimator, ratio estimator, regression estimator, and the Gupta et al. (2012) estimator. It is shown that the improved difference-cum-exponential estimator can produce further improvement relative to other estimators previously proposed (Gupta et al., 2013).
- In **Chapter 6** we extend the ratio and regression estimators to the stratified sampling setting. Although both the ratio and regression estimators perform better than the ordinary RRT mean estimator, the improvement is much larger with the regression estimator. The results agree with the findings of Sousa et al. (2010) and Gupta et al. (2012) in simple random sampling. We show that the advantage of using the RRT in the presence of auxiliary information still holds in the context of stratified sampling (Sousa et al., 2013).
- In **Chapter 7** we present a general discussion referring to the main results and conclusions. We present a study with a real dataset and we show the numerical results for the *Bias* and *MSE*, as well as graphic evidence which illustrates the performance of the main study estimators.

In the last part of each chapter we attach the R routines developed for the simulation studies and for the numerical examples.

References

- EDWARDS, A. L. 1957. *The social desirability variable in personality assessment and research*, New York: Dryden, Praeger.
- EICHHORN, B. H. & HAYRE, L. S. 1983. Scrambled randomized response methods for obtaining sensitive quantitative data. *Journal of Statistical Planning and Inference*, 7, 307-316.
- GREENBERG, B., ABDUL-ELA, A., SIMMONS, W. & HORVITZ, D. 1969. The unrelated question randomized response model: theoretical framework. *Journal of American Statistical Association*, 520-539.
- GUPTA, S. N., GUPTA, B. C. & SINGH, S. 2002. Estimation of sensitivity level of personal interview survey questions. *Journal of Statistical Planning and Inference*, 100, 239-247.
- GUPTA, S., SHABBIR, J., SOUSA, R. & REAL, P. C. 2012. Estimation of the Mean of a

Sensitive Variable in the Presence of Auxiliary Information. *Communications in Statistics - Theory and Methods*, 41(13-14), 2394-2404.

GUPTA, S., SHABBIR, J., SOUSA, R. & REAL, P. C. 2013. Improved exponential type estimators of the mean of a sensitive variable in the presence of non-sensitive auxiliary information. (*submitted*)

GUPTA, S. & THORNTON, B. 2002. Circumventing social desirability response bias in personal interview surveys. *American Journal of Mathematical and Management Sciences*, 22, 369-383.

JONES, E. E. & SIGALL, H. 1971. The Bogus Pipeline: A New Paradigm for Measuring Affect and Attitude. *Psychological Bulletin*, 76, 349-364.

KOYUNCU, N., GUPTA, S. & SOUSA, R. 2013. Exponential type estimators of the mean of a sensitive variable in the presence of non-sensitive auxiliary information. *Communications in Statistics - Simulation and Computation*. (*accepted*).

MANGAT, N. S. & SINGH, R. 1990. An Alternative Randomized Response Procedure. *Biometrika*, 77, 439-442.

SOUSA, R., GUPTA, S., SHABBIR, J. & REAL, P. C. 2013. Improved Mean Estimation of a Sensitive Variable Using Auxiliary Information in Stratified Sampling. *Journal of Statistics and Management Systems*. (*submitted*).

SOUSA, R., SHABBIR, J., REAL, P. C. & GUPTA, S. 2010. Ratio estimation of the mean of a sensitive variable in the presence of auxiliary information. *Journal of Statistical Theory and Practice*, 4(3), 495-507.

WARNER, S. L. 1965. Randomized response: a survey technique for elimination evasive answer bias. *Journal of American Statistical Association*, 60, 63-69.

[1] The R Project for Statistical Computing: www.r-project.org.



Ratio Estimation of the Mean of a Sensitive Variable in the Presence of Auxiliary Information

Abstract

We propose a ratio estimator for the mean of sensitive variable utilizing information from a non-sensitive auxiliary variable. Expressions for the *Bias* and Mean Square Error (*MSE*) of the proposed estimator (correct up to first and second order approximations) are derived. We show that the proposed estimator does better than the ordinary Randomized Response Technique (RRT) mean estimator that does not utilize the auxiliary information. We also show that there is hardly any difference in the first order and second order approximations for *MSE* even for small sample sizes. We also generalize the proposed estimator to the case of transformed ratio estimators but these transformations do not result in any significant reduction in *MSE*. An extensive simulation study is presented to evaluate the performance of the proposed estimator. The procedure is also applied to some financial data (purchase orders (sensitive variable) and gross turnover (non-sensitive variable)) in 2009 for 5090 companies in Portugal from a survey on Information and Communication Technologies (ICT) usage.

Published as: SOUSA, R., SHABBIR, J. REAL, P. C. & GUPTA, S. (2010). Ratio estimation of the mean of a sensitive variable in the presence of auxiliary information. *Journal of Statistical Theory and Practice*, 4(3), 495-507.

2.1 Introduction

In survey research, there are many situations when the primary variable of interest (Y) is sensitive and direct observation on this variable may not be possible. However, we may be able to directly observe a highly correlated auxiliary variable (X). For example, Y may be the number of abortions a woman might have had in her life and X may be her age. Similarly Y may be the total purchase orders in a year for a company and X may be the total turn-over for that company in that year. In such cases, one will be able to observe X directly but will have to rely on some Randomized Response Technique (RRT) to collect information on Y . In such situations, mean of Y can be estimated by using one of many randomized response techniques but this estimator can be improved considerably by utilizing information from the auxiliary variable X . Many authors have presented ratio estimators when both Y and X are directly observable. These include Kadilar and Cingi (2006), Turgut and Cingi (2008), Singh and Vishwakarma (2008), Koyuncu and Kadilar (2009) and Shabbir and Gupta (2010).

Also, many authors have estimated the mean of a sensitive variable when the primary variable is sensitive and there is no auxiliary variable available. These include Eichhorn and Hayre (1983), Gupta and Shabbir (2004), Gupta et al. (2002), Saha (2008) and Gupta et al. (2010).

In this paper, we propose a ratio estimator where the RRT estimator of the mean of Y is further improved by using information on an auxiliary variable X . Expressions for the *Bias* and *MSE* for the proposed estimator are derived, correct up to both the first order and second order approximations. It is shown that the two approximations are very similar even for moderate sample size. We also observe that there is considerable reduction in *MSE* when auxiliary information is used, particularly when the correlation between the study variable and the auxiliary variable is high.

2.2 Terminology

Let Y be the study variable, a sensitive variable which cannot be observed directly. Let X be a non-sensitive auxiliary variable which is strongly correlated with Y . Let S be a scrambling variable independent of Y and X . The respondent is asked to report a scrambled response for Y given by $Z = Y + S$ but is asked to provide a true response for X . Let a random sample of size n be drawn without replacement from a finite population $U = (U_1, U_2, \dots, U_N)$. For the i^{th} unit ($i = 1, 2, \dots, N$), let y_i and x_i respectively be the values of the study variable Y and auxiliary variable X . Moreover, let $\bar{y} = \frac{\sum_{i=1}^n y_i}{n}$, $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$ and $\bar{z} = \frac{\sum_{i=1}^n z_i}{n}$ be the sample means and $\bar{Y} = E(Y)$, $\bar{X} = E(X)$ and $\bar{Z} = E(Z)$ be the population means for Y , X and Z , respectively. We assume that \bar{X} is known and $\bar{S} = E(S) = 0$. Thus, $E(Z) = E(Y)$. Let us also define $\delta_z = \frac{\bar{z} - \bar{Z}}{\bar{Z}}$ and $\delta_x = \frac{\bar{x} - \bar{X}}{\bar{X}}$, such that

$$E(\delta_i) = 0, i = z, x.$$

If information on X is ignored, then an unbiased estimator of μ_Y is the ordinary sample mean (\bar{z}) given by (2.1) below

$$\hat{\mu}_Y = \bar{z}. \quad (2.1)$$

The mean square error (MSE) of $\hat{\mu}_Y$ is given by

$$MSE(\hat{\mu}_Y) = \frac{1-f}{n} (S_y^2 + S_s^2), \quad (2.2)$$

where

$$f = n/N, S_y^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{Y})^2, S_x^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{X})^2 \text{ and } S_s^2 = \frac{1}{N-1} \sum_{i=1}^N (s_i - \bar{S})^2.$$

2.3 The Proposed Estimator

We propose the following ratio estimator for estimating the population mean of the study variable Y using the auxiliary variable X :

$$\begin{aligned} \hat{\mu}_R &= \bar{z} \left(\frac{\bar{X}}{\bar{x}} \right) \\ &= \bar{Z} (1 + \delta_z) (1 + \delta_x)^{-1}. \end{aligned} \quad (2.3)$$

Using Taylor's approximation and retaining terms of order up to 4, (2.3) can be rewritten as

$$\hat{\mu}_R - \bar{Z} \cong \bar{Z} \{ \delta_z - \delta_x - \delta_z \delta_x + \delta_x^2 - \delta_x^3 + \delta_x^4 + \delta_z \delta_x^2 - \delta_z \delta_x^3 \}. \quad (2.4)$$

Under the assumption of bivariate normality (see Sukhatme and Sukhatme, 1984), we have $E(\delta_z^2) = \frac{1-f}{n} C_z^2$, $E(\delta_x^2) = \frac{1-f}{n} C_x^2$, $E(\delta_x \delta_z) = \frac{1-f}{n} C_{zx}$, where $C_{zx} = \rho_{zx} C_z C_x$ and C_z and C_x are the coefficients of variation of Z and X , respectively. Also we have:

$$\begin{aligned} E(\delta_z \delta_x^3) &= \left(\frac{1-f}{n} \right)^2 3\rho_{zx} C_z C_x^3, & E(\delta_z^2 \delta_x^2) &= \left(\frac{1-f}{n} \right)^2 (1 + 2\rho_{zx}^2) C_z^2 C_x^2, \\ E(\delta_x^4) &= \left(\frac{1-f}{n} \right)^2 3C_x^4, & E(\delta_z \delta_x^2) &= E(\delta_z^2 \delta_x) = E(\delta_x^3) = 0, \end{aligned}$$

and

$$C_z^2 = C_y^2 + \frac{S_s^2}{\bar{Y}^2}, \rho_{zx} = \frac{\rho_{yx}}{\sqrt{1 + \frac{S_s^2}{S_y^2}}}.$$

Recognizing that $\bar{Z} = \bar{Y}$ in Equation (2.4), we can get expressions for the Bias of $\hat{\mu}_R$,

correct up to second order of approximation, as given by

$$Bias^{(2)}(\hat{\mu}_R) \cong Bias^{(1)}(\hat{\mu}_R) + 3 \left(\frac{1-f}{n} \right)^2 \bar{Y} [C_x^4 - \rho_{yx} C_y C_x^3], \quad (2.5)$$

where

$$Bias^{(1)}(\hat{\mu}_R) = \left(\frac{1-f}{n} \right) \bar{Y} [C_x^2 - \rho_{yx} C_y C_x] \quad (2.6)$$

is the *Bias* corresponding to first order of approximation.

Similarly from (2.4), *MSE* of $\hat{\mu}_R$, correct up to second order of approximation, is given by

$$MSE^{(2)}(\hat{\mu}_R) = E(\hat{\mu}_R - \bar{Z})^2 \cong \bar{Z}^2 E\{\delta_z - \delta_x - \delta_z \delta_x + \delta_x^2 - \delta_x^3 + \delta_x^4 + \delta_z \delta_x^2 - \delta_z \delta_x^3\}^2$$

or

$$MSE^{(2)}(\hat{\mu}_R) \cong \bar{Z}^2 E\{\delta_z^2 + \delta_x^2 - 2\delta_z \delta_x + 3\delta_z^2 \delta_x^2 + 3\delta_x^4 - 6\delta_z \delta_x^3 - 2\delta_z^2 \delta_x + 4\delta_z \delta_x^2 - 2\delta_x^3\}.$$

Since $\bar{Z} = \bar{Y}$, we have

$$MSE^{(2)}(\hat{\mu}_R) \cong MSE^{(1)}(\hat{\mu}_R) + 3\bar{Y}^2 \left(\frac{1-f}{n} \right)^2 C_x^2 [(1 + 2\rho_{yx}^2)C_y^2 + 3C_x^2 - 6\rho_{yx} C_y C_x], \quad (2.7)$$

where

$$MSE^{(1)}(\hat{\mu}_R) \cong \left(\frac{1-f}{n} \right) \bar{Y}^2 (C_y^2 + C_x^2 - 2\rho_{yx} C_y C_x) \quad (2.8)$$

is the *MSE* corresponding to the first order approximation. The difference between the two approximations for *MSE* is given by

$$3\bar{Y}^2 \left(\frac{1-f}{n} \right)^2 C_x^2 [(1 + 2\rho_{yx}^2)C_y^2 + 3C_x^2 - 6\rho_{yx} C_y C_x],$$

and it converges to zero as $n \rightarrow N$. Our simulation results in Section 2.4 will also confirm this pattern.

According to the first order of approximation, $MSE^{(1)}(\hat{\mu}_R) < MSE(\hat{\mu}_Y)$ if

$$\left(\rho_{yx} - \frac{1}{2} \frac{C_x}{C_y} \right) > 0. \quad (2.9)$$

If second order approximation is used, we can easily see that $MSE^{(2)}(\hat{\mu}_R) < MSE(\hat{\mu}_Y)$ if

$$2\rho_{yx} \frac{C_x}{C_y} + 3 \left(\frac{1-f}{n} \right) [6\rho_{yx} C_x C_y - 3C_x^2 - (1 + 2\rho_{yx}^2)C_y^2] > 1. \quad (2.10)$$

2.4 A Simulation Study

In this section, we conduct a simulation study with particular focus on the following two issues:

- How does the ratio estimator $\hat{\mu}_R$ compare with $\hat{\mu}_Y$ the RRT mean estimator $\hat{\mu}_Y$;
- How do the *Bias* and *MSE* for the ratio estimator, correct up to second order of approximation, compare with the *Bias* and *MSE* expressions correct up to first order of approximation.

We considered 3 bivariate normal populations with different covariance matrices to represent the distribution of (Y, X) . The scrambling variable S is taken to be a normal distribution with mean equal to zero and standard deviation equal to 10% of the standard deviation of X . The reported response is given by $Z = Y + S$.

All of the simulated populations have theoretical mean of $[Y, X]$ as $\mu = [2, 2]$ and covariance matrices as given below.

Population 1

$$N = 1000$$

$$\Sigma = \begin{bmatrix} 9 & 1.9 \\ 1.9 & 4 \end{bmatrix}, \rho_{XY} = 0.3167.$$

Population 2

$$N = 1000$$

$$\Sigma = \begin{bmatrix} 10 & 3 \\ 3 & 2 \end{bmatrix}, \rho_{XY} = 0.6708.$$

Population 3

$$N = 1000$$

$$\Sigma = \begin{bmatrix} 6 & 3 \\ 3 & 2 \end{bmatrix}, \rho_{XY} = 0.8660.$$

For each population we considered five sample sizes: $n = 20, 50, 100, 200$ and 300 .

The Absolute Relative Bias (*ARB*) for the two estimators is given by $|Bias(\hat{\mu}_Y)/\bar{Y}|$ and $|Bias(\hat{\mu}_R)/\bar{Y}|$. We estimate the *ARB* using 5000 samples of size n selected from each population. The empirical *ARB* values for both estimators are given in Table 2.1. As expected, the *ARB* generally decreases as the sample size increases, with some exceptions due to random fluctuations.

The RRT mean estimator should generally perform better than the ratio estimator because this is an unbiased estimator. Nevertheless, the ratio estimator produces fairly good results.

Table 2.1: Empirical *ARB* for RRT mean estimator and ratio estimator (bold).

Population		Empirical <i>ARB</i>				
<i>N</i>	ρ_{XY}	<i>n</i> = 20	<i>n</i> = 50	<i>n</i> = 100	<i>n</i> = 200	<i>n</i> = 300
1000	0.3549	0.0021	0.0011	0.0010	0.0018	0.0016
		0.0223	0.0071	0.0057	0.0006	0.0009
	0.6965	0.0010	0.0021	0.0014	0.0014	0.0011
		0.0193	0.0061	0.0015	0.0029	0.0021
	0.8783	0.0012	0.0013	0.0008	0.0013	0.0011
		0.0181	0.0063	0.0023	0.0026	0.0020

The theoretical *ARB* results for the ratio estimator, correct up to first and second degree of approximation, are presented in Table 2.2.

One can see that second order approximation as compared to first order approximation does not result in major difference in *ARB* even for modest sample size of *n* = 20 and 50.

Table 2.2: Theoretical *ARB* for ratio estimator based on 1st and 2nd order (bold) approximation.

Population		Theoretical <i>ARB</i>				
<i>N</i>	ρ_{XY}	<i>n</i> = 20	<i>n</i> = 50	<i>n</i> = 100	<i>n</i> = 200	<i>n</i> = 300
1000	0.3549	0.0224	0.0087	0.0041	0.0018	0.0011
		0.0258	0.0092	0.0042	0.0019	0.0011
	0.6965	0.0155	0.0060	0.0029	0.0013	0.0007
		0.0167	0.0062	0.0029	0.0013	0.0007
	0.8783	0.0142	0.0055	0.0026	0.0012	0.0007
		0.0153	0.0057	0.0026	0.0012	0.0007

Table 2.3 below gives empirical and theoretical *MSE*'s for the ratio estimator based on both the first order and second order approximations. As we see from the table, there is hardly a difference between the two approximations even for small samples. Hence the Percent Relative Efficiency (*PRE*) is calculated based on first order of approximation only. We use the following expression to find the *PRE* of ratio estimator as compared to the RRT mean estimator:

$$PRE = \frac{MSE(\hat{\mu}_Y)}{MSE(\hat{\mu}_R)} \times 100.$$

All the percent relative efficiencies are greater than 100 indicating that the ratio estimator is better than the RRT mean estimator.

There are small differences between MSE values based on first and second order approximation for smaller sample sizes ($n=20$ and 50) but the MSE values are very similar when the sample size is larger. We can also note that the ratio estimator gets more and more efficient as the coefficient of correlation between X and Y increases. We can further note that for small correlation values, the ratio estimator may not be better than the RRT mean estimator, particularly so if sample size is small.

Table 2.3: MSE correct up to 1^{st} and 2^{nd} order approximations and PRE for the ratio estimator relative to the RRT mean estimator.

Population		MSE Estimation			MSE Condition		PRE		
N	ρ_{XY}	n	Empirical	1^{st} Order	2^{nd} Order	1^{st} Order ¹	2^{nd} Order ²	1^{st} Order	2^{nd} Order
1000	0.3549	20	0.5782	0.4462	0.5249		0.6947		89.15
		50	0.1837	0.1730	0.1848		0.9464		98.15
		100	0.0819	0.0820	0.0846	0.0340	1.0304	104.86	101.57
		200	0.0358	0.0364	0.0370		1.0723		103.37
		500	0.0219	0.0219	0.0214		1.0863		103.398
	0.6965	20	0.3434	0.3036	0.3327		2.9075		156.65
		50	0.1202	0.1177	0.1221		3.0887		165.49
		100	0.0548	0.0558	0.0568	0.4785	3.1492	171.63	168.67
		200	0.0248	0.0248	0.0250		3.1794		170.30
		500	0.0152	0.0145	0.0145		3.1894		170.85
	0.8783	20	0.1178	0.1012	0.1139		2.9424		274.89
		50	0.0406	0.0392	0.0412		3.0185		294.99
		100	0.0183	0.0186	0.0190	0.5919	3.0439	309.31	302.36
		200	0.0083	0.0083	0.0083		3.0565		306.18
		500	0.0050	0.0048	0.0048		3.0608		307.48

¹ MSE comparison condition based on 1^{st} order approximation given in expression (2.9).

² MSE comparison condition based on 2^{nd} order approximation given in expression (2.10).

2.5 Numerical Example

We now compare the RRT mean estimator and the ratio estimator using a real data set. The data come from a sample from the survey on Information and Communication Technologies (ICT) usage in enterprises in 2009 with seat in Portugal (Smilhily and Storm, 2010). This survey intends to promote the development of the national statistical system in the information society and to contribute to a deeper knowledge about the usage of ICT by enterprises. The target population covers all industries with one and more persons employed in the sections of economic activity C (Manufacturing) to N (Administrative and support service activities) and S (Other service activities), from NACE¹ Rev. 2 (Eurostat, 2008). The data are essentially collected using Electronic Data Interchange, applying direct connection between information systems at the respondent and the National Statistics Institute. For some enterprises the paper questionnaire is still used. The

¹NACE is derived from the French title "Nomenclature générale des Activités économiques dans les Communautés Européennes" (Statistical classification of economic activities in the European Communities).

questions in the structural business surveys mainly deal with characteristics that can be found in the organisations' annual reports and financial statements, such as employment, turnover and investment.

In our application the study variable Y is the purchase orders in 2009, collected by the ICT survey in that year. This is typically a confidential variable for enterprises, only known from business surveys. The auxiliary variable X is the turnover of each enterprise. This information can be easily obtained from enterprise records available in the public domain, as administrative information. In 2009 the population survey contained approximately 278000 enterprises and we know the value of X for all these enterprises. The purchase orders information was collected in the ICT survey and we have the values of Y for 5090 enterprises (which answered this question in the ICT survey in 2009). For this study, these 5090 enterprises are considered as our population. The scrambling variable S is taken to be a normal random variable with mean equal to zero and standard deviation equal to 10% of the standard deviation of X , that is $\sigma_S = 0.1\sigma_X$. The reported response is given by $Z = Y + S$ (the purchase order value plus a random quantity). The variables Y and X are strongly correlated so we can take advantage of this correlation by using the ratio estimator. In the next tables we present the results for the RRT mean estimator and for the ratio estimator for different sample sizes.

Population Characteristics:

$N = 5090, \rho_{XY} = 0.9832$
$\mu_X = 32.53, \mu_Y = 26.06, \sigma_X = 183.42, \sigma_Y = 67.07$ (in millions of Euros)
$\gamma_1^X = 31.54, \gamma_1^Y = 36.12, \gamma_2^X = 1481.08, \gamma_2^Y = 1839.13$

where γ_1 and γ_2 are the coefficients of *skewness* and *kurtosis*, respectively. We use the following samples sizes in our simulation study: $n = 100, 200, 300, 400, 500, 1000, 1500$ and 2000.

The empirical *ARB* values for both estimators, based on 5000 iterations, are given in Table 2.4. As expected, the bias decreases as the sample size increases, except for some random fluctuation. We expect the RRT mean estimator to perform better than the ratio estimator because this is an unbiased estimator, however, we don't see major differences between the two for larger samples.

Table 2.4: Empirical *ARB* for the RRT mean estimator and the ratio estimator (bold).

Population		Empirical <i>ARB</i>							
N	ρ_{XY}	$n = 100$	$n = 200$	$n = 300$	$n = 400$	$n = 500$	$n = 1000$	$n = 1500$	$n = 2000$
5090	0.9832	0.0219	0.0002	0.0096	0.0107	0.0163	0.0145	0.0106	0.0096
		0.0284	0.0198	0.0171	0.0183	0.0166	0.0149	0.0127	0.0121

The theoretical *ARB* results for the ratio estimator, correct up to first degree of approximation, are presented in Table 2.5. We use only the first order approximations from here on since the first and second order approximations are very similar, as we have seen earlier.

Table 2.5: Theoretical *ARB* for the RRT mean estimator and the ratio estimator.

Population		Theoretica <i>ARB</i>							
<i>N</i>	ρ_{XY}	<i>n</i> = 100	<i>n</i> = 200	<i>n</i> = 300	<i>n</i> = 400	<i>n</i> = 500	<i>n</i> = 1000	<i>n</i> = 1500	<i>n</i> = 2000
5090	0.9832	0.0368	0.0180	0.0118	0.0086	0.0068	0.0030	0.0018	0.0011

Table 2.6 presents the results for the empirical *MSE* estimates, the theoretical estimates, correct up to first degree of approximation and the *PRE* of ratio estimator relative to the RRT mean estimator.

Table 2.6: *MSE*, corrected to 1st order approximation, and *PRE* for the ratio estimator related to the RRT mean estimator.

Population		<i>n</i>	<i>MSE</i> Estimation		<i>PRE</i>
<i>N</i>	ρ_{XY}		Empirical	Theoretical	
5090	0.9832	100	12.8924	15.2630	2286.36
		200	6.4608	7.4786	
		300	4.4498	4.8838	
		400	3.5279	3.5864	
		500	2.7380	2.8079	
		1000	1.4117	1.2510	
		1500	0.8805	0.7321	
		2000	0.6033	0.4726	

Clearly the ratio estimator performs better than the RRT mean estimator for the real data also. The effect of sample size on the *PRE* calculation is neutralized when first order approximation is used, as can be seen from Equations (2.2) and (2.8).

2.6 Transformed Ratio Estimators

Now consider the transformed ratio estimator:

$$\hat{\mu}_{TR} = \bar{z} \left(\frac{c\bar{X} + d}{c\bar{x} + d} \right), \quad (2.11)$$

where *c* and *d* are the unit-free parameters, which may be quantities such as the coefficient of *skewness* and coefficient of *kurtosis* for *X*. Many researchers have used transformed ratio estimators. These include Sisodia and Dwivedi (1981), Singh et al. (1973), Kulkarni (1977), Upadhyaya and Singh (1999), Upadhyaya et al. (2000) and Chandra and Singh (2005).

We can rewrite (2.11) using relative error terms in the form

$$\hat{\mu}_{TR} = \bar{z} (1 + \delta_z) (1 + \eta\delta_x)^{-1}, \quad (2.12)$$

where $\eta = \frac{c\bar{X}}{c\bar{X} + d}$.

Expanding (2.12), the *Bias*, correct up to first order of approximation, is given by

$$Bias^{(1)}(\hat{\mu}_{TR}) \cong \left(\frac{1-f}{n} \right) \bar{Y} \{ \eta^2 C_x^2 - \eta \rho_{yx} C_y C_x \}. \quad (2.13)$$

By (2.6) and (2.13) $Bias^{(1)}(\hat{\mu}_{TR}) < Bias^{(1)}(\hat{\mu}_R)$ if

$$(\eta - 1) \left\{ \rho_{yx} - \frac{(\eta + 1)C_x}{C_y} \right\} > 0. \quad (2.14)$$

Similarly *MSE* of $\hat{\mu}_{TR}$, to first order of approximation, is given by

$$MSE^{(1)}(\hat{\mu}_{TR}) \cong \left(\frac{1-f}{n} \right) \bar{Y}^2 (C_y^2 + \eta^2 C_x^2 - 2\eta \rho_{yx} C_y C_x). \quad (2.15)$$

By (2.8) and (2.15) $MSE^{(1)}(\hat{\mu}_{TR}) < MSE^{(1)}(\hat{\mu}_R)$ if

$$(\eta - 1) \left\{ \rho_{yx} - \frac{(\eta + 1)C_x}{2C_y} \right\} > 0. \quad (2.16)$$

Now we conduct a simulation study with particular focus on the comparison between the ratio estimator $\hat{\mu}_R$ and the transformed ratio estimator $\hat{\mu}_{TR}$. We considered the same three bivariate normal populations as in the previous simulation study (Section 2.4).

The scrambling variable S is taken to be a normal random variable with mean equal to zero and the standard deviation equal to 10% of the standard deviation of X . The reported response is given by $Z = Y + S$. To compare these estimators, we present the results for the RRT mean estimator ($\hat{\mu}_Y$), the ratio estimator ($\hat{\mu}_R$) and for transformed ratio estimator $\hat{\mu}_{TRi}$ ($i = 1, 2, 3, 4$) with four different combinations of parameters c and d :

1. $\hat{\mu}_{TR1} = \bar{z} \left(\frac{c\bar{X} + d}{c\bar{x} + d} \right)$,
where $c = 1$ and $d =$ coefficient of *skewness*;
2. $\hat{\mu}_{TR2} = \bar{z} \left(\frac{c\bar{X} + d}{c\bar{x} + d} \right)$,
where $c = 1$ and $d =$ coefficient of *kurtosis*;
3. $\hat{\mu}_{TR3} = \bar{z} \left(\frac{c\bar{X} + d}{c\bar{x} + d} \right)$,
where $c =$ coefficient of *skewness* and $d =$ coefficient of *kurtosis*;

$$4. \hat{\mu}_{TRA} = \bar{z} \left(\frac{c\bar{X} + d}{c\bar{x} + d} \right),$$

where c = coefficient of *kurtosis* and d = coefficient of *skewness*.

The empirical *ARB* values for these six estimators are given in Table 2.7.

Table 2.7: Empirical *ARB* for the RRT mean estimator, the ratio estimator and for the transformed ratio estimators.

Population			Empirical <i>ARB</i>					
N	ρ_{XY}	n	$\hat{\mu}_Y$	$\hat{\mu}_R$	$\hat{\mu}_{TR1}$	$\hat{\mu}_{TR2}$	$\hat{\mu}_{TR3}$	$\hat{\mu}_{TR4}$
1000	0.3209	20	0.0002	0.0337	0.0435	0.0006	0.0026	0.0366
		50	0.0007	0.0118	0.0146	0.0002	0.0019	0.0126
		100	0.0003	0.0052	0.0065	0.0000	0.0009	0.0056
		150	0.0000	0.0032	0.0040	0.0001	0.0003	0.0035
		200	0.0012	0.0025	0.0030	0.0008	0.0016	0.0027
		300	0.0020	0.0041	0.0045	0.0023	0.0021	0.0043
	0.6746	20	0.0011	0.0122	0.0113	0.0111	0.0018	0.0119
		50	0.0004	0.0042	0.0038	0.0037	0.0015	0.0041
		100	0.0001	0.0022	0.0021	0.0021	0.0005	0.0022
		150	0.0005	0.0016	0.0015	0.0016	0.0002	0.0016
		200	0.0010	0.0005	0.0005	0.0001	0.0013	0.0005
		300	0.0015	0.0013	0.0014	0.0011	0.0016	0.0014
	0.8684	20	0.0006	0.0120	0.0115	0.0108	0.0013	0.0119
		50	0.0005	0.0041	0.0039	0.0036	0.0012	0.0040
		100	0.0001	0.0018	0.0017	0.0017	0.0005	0.0018
		150	0.0002	0.0010	0.0010	0.0012	0.0001	0.0010
		200	0.0009	0.0004	0.0004	0.0001	0.0011	0.0004
		300	0.0014	0.0010	0.0011	0.0010	0.0015	0.0010

The empirical *ARB* results in the Table 2.7 and the theoretical *ARB* results, to first degree of approximation, in the Table 2.8 indicate that the transformed ratio estimators do not produce major reductions in *ARB* as compared to the ratio estimator when sample size is large. Some reduction is observed for small sample size when using transformations where the additive parameter (d) is the *kurtosis*.

Table 2.8: Theoretical ARB to 1^{st} order approximation for the RRT mean estimator, the ratio estimator and for the transformed ratio estimators.

Population			Theoretical ARB (1^{st} Order)				
N	ρ_{XY}	n	$\hat{\mu}_R$	$\hat{\mu}_{TR1}$	$\hat{\mu}_{TR2}$	$\hat{\mu}_{TR3}$	$\hat{\mu}_{TRA}$
1000	0.3209	20	0.0248	0.0310	0.0017	0.0031	0.0267
		50	0.0096	0.0120	0.0006	0.0012	0.0103
		100	0.0046	0.0057	0.0003	0.0006	0.0049
		150	0.0029	0.0036	0.0002	0.0004	0.0031
		200	0.0020	0.0025	0.0001	0.0003	0.0022
		300	0.0012	0.0015	0.0001	0.0001	0.0013
	0.6746	20	0.0124	0.0116	0.0108	0.0031	0.0121
		50	0.0048	0.0045	0.0042	0.0012	0.0047
		100	0.0023	0.0021	0.0020	0.0006	0.0022
		150	0.0014	0.0013	0.0012	0.0004	0.0014
		200	0.0010	0.0009	0.0009	0.0003	0.0010
		300	0.0006	0.0006	0.0005	0.0001	0.0006
	0.8684	20	0.0123	0.0118	0.0108	0.0020	0.0121
		50	0.0048	0.0046	0.0042	0.0008	0.0047
		100	0.0023	0.0022	0.0020	0.0004	0.0022
		150	0.0014	0.0014	0.0013	0.0002	0.0014
		200	0.0010	0.0010	0.0009	0.0002	0.0010
		300	0.0006	0.0006	0.0005	0.0001	0.0006

Table 2.9 presents the results for the empirical MSE estimates and for the theoretical estimates, correct up to first order of approximation. Both results indicate that modest gains can be achieved by using transformations where the additive parameter (d) is the coefficient of *skewness*.

Table 2.9: Empirical *MSE* and theoretical (bold) *MSE* to 1st order of approximation for the RRT mean estimator, the ratio estimator and for the transformed ratio estimators.

Population			MSE				
<i>N</i>	ρ_{XY}	<i>n</i>	$\hat{\mu}_R$	$\hat{\mu}_{TR1}$	$\hat{\mu}_{TR2}$	$\hat{\mu}_{TR3}$	$\hat{\mu}_{TR4}$
1000	0.3209	20	0.5799 0.4496	0.6584 0.4672	0.4097 0.3994	0.4686 0.4663	0.6010 0.4548
		50	0.1881 0.1743	0.2000 0.1811	0.1546 0.1549	0.1793 0.1808	0.1915 0.1763
		100	0.0872 0.0826	0.0914 0.0858	0.0750 0.0734	0.0875 0.0857	0.0884 0.0835
		150	0.0546 0.0520	0.0571 0.0540	0.0475 0.0462	0.0551 0.0539	0.0554 0.0526
		200	0.0395 0.0367	0.0412 0.0381	0.0343 0.0326	0.0394 0.0381	0.0400 0.0371
		300	0.0223 0.0214	0.0232 0.0222	0.0196 0.0190	0.0227 0.0222	0.0225 0.0217
		20	0.3226 0.2939	0.3197 0.2884	0.3956 0.3885	0.5167 0.5162	0.3216 0.2961
	50	0.1165 0.1140	0.1148 0.1118	0.1497 0.1507	0.1979 0.2002	0.1159 0.1132	
	100	0.0558 0.0540	0.0548 0.0530	0.0728 0.0714	0.0967 0.0948	0.0554 0.0536	
	150	0.0352 0.0340	0.0346 0.0333	0.0460 0.0449	0.0608 0.0597	0.0350 0.0338	
	200	0.0254 0.0240	0.0250 0.0235	0.0330 0.0317	0.0434 0.0421	0.0253 0.0238	
	300	0.0145 0.0140	0.0142 0.0137	0.0189 0.0185	0.0250 0.0246	0.0144 0.0139	
	0.8684	20	0.1117 0.0984	0.1083 0.0947	0.1973 0.1920	0.3119 0.3113	0.1106 0.0971
	50	0.0396 0.0381	0.0382 0.0367	0.0743 0.0744	0.1195 0.1207	0.0391 0.0377	
100	0.0188 0.0181	0.0181 0.0174	0.0361 0.0353	0.0584 0.0572	0.0186 0.0178		
150	0.0118 0.0114	0.0114 0.0110	0.0228 0.0222	0.0367 0.0360	0.0117 0.0112		
200	0.0086 0.0080	0.0083 0.0077	0.0164 0.0157	0.0263 0.0254	0.0085 0.0079		
300	0.0049 0.0047	0.0047 0.0045	0.0094 0.0091	0.0151 0.0148	0.0048 0.0046		

Table 2.10 gives the *PRE* of various transformed ratio estimators relative to the ratio estimator based on first order approximation.

We can observe that the transformed ratio estimators that utilize the parameter *d* as coefficient of *skewness* result in higher *PRE* as compared to the ratio estimator when the correlation is larger. This was expected based on Condition (2.16) and Table 2.11 below.

Table 2.10: *PRE* for the transformed ratio estimator related to the ratio estimator based on 1st order of approximation.

Population		<i>PRE</i> (1 st Order)			
N	ρ_{XY}	$\hat{\mu}_{TR1}$	$\hat{\mu}_{TR2}$	$\hat{\mu}_{TR3}$	$\hat{\mu}_{TRA}$
	0.3209	96.24	112.56	96.41	98.86
1000	0.6746	101.93	75.65	56.94	100.64
	0.8684	103.87	51.24	31.60	101.28

Note that the transformed ratio estimator performs better than the ratio estimator when the condition in (2.16) is satisfied.

Table 2.11: Calculations for the expression in (2.16).

Population		Condition (<i>MSE</i> - 1 st Order)			
N	ρ_{XY}	$\hat{\mu}_{TR1}$	$\hat{\mu}_{TR2}$	$\hat{\mu}_{TR3}$	$\hat{\mu}_{TRA}$
	0.3209	-0.0299	0.0855	-0.0285	-0.0088
1000	0.6746	0.0127	-0.2160	-0.5074	0.0043
	0.8684	0.0108	-0.2751	-0.6259	0.0037

2.7 Conclusions

We can observe from this study that the estimation of the mean of a sensitive variable can be improved by using a non-sensitive auxiliary variable. The ratio estimators, in spite of being biased, can have much better *PRE* as compared to the RRT mean estimator. Our simulation study and the numerical example show that this improvement can be quite substantial if the correlation between the study variable and the auxiliary variable is high. We also note that there is hardly any difference in the *Bias* or *MSE* of the proposed estimator when using first or second order approximation. It is further noticed that the transformed ratio estimators produce very minimal gain over the ordinary ratio estimator.

References

- CHANDRA, P. & SINGH, H. P. 2005. A family of estimators for population variance using knowledge of kurtosis of an auxiliary variable in sample surveys. *Statistics in Transition*, 7(1), 27-34.
- EICHHORN, B. H. & HAYRE, L. S. 1983. Scrambled randomized response methods for obtaining sensitive quantitative data. *Journal of Statistical Planning and Inference*, 7, 307-316.
- EUROSTAT. 2008. NACE Rev. 2 - Statistical classification of economic activities in the

- European Community. *Official Publications of the European Communities*, 112-285 and 306-311.
- GUPTA, S. N., GUPTA, B. C. & SINGH, S. 2002. Estimation of sensitivity level of personal interview survey questions. *Journal of Statistical Planning and Inference*, 100, 239-247.
- GUPTA, S. & SHABBIR, J. 2004. Sensitivity estimation for personal interview survey questions. *Statistica*, 64, 643-653.
- GUPTA, S., SHABBIR, J. & SEHRA, S. 2010. Mean and sensitivity estimation in optional randomized response models. *Journal of Statistical Planning and Inference*, 140(10), 2870-2874.
- KADILAR, C. & CINGI, H. 2006. Improvement in estimating the population mean in simple random sampling. *Applied Mathematics Letters*, 19(1), 75-79.
- KOYUNCU, N. & KADILAR, C. 2009. Efficient estimators for the population mean, *Hacettepe Journal of Mathematics and Statistics*, 38(2), 217-225.
- KULKARNI, S. P. 1977. A note on modified ratio estimator using transformation, *Journal of the Indian Society of Agricultural Statistics*, 30(2), 125-128.
- SAHA, A. 2008. A randomized response technique for quantitative data under unequal probability sampling. *Journal of statistical Theory and Practice*, 2(4), 589-596.
- SHABBIR, J. & GUPTA, S. 2010. Estimation of the finite population mean in two-phase sampling when auxiliary variables are attribute, *Hacettepe Journal of Mathematics and Statistics*, 39(1), 121-129.
- SINGH, J., PANDAY, B. N. & HIRANO, K. 1973. On the utilization of a known coefficient of kurtosis in the estimation procedure of variance. *Annals of the Institute of Statistical Mathematics*, 25, 51-55.
- SINGH, H. P. & VISHWAKARMA, G. 2008. Some families of estimators of variance of stratified random sample mean using auxiliary information. *Journal of Statistical Theory and Practice*, 2(1), 21-43.
- SISODIA, B. V. S. & DWIDEDI, V. K. 1981. A modified ratio estimator using coefficient of variation of auxiliary variable. *Journal of the Indian Society of Agricultural Statistics*, 33, 13-18.
- SMILHILY, M. & STORM, H. 2010. ICT usage in enterprises - 2009, *Eurostat Publications*, Issue 1.
- SUKHATME, P.V. & SUKHATME, B.V. 1984. *Sampling theory of surveys with applications*, 3rd Ed., Ames, Iowa, Iowa State University Press.
- TURGUT, Y. & CINGI, H. 2008. New generalized estimators for the population variance

using auxiliary information. *Hacettepe Journal of Mathematics and Statistics*, 37(2), 177-184.

UPADHYAYA, L. N. & SINGH, H. P. 1999. Use of transformed auxiliary variable in estimating the finite population mean. *Biometrical Journal*, 41(5), 627-636.

UPADHYAYA L. N., SINGH, G. N. & SINGH, H. P. 2000. Use of transformed auxiliary variable in the estimation of population ratio in sample surveys. *Statistics in Transition*, 4(6), 1019-1027.

Appendix A - R Routines

Listing 2.1: R Code for Simulation Study of Proposed Estimator in Chapter 2

```

1
2 proj1 <- function(N, sigma, mu)
3 {
4   set.seed(100)
5   #Generation of a bivariate normal population
6   data_yx <- mvrnorm(N, mu, sigma)
7
8   #Study variable
9   Y <- data_yx[,1]
10  #Auxiliary variable, correlated with Y
11  X <- data_yx[,2]
12
13  #Coefficient of correlation between Y and X
14  Ro_YX <- cor(Y,X)
15
16  #Scrambling variable independent of Y and X, with mean=0
17  S <- rnorm(N, mean=0, sd=0.1*sd(X))
18  #Scrambled response
19  Z <- Y+S
20
21  #Coefficient of correlation between Z and X
22  Ro_ZX <- Ro_YX/sqrt(1+(var(S)/var(Y)))
23
24  #population
25  univ <- data.frame(cbind(Y=Y, S=S, Z=Z, X=X, NRAND=runif(N)))
26  univ <- univ[order(univ$NRAND),]
27
28  #Mean of Y
29  my <- mean(univ$Y)
30  mz <- mean(univ$Z)
31  mx <- mean(univ$X)
32  ms <- mean(univ$S)
33
34  #Sample dimension
35  dim_samp <- c(20, 50, 100, 200, 300)
36
37  res <- NULL
38  for (i in 1:length(dim_samp))
39  {
40    #sample dimension
41    n <- dim_samp[i]
42    #sample
43    samp <- univ[1:n,]
44    #sampling rate
45    f <- n/N
46

```

2. RATIO ESTIMATION OF THE MEAN OF A SENSITIVE VARIABLE IN THE PRESENCE OF AUXILIARY INFORMATION
Appendix A - R Routines

```

47  #estimators
48  est1 <- mean(samp$Z)
49  est2 <- mean(samp$Z) * (mean(univ$X) / mean(samp$X))
50
51  #Ratio
52  R <- mean(univ$X) / mean(samp$X)
53
54  #Mean Square Error of 1st estimator
55  mse1 <- ((1-f)/n) * var(univ$Z)
56
57  #Coefficient of variation
58  c_x <- sd(univ$X) / mx
59  c_y <- sd(univ$Y) / my
60  c2_x <- c_x^2
61  c2_y <- c_y^2
62  c2_z <- c2_y + (var(univ$S) / (my^2))
63  c_z <- sqrt(c2_z)
64
65  #Bias of ratio estimator - 1st degree approximation
66  bias2i <- ((1-f)/n) * my * (c2_x - Ro_ZX * c_z * c_x)
67  #Bias of ratio estimator - 2nd degree approximation
68  bias2ii <- bias2i * (1 + ((1-f)/n) * 3 * c2_x)
69
70  #Mean Square Error of ratio estimator - 1st degree approximation
71  mse2i <- ((1-f)/n) * (my^2) * (c2_z + c2_x - 2 * Ro_ZX * c_z * c_x)
72  #Mean Square Error of ratio estimator - 2nd degree approximation
73  mse2ii <- mse2i + 3 * (my^2) * (((1-f)/n)^2) * c2_x * ((1 + 2 * (Ro_ZX^2)) * c2_z
74  + 3 * c2_x - 6 * Ro_ZX * c_z * c_x)
75
76  aux_bias <- (c_x - Ro_ZX * c_z)
77
78  aux_mse1 <- (Ro_ZX - (1/2) * (c_x / c_z))
79  aux_mse2 <- 2 * Ro_ZX * (c_z / c_x) - 3 * ((1-f)/n) * ((1 + 2 * (Ro_ZX^2)) * c2_z
80  + 3 * c2_x - 6 * Ro_ZX * c_z * c_x)
81
82  emp <- NULL
83
84  #Empirical results
85  #Simulation of 5000 replicas of estimates
86  ...
87
88  #Results
89  res <- rbind(res, c(N, n, Ro_YX, Ro_ZX, R, my, mz, ms,
90  med_est1, med_est2, bias2i, bias2ii, emp_mse1, mse1,
91  emp_mse2, mse2i, mse2ii, aux_bias, aux_mse1, aux_mse2))
92 }
93 colnames(res) <- c("N", "n", "RhoXY", "RhoZX", "R", "mY", "mZ", "mS",
94 "Est1", "Est2", "BIAS2I", "BIAS2II", "EMP_MSE1", "MSE1",
95 "EMP_MSE2", "MSE2I", "MSE2II", "AUX_BIAS", "AUX_MSE1", "AUX_MSE2")
96 return(res)

```



```
97 }
98
99 #Package for generation
100 require(MASS)
101 N<-1000
102 #Parameters
103 sigma1 <- matrix(c(9,1.9,1.9,4),2,2)
104 sigma2 <- matrix(c(10,3,3,2),2,2)
105 sigma3 <- matrix(c(6,3,3,2),2,2)
106 mu <- c(2,2)
107
108 res <- NULL
109 for (i in 1:length(N))
110 {
111   res <- rbind(res,proj1(N[i],sigma1,mu))
112   res <- rbind(res,proj1(N[i],sigma2,mu))
113   res <- rbind(res,proj1(N[i],sigma3,mu))
114 }
115 write.table(res,"chapter2_ss_results1.txt",sep="\t",dec="," , row.names=FALSE)
```

Listing 2.2: R Code for Simulation Study of Transformed Ratio Estimators in Chapter 2

```
1
2 mykurtosis <- function(x)
3 {
4   m4 <- mean((x-mean(x))^4)
5   kurt <- m4/(sd(x)^4)
6   return(kurt)
7 }
8 myskewness <- function(x)
9 {
10  m3 <- mean((x-mean(x))^3)
11  skew <- m3/(sd(x)^3)
12  return(skew)
13 }
14 projl_transf <- function(N,sigma,mu)
15 {
16
17   #Generation of a bivariate normal population
18   data_yx <- mvrnorm(N, mu, sigma)
19
20   #Study variable
21   Y <- data_yx[,1]
22   #Auxiliary variable, correlated with Y
23   X <- data_yx[,2]
24
25   #Coefficient of correlation between Y and X
26   Ro_YX <- cor(Y,X)
27
28   #Scrambling variable independent of Y and X, with mean=0
29   S <- rnorm(N,mean=0,sd=0.1*sd(X))
30   #Scrambled response
31   Z <- Y+S
32
33   #Coefficient of correlation between Z and X
34   Ro_ZX <- Ro_YX/sqrt(1+(var(S)/var(Y)))
35
36   #population
37   univ <- data.frame(cbind(Y=Y,S=S,Z=Z,X=X,NRAND=runif(N)))
38   univ <- univ[order(univ$NRAND),]
39
40   #Mean of Y
41   my <- mean(univ$Y)
42   mz <- mean(univ$Z)
43   mx <- mean(univ$X)
44
45   #Samples dimension
46   dim_samp <- c(20,50,100,150,200,300)
47
```

```

48 res <- NULL
49 for (i in 1:length(dim_samp))
50 {
51   #sample dimension
52   n <- dim_samp[i]
53   #sample
54   samp <- univ[1:n,]
55   #sampling rate
56   f <- n/N
57
58   #Ratio
59   R <- mean(univ$X)/mean(samp$X)
60
61   #Ordinary mean
62   est1 <- mean(samp$Z)
63   #Ratio estimator
64   est2 <- mean(samp$Z) * (mean(univ$X)/mean(samp$X))
65
66   #Coefficient of variation
67   c_x <- sd(univ$X)/mx
68   c_y <- sd(univ$Y)/my
69   c2_x <- c_x^2
70   c2_y <- c_y^2
71   c2_z <- c2_y + (var(univ$S) / (my^2))
72   c_z <- sqrt(c2_z)
73
74   #Bias of ratio estimator - 1st degree approximation
75   bias2i <- ((1-f)/n) * my * (c2_x - Ro_ZX * c_z * c_x)
76   #Bias of ratio estimator - 2nd degree approximation
77   bias2ii <- bias2i * (1 + ((1-f)/n) * 3 * c2_x)
78
79   #Mean Square Error of 1st estimator (ordinal mean)
80   mse1 <- ((1-f)/n) * (var(univ$Y) + var(univ$S))
81
82   #Mean Square Error of ratio estimator - 1st degree approximation
83   mse2i <- ((1-f)/n) * (my^2) * (c2_z + c2_x - 2 * Ro_ZX * c_z * c_x)
84   #Mean Square Error of ratio estimator - 2nd degree approximation
85   mse2ii <- mse2i + 3 * (my^2) * (((1-f)/n)^2) * c2_x * ((1 + 2 * (Ro_ZX^2)) * c2_z
86     + 3 * c2_x - 6 * Ro_ZX * c_z * c_x)
87
88   nu <- 1
89   aux_m <- c2_x - 2 * Ro_ZX * c_z * c_x
90
91   s <- myskewness(univ$X)
92   k <- mykurtosis(univ$X)
93
94   vc <- c(1, 1, 1, s, k)
95   vd <- c(s, k, Ro_YX, k, s)
96
97   #Initialize the variables est3, mse3i, ...

```

2. RATIO ESTIMATION OF THE MEAN OF A SENSITIVE VARIABLE IN THE PRESENCE OF AUXILIARY INFORMATION
Appendix A - R Routines

```

98
99   for (i in 1:length(vc))
100   {
101     nu <- (vc[i]*mean(univ$X))/(vc[i]*mean(univ$X)+vd[i])
102     vnu <- c(vnu,nu)
103
104     aux_bias1 <- (nu-1)*(Ro_ZX-(nu+1)*c_x/c_z)
105     aux_mse1 <- (nu-1)*(Ro_ZX-(nu+1)*c_x/(2*c_z))
106
107     vb1 <- c(vb1,aux_bias1)
108     vm1 <- c(vm1,aux_mse1)
109
110     #Transformed ratio estimator
111     est3 <- c(est3,mean(samp$Z)*(vc[i]*mean(univ$X)+vd[i])
112               / (vc[i]*mean(samp$X)+vd[i]))
113
114     #Mean Square Error of transformed ratio estimator
115     #1st degree approximation
116     mse3i <- c(mse3i,((1-f)/n)*(my^2)*(c2_z+(nu^2)*c2_x-2*nu*Ro_ZX*c_z*c_x))
117
118     #Mean Square Error of transformed ratio estimator
119     #1st degree approximation
120     mse3ii <- c(mse3ii,mse3i[i]+3*(my^2)
121                *(((1-f)/n)^2)*c2_x*((nu^2)*(1+2*(Ro_ZX^2))
122                *c2_z+3*(nu^4)*c2_x-6*(nu^3)*Ro_ZX*c_z*c_x))
123
124     #Bias of transformed ratio estimator - 1st degree approximation
125     bias3i <- c(bias3i,((1-f)/n)*my*((nu^2)*c2_x-nu*Ro_ZX*c_z*c_x))
126     #Bias of transformed ratio estimator - 2nd degree approximation
127     bias3ii <- c(bias3ii,bias3i[i]
128                 +(((1-f)/n)^2)*3*my*((nu^4)*(c2_x^2)-(nu^3)
129                 *Ro_ZX*c_z*(c_x^3)))
130   }
131
132   #Empirical results
133   #Simulation of 5000 replicas of estimates
134   ...
135
136   #Results
137   res <- rbind(res,c(N,n,Ro_YX,Ro_ZX,R,my,med_est1,med_est2,
138                  med_est3,
139                  bias2i,bias2ii,
140                  bias3i,
141                  bias3ii,
142                  emp_mse1,mse1,emp_mse2,mse2i,mse2ii,
143                  emp_mse3,
144                  mse3i,
145                  mse3ii,
146                  vnu,
147                  vb1,

```

```

148         vm1))
149     }
150     colnames(res) <- c("N", "n", "RhoXY", "RhoZX", "R", "mY", "Est1", "Est2",
151         paste("Est3_", 1:length(vc), sep=""),
152         "BIAS2I", "BIAS2II",
153         paste("BIAS3I_", 1:length(vc), sep=""),
154         paste("BIAS3II_", 1:length(vc), sep=""),
155         "EMP_MSE1", "MSE1", "EMP_MSE2", "MSE2I", "MSE2II",
156         paste("EMP_MSE3_", 1:length(vc), sep=""),
157         paste("MSE3I_", 1:length(vc), sep=""),
158         paste("MSE3II_", 1:length(vc), sep=""),
159         paste("NU_", 1:length(vc), sep=""),
160         paste("AUX3_BIAS1_", 1:length(vc), sep=""),
161         paste("AUX3_MSE1_", 1:length(vc), sep=""))
162     return(res)
163 }
164 #Package for generation
165 require(MASS)
166 N <- 1000
167
168 #Parameters
169 sigma1 <- matrix(c(9, 1.9, 1.9, 4), 2, 2)
170 sigma2 <- matrix(c(10, 3, 3, 2), 2, 2)
171 sigma3 <- matrix(c(6, 3, 3, 2), 2, 2)
172 mu <- c(2, 2)
173
174 res <- NULL
175 for (i in 1:length(N))
176 {
177     res <- rbind(res, proj1_transf(N[i], sigma1, mu))
178     res <- rbind(res, proj1_transf(N[i], sigma2, mu))
179     res <- rbind(res, proj1_transf(N[i], sigma3, mu))
180 }
181 write.table(res, "chapter2_ss_results2.txt", sep="\t", dec=".", row.names=FALSE)

```

Listing 2.3: R Code for Numerical Example of Proposed Estimator in Chapter 2

```
1
2 proj1_real <- function(Y,X,N)
3 {
4
5   #Coefficient of correlation between Y and X
6   Ro_YX <- cor(Y,X)
7
8   #Scrambling variable independent of Y and X, with mean=0
9   S <- rnorm(N,mean=0,sd=sd(X)*0.1)
10  #Scrambled response
11  Z <- Y+S
12
13  #Coefficient of correlation between Z and X
14  Ro_ZX <- Ro_YX/sqrt(1+(var(S)/var(Y)))
15
16  #population
17  univ <- data.frame(cbind(Y=Y,S=S,Z=Z,X=X,NRAND=runif(N)))
18  univ <- univ[order(univ$NRAND),]
19
20  #Mean of Y
21  my <- mean(univ$Y)
22  mz <- mean(univ$Z)
23  mx <- mean(univ$X)
24
25  #Samples dimension
26  dim_samp <- c(100,200,300,400,500,1000,1500,2000)
27
28  res <- NULL
29  for (i in 1:length(dim_samp))
30  {
31    #sample dimension
32    n <- dim_samp[i]
33    #sample
34    samp <- univ[1:n,]
35    #Sampling rate
36    f <- n/N
37
38    #estimators
39    est1 <- mean(samp$Z)
40    est2 <- mean(samp$Z) * (mean(univ$X)/mean(samp$X))
41
42    #Ratio
43    R <- mean(univ$X)/mean(samp$X)
44
45    #Mean Square Error of 1st estimator
46    mse1 <- ((1-f)/n) * (var(univ$Y)+var(univ$S))
47
```

```

48   #Coefficient of variation
49   c_x <- sd(univ$X)/mx
50   c_y <- sd(univ$Y)/my
51   c2_x <- c_x^2
52   c2_y <- c_y^2
53   c2_z <- c2_y+(var(univ$S)/(my^2))
54   c_z <- sqrt(c2_z)
55
56   #Bias of ratio estimator - 1st degree approximation
57   bias2i <- ((1-f)/n)*my*(c2_x-Ro_ZX*c_z*c_x)
58   #Bias of ratio estimator - 2nd degree approximation
59   bias2ii <- bias2i*(1+((1-f)/n)*3*c2_x)
60
61   #Mean Square Error of ratio estimator - 1st degree approximation
62   mse2i <- ((1-f)/n)*(my^2)*(c2_z+c2_x-2*Ro_ZX*c_z*c_x)
63   #Mean Square Error of ratio estimator - 2nd degree approximation
64   mse2ii <- mse2i+3*(my^2)*(((1-f)/n)^2)*c2_x*((1+2*(Ro_ZX^2))*c2_z
65     +3*c2_x-6*Ro_ZX*c_z*c_x)
66
67   aux_bias <- (c_x-Ro_ZX*c_z)
68
69   aux_mse1 <- (Ro_ZX-(1/2))*(c_x/c_z)
70   aux_mse2 <- 2*Ro_ZX*(c_z/c_x)-3*((1-f)/n)*((1+2*(Ro_ZX^2))*c2_z
71     +3*c2_x-6*Ro_ZX*c_z*c_x)
72
73   #Empirical results
74   #Simulation of 5000 replicas of estimates
75   ...
76
77   #Results
78   res <- rbind(res,c(N,n,Ro_YX,Ro_ZX,R,my,med_est1,med_est2,
79     bias2i,bias2ii,emp_mse1,mse1,emp_mse2,mse2i,mse2ii,
80     aux_bias,aux_mse1,aux_mse2))
81 }
82 colnames(res) <- c("N","n","RhoXY","RhoZX","R","mY","Est1","Est2",
83   "BIAS2I","BIAS2II","EMP_MSE1","MSE1","EMP_MSE2",
84   "MSE2I","MSE2II","AUX_BIAS","AUX_MSE1","AUX_MSE2")
85 return(res)
86 }
87
88 #Package for generation
89 require(MASS)
90
91 #Import data
92 data_yx <- read.table("IUTICE09.dat",sep="\t",dec=".",header = T)
93 #Study variable (purchase, millions of euros)
94 Y <- data_yx[,3]
95 #Auxiliary variable, correlated with Y (turnover, millions of euros)
96 X <- data_yx[,2]
97

```

2. RATIO ESTIMATION OF THE MEAN OF A SENSITIVE VARIABLE IN THE PRESENCE OF AUXILIARY INFORMATION

Appendix A - R Routines

```
98 #Data application
99 N <- dim(data_yx)[1]
100 res <- proj1_real(Y,X,N)
101
102 #Export data
103 write.table(res, "chapter2_ne_results", sep="\t", dec=", ", row.names=FALSE)
```




Estimation of the Mean of a Sensitive Variable in the Presence of Auxiliary Information

Abstract

Sousa et al. (2010) introduced a ratio estimator for the mean of a sensitive variable and showed that this estimator performs better than the ordinary mean estimator based on a Randomized Response Technique (RRT). In this paper, we introduce a regression estimator that performs better than the ratio estimator even for modest correlation between the primary and the auxiliary variables. The underlying assumption is that the primary variable is sensitive in nature but a non-sensitive auxiliary variable exists that is positively correlated with the primary variable. Expressions for the *Bias* and Mean Square Error (*MSE*) are derived based on the first order of approximation. It is shown that the proposed regression estimator performs better than the ratio estimator and the ordinary RRT mean estimator (that does not utilize the auxiliary information). We also consider a generalized regression-cum-ratio estimator that has even smaller *MSE*. An extensive simulation study is presented to evaluate the performances of the proposed estimators in relation to other estimators in the study.

Published as: GUPTA, S., SHABBIR, J., SOUSA, R. & REAL, P.C. 2012. Estimation of the Mean of a Sensitive Variable in the Presence of Auxiliary Information. *Communications in Statistics - Theory and Methods*, 41(13-14), 2394-2404.

The procedure is also applied to some financial data: purchase orders (a sensitive variable) and gross turnover (a non-sensitive variable) in 2009 for a population of 5336 companies in Portugal from a survey on Information and Communication Technologies (ICT) usage.

3.1 Introduction

In survey research, direct reliable observation on the variable of interest (Y) is sometimes not possible because the variable may be sensitive in nature. In this paper we focus on estimating the mean of a sensitive variable Y using an auxiliary variable (X) that can be directly observed and that is correlated with the variable of the interest. For example, Y may be the total number of abortions a woman of child bearing age might have had and X may be her current age. Similarly, Y may be the total value of purchase orders in a year for a company and X may be the total turnover for that company in that year. In such situations, mean of Y can be estimated by using one of many randomized response techniques if the auxiliary information is to be ignored.

Many authors have estimated the mean of a sensitive variable when the primary variable is sensitive and there is no auxiliary variable available. These include Eichhorn and Hayre (1983), Gupta and Shabbir (2004), Gupta et al. (2002, 2010), Wu et al. (2008), Saha (2008) and Perri (2008). Also, many authors have presented ratio and regression estimators when both Y and X are directly observable. These include Kadilar and Cingi (2005), Kadilar et al. (2007), Shabbir and Gupta (2007, 2010) and Nangsue (2009).

In this paper, we propose a regression estimator where the RRT estimator of the mean of Y is further improved by using an auxiliary variable X . We also consider a generalized regression-cum-ratio estimator under the same conditions. Expressions for the *Bias* and *MSE* for the proposed estimators are derived, correct up to first order of approximation. We compare the performances of the proposed estimators with those of the ratio and the ordinary RRT mean estimators. We observe that there is considerable reduction in *MSE*, particularly when the correlation between the study variable and the auxiliary variable is high.

3.2 Terminology

Let Y be the study variable, a sensitive variable which cannot be observed directly due to respondent bias. Let X be a non-sensitive auxiliary variable which has a positive correlation with Y . Let S be a scrambling variable independent of Y and X . The respondent is asked to report a scrambled response for Y given by $Z = Y + S$ but is asked to provide a true response for X . Let a random sample of size n be drawn without replacement from

a finite population $U = (U_1, U_2, \dots, U_N)$. For the i^{th} unit ($i = 1, 2, \dots, N$), let y_i and x_i , respectively, be the values of the study variable Y and auxiliary variable X . Let $\bar{y} = \frac{\sum_{i=1}^n y_i}{n}$, $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$ and $\bar{z} = \frac{\sum_{i=1}^n z_i}{n}$ be the sample means and $\bar{Y} = E(Y)$, $\bar{X} = E(X)$, and $\bar{Z} = E(Z)$ be the corresponding population means for Y , X and Z , respectively. We assume that \bar{X} is known and $\bar{S} = E(S) = 0$. Thus, $E(Z) = E(Y)$ and $C_z^2 = C_y^2 + \frac{S_s^2}{\bar{Y}^2}$, where C_z and C_y are the coefficients of variation of z and y , respectively. If $e_0 = \frac{\bar{z} - \bar{Z}}{\bar{Z}}$, $e_1 = \frac{\bar{x} - \bar{X}}{\bar{X}}$, $e_2 = \frac{s_x^2 - S_x^2}{S_x^2}$, and $e_3 = \frac{s_{zx} - S_{zx}}{S_{zx}}$, then we have $E(e_i) = 0$, $i = 0, 1, 2, 3$.

If information on X is ignored, then an unbiased estimator of μ_Y is the ordinary sample mean (\bar{z}) given by (3.1) below.

$$\hat{\mu}_Y = \bar{z} \quad (3.1)$$

The mean square error (MSE) of $\hat{\mu}_Y$ is given by

$$MSE(\hat{\mu}_Y) = \frac{1-f}{N} (S_y^2 + S_s^2), \quad (3.2)$$

where $f = n/N$, $S_y^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{Y})^2$ and $S_s^2 = \frac{1}{N-1} \sum_{i=1}^N (s_i - \bar{S})^2$.

3.3 The Ratio Estimator

Sousa et al. (2010) proposed a ratio estimator for the mean of sensitive variable (Y) utilizing information from a non-sensitive auxiliary variable (X). This estimator is given by

$$\hat{\mu}_R = \bar{z} \left(\frac{\bar{X}}{\bar{x}} \right). \quad (3.3)$$

Bias and *MSE* of $\hat{\mu}_R$, correct up to first order of approximation, are given by

$$Bias(\hat{\mu}_R) \cong \left(\frac{1-f}{n} \right) \bar{Y} (C_x^2 - \rho_{zx} C_z C_x) \quad (3.4)$$

and

$$MSE(\hat{\mu}_R) \cong \left(\frac{1-f}{n} \right) \bar{Y}^2 (C_z^2 + C_x^2 - 2\rho_{zx} C_z C_x). \quad (3.5)$$

It can be observed that $MSE(\hat{\mu}_R) < MSE(\hat{\mu}_Y)$ if

$$\rho_{yx} > \frac{1}{2} \frac{C_x}{C_y} \sqrt{1 + \frac{S_s^2}{S_y^2}}. \quad (3.6)$$

3.4 Ordinary Regression Estimator

Assuming linear relationship between Y and X , we propose the following regression estimator for the population mean of Y

$$\hat{\mu}_{Reg} = \bar{z} + \hat{\beta}_{zx} (\bar{X} - \bar{x}), \quad (3.7)$$

where $\hat{\beta}_{zx}$ is the sample regression coefficient between Z and X and $Z = Y + S$ is the scrambled response on Y . Using Taylor's approximation and retaining terms of order up to 2, (3.7) can be rewritten as

$$\hat{\mu}_{Reg} - \bar{Z} \cong \bar{Z}e_0 - \beta_{zx}\bar{X} [e_1 + e_1e_3 - e_1e_2]. \quad (3.8)$$

From Mukhopadhyay (1998, p. 123), we have:

$$E(e_0^2) = \frac{1-f}{n} C_z^2, E(e_1^2) = \frac{1-f}{n} C_x^2, E(e_{12}) = \frac{1-f}{n} \frac{1}{\bar{X}} \frac{\mu_{03}}{\mu_{02}}, E(e_{13}) = \frac{1-f}{n} \frac{1}{\bar{X}} \frac{\mu_{12}}{\mu_{11}},$$

where $\mu_{rs} = \frac{1}{N-1} \sum_{i=1}^N (z_i - \bar{Z})^r (x_i - \bar{X})^s$ and C_x, C_z are the coefficients of variation of x and z , respectively.

$$\text{Also we have: } \beta_{zx} = \frac{S_{zx}}{S_x^2} = \frac{S_{yx}}{S_x^2} = \rho_{yx} \frac{S_y}{S_x} = \beta_{yx}, \rho_{zx} = \frac{\rho_{yx}}{\sqrt{1 + \frac{S_s^2}{S_y^2}}},$$

where ρ_{yx} and ρ_{zx} are the coefficients of correlation between y and x , and between z and x , respectively.

Recognizing that $\bar{Z} = \bar{Y}$ in Equation (3.8), the *Bias* and *MSE* of $\hat{\mu}_{Reg}$, to first order of approximation, are given by

$$\text{Bias}(\hat{\mu}_{Reg}) \cong -\beta_{zx} \left(\frac{1-f}{n} \right) \left\{ \frac{\mu_{12}}{\mu_{11}} - \frac{\mu_{03}}{\mu_{02}} \right\} \quad (3.9)$$

and

$$\text{MSE}(\hat{\mu}_{Reg}) \cong \left(\frac{1-f}{n} \right) \bar{Y}^2 C_z^2 (1 - \rho_{zx}^2) = \left(\frac{1-f}{n} \right) S_y^2 \left\{ \left(1 + \frac{S_s^2}{S_y^2} \right) - \rho_{yx}^2 \right\}. \quad (3.10)$$

It can be verified easily that

(i) $\text{MSE}(\hat{\mu}_{Reg}) < \text{MSE}(\hat{\mu}_Y)$ if

$$\rho_{yx}^2 > 0; \quad (3.11)$$

(ii) $\text{MSE}(\hat{\mu}_{Reg}) < \text{MSE}(\hat{\mu}_R)$ if

$$(C_x - C_z \rho_{zx})^2 > 0. \quad (3.12)$$

These conditions will always hold true indicating that up to first order of approximation, the regression estimator performs better than $\hat{\mu}_Y$ and $\hat{\mu}_R$.

3.5 Generalized Regression-cum-ratio Estimator

Many authors have used regression-cum-ratio estimators that combine both the regression estimator and ratio estimator. These include Ray and Singh (1981), Perri (2004), and Kadilar and Cingi (2004, 2006). We consider a similar hybrid estimator, as a generalized regression-cum-ratio estimator with constant coefficients whose values are to be determined later from optimality considerations. The main idea is to see if further gains can be achieved by using a generalized regression-cum-ratio estimator, as compared to regression estimator given by (3.7). This estimator is given by:

$$\hat{\mu}_{GRR} = [k_1\bar{z} + k_2(\bar{X} - \bar{x})] \left(\frac{\bar{X}}{\bar{x}} \right), \quad (3.13)$$

where k_1 and k_2 are constants.

Solving (3.13) using Taylor's approximation and retaining terms of order up to 2, we have

$$\hat{\mu}_{GRR} - \bar{Y} \cong (k_1 - 1)\bar{Y} + k_1\bar{Y}(e_0 - e_1 - e_0e_1 + e_1^2) - k_2\bar{X}(e_1 - e_1^2). \quad (3.14)$$

From (3.14), the *Bias* and *MSE* of $\hat{\mu}_{GRR}$ to first order of approximation are given by

$$Bias(\hat{\mu}_{GRR}) \cong (k_1 - 1)\bar{Y} + k_1\bar{Y} \left(\frac{1-f}{n} \right) \{C_x^2 - \rho_{zx}C_zC_x\} + k_2\bar{X} \left(\frac{1-f}{n} \right) C_x^2 \quad (3.15)$$

and

$$\begin{aligned} MSE(\hat{\mu}_{GRR}) \cong & (k_1 - 1)^2\bar{Y}^2 + k_1^2\bar{Y}^2 \left(\frac{1-f}{n} \right) \{C_z^2 + 3C_x^2 - 4\rho_{zx}C_zC_x\} \\ & + k_2^2\bar{X}^2 \left(\frac{1-f}{n} \right) C_x^2 - 2k_1\bar{Y}^2 \left(\frac{1-f}{n} \right) \{C_x^2 - \rho_{zx}C_zC_x\} \\ & - 2k_2\bar{Y}\bar{X} \left(\frac{1-f}{n} \right) C_x^2 - 2k_1k_2\bar{Y}\bar{X} \left(\frac{1-f}{n} \right) \{\rho_{zx}C_zC_x - 2C_x^2\}. \end{aligned} \quad (3.16)$$

Differentiating (3.16) with respect to k_1 and k_2 we get the following optimum values:

$$k_{1(opt)} = \frac{1 - \left(\frac{1-f}{n} \right) C_x^2}{1 - \left(\frac{1-f}{n} \right) \{C_x^2 - C_z^2(1 - \rho_{zx}^2)\}} \quad (3.17)$$

and

$$k_{2(opt)} = \frac{\bar{Y}}{\bar{X}} \left\{ 1 + k_{1(opt)} \left(\frac{\rho_{zx}C_z}{C_x} - 2 \right) \right\}. \quad (3.18)$$

which minimize the *MSE*.

Substituting the optimum values of k_1 and k_2 in (3.16), we get

$$MSE(\hat{\mu}_{GRR})_{min} \cong \frac{\bar{Y}^2 C_z^2 (1 - \rho_{zx}^2) \left(\frac{1-f}{n}\right) \left\{1 - \left(\frac{1-f}{n}\right) C_x^2\right\}}{C_z^2 (1 - \rho_{zx}^2) \left(\frac{1-f}{n}\right) + \left\{1 - \left(\frac{1-f}{n}\right) C_x^2\right\}}. \quad (3.19)$$

It can be verified that:

(i) $MSE(\hat{\mu}_{GRR})_{min} < MSE(\hat{\mu}_Y)$ if

$$\left(\frac{1-f}{n}\right) \{S_y^2 + S_s^2\} > 0, \quad (3.20)$$

which is always true.

(ii) $MSE(\hat{\mu}_{GRR})_{min} < MSE(\hat{\mu}_R)$ if

$$\left(\frac{C_x}{C_z} - \rho_{zx}\right)^2 + \frac{\left(\frac{1-f}{n}\right) C_z^2 (1 - \rho_{zx}^2)^2}{\left(\frac{1-f}{n}\right) C_z^2 (1 - \rho_{zx}^2) + \left(1 - \left(\frac{1-f}{n}\right) C_x^2\right)} > 0. \quad (3.21)$$

(iii) $MSE(\hat{\mu}_{GRR})_{min} < MSE(\hat{\mu}_{Reg})$ if

$$\left(\frac{1-f}{n}\right) C_z^2 (1 - \rho_{zx}^2) > 0, \quad (3.22)$$

which is always true.

From these conditions we can conclude that the generalized estimator in (3.13) with optimal coefficients is always better than $\hat{\mu}_Y$, and $\hat{\mu}_{Reg}$. It is also better than the ratio and regression estimators if

$$1 - \left(\frac{1-f}{n}\right) C_x^2 > 0$$

or

$$\left(\frac{1-f}{n}\right) C_x^2 < 1. \quad (3.23)$$

The last condition is very likely to hold true. So, with this generalized regression-cum-ratio estimator, we may be able to achieve further gain in terms of MSE , as can be observed from the simulation results in the next section.

3.6 The Simulation Study

In this section, we present results of a simulation study with particular focus on the performance for the regression estimator $\hat{\mu}_{Reg}$ and the proposed generalized regression-cum-ratio estimator $\hat{\mu}_{GRR}$ as compared to the RRT mean estimator $\hat{\mu}_Y$, and the ratio estimator $\hat{\mu}_R$.

We consider three finite sub-populations of size 1000 each from bivariate normal populations with different covariance matrices to represent the distribution of (Y, X) . The scrambling variable S is taken to be a normal variate with mean equal to zero and standard deviation equal to 10% of the standard deviation of X . The reported response is given by $Z = Y + S$.

All of the simulated populations have theoretical mean of $[Y, X]$ as $\mu = [2, 2]$. The covariance matrices (Σ) are as given below.

Population 1

$$\Sigma = \begin{bmatrix} 9 & 1.9 \\ 1.9 & 4 \end{bmatrix}, \rho_{XY} = 0.3209.$$

Population 2

$$\Sigma = \begin{bmatrix} 10 & 3 \\ 3 & 2 \end{bmatrix}, \rho_{XY} = 0.6746.$$

Population 3

$$\Sigma = \begin{bmatrix} 6 & 3 \\ 3 & 2 \end{bmatrix}, \rho_{XY} = 0.8684.$$

For each population, we consider five sample sizes: $n = 50, 100, 200$ and 300 .

Table 3.1 below gives empirical and theoretical MSE 's for various estimators based on the first order approximation. We estimate the empirical MSE using 5000 samples of various sizes selected from each population. We use the following expression to find the PRE of ratio, regression and generalized regression-cum-ratio estimators as compared to the RRT mean estimator:

$$PRE = \frac{MSE(\hat{\mu}_Y)}{MSE(\hat{\mu}_\alpha)} \times 100,$$

where $\alpha = R, Reg, GRR$.

3. ESTIMATION OF THE MEAN OF A SENSITIVE VARIABLE IN THE PRESENCE OF AUXILIARY INFORMATION

3.6. The Simulation Study

Table 3.1: *MSE* correct up to 1st order approximation and *PRE* for the ratio estimator ($\hat{\mu}_R$), the regression estimator ($\hat{\mu}_{Reg}$) and the generalized regression-cum-ratio estimator ($\hat{\mu}_{GRR}$) relative to the RRT mean estimator.

Population			MSE Estimation				
<i>N</i>	ρ_{XY}	<i>n</i>	Estimator	Empirical	Theoretical	<i>PRE</i>	Condition ¹
1000	0.3209	50	$\hat{\mu}_R$	0.1885	0.1743	98.62	0.0186
			$\hat{\mu}_{Reg}$	0.1557	0.1543	111.43	
			$\hat{\mu}_{GRR}$	0.1571	0.1485	115.77	
		100	$\hat{\mu}_R$	0.0877	0.0826	98.62	0.0088
			$\hat{\mu}_{Reg}$	0.0733	0.0731	111.43	
			$\hat{\mu}_{GRR}$	0.0736	0.0718	113.46	
		200	$\hat{\mu}_R$	0.0382	0.0367	98.62	0.0039
			$\hat{\mu}_{Reg}$	0.0331	0.0325	111.43	
			$\hat{\mu}_{GRR}$	0.0331	0.0322	112.33	
300	$\hat{\mu}_R$	0.0222	0.0214	98.62	0.0023		
	$\hat{\mu}_{Reg}$	0.0194	0.0189	111.43			
	$\hat{\mu}_{GRR}$	0.0194	0.0189	111.95			
1000	0.6746	50	$\hat{\mu}_R$	0.1181	0.1140	167.23	0.0094
			$\hat{\mu}_{Reg}$	0.1041	0.1040	183.20	
			$\hat{\mu}_{GRR}$	0.1035	0.1014	187.97	
		100	$\hat{\mu}_R$	0.0548	0.0540	167.23	0.0044
			$\hat{\mu}_{Reg}$	0.0505	0.0493	183.20	
			$\hat{\mu}_{GRR}$	0.0502	0.0487	185.45	
		200	$\hat{\mu}_R$	0.0246	0.0240	167.23	0.0020
			$\hat{\mu}_{Reg}$	0.0224	0.0219	183.20	
			$\hat{\mu}_{GRR}$	0.0224	0.0218	184.20	
300	$\hat{\mu}_R$	0.0143	0.0140	167.23	0.0012		
	$\hat{\mu}_{Reg}$	0.0131	0.0128	183.20			
	$\hat{\mu}_{GRR}$	0.0131	0.0127	183.78			
1000	0.8684	50	$\hat{\mu}_R$	0.0399	0.0381	300.26	0.0094
			$\hat{\mu}_{Reg}$	0.0283	0.0284	402.62	
			$\hat{\mu}_{GRR}$	0.0288	0.0282	405.49	
		100	$\hat{\mu}_R$	0.0186	0.0181	300.26	0.0045
			$\hat{\mu}_{Reg}$	0.0141	0.0135	402.62	
			$\hat{\mu}_{GRR}$	0.0141	0.0134	403.97	
		200	$\hat{\mu}_R$	0.0083	0.0080	300.26	0.0020
			$\hat{\mu}_{Reg}$	0.0061	0.0060	402.62	
			$\hat{\mu}_{GRR}$	0.0061	0.0060	403.22	
300	$\hat{\mu}_R$	0.0048	0.0047	300.26	0.0012		
	$\hat{\mu}_{Reg}$	0.0036	0.0035	402.62			
	$\hat{\mu}_{GRR}$	0.0036	0.0035	402.97			

¹ MSE comparison base on 1st order approximation given in expression 3.23.

For the regression and the generalized regression-cum-ratio estimators all the percent relative efficiencies are greater than 100 indicating that all these estimators are better than the RRT mean estimator. The same cannot be said about the ratio estimator because it is better than RRT mean estimator only for larger correlation values between *X* and *Y*. As expected, the generalized regression-cum-ratio estimator presents larger percent relative efficiencies although the improvement over the ordinary regression estimator is only modest.

3.7 Numerical Example

We now compare the performances of different estimators using a real data set. We focus on the ratio estimator, regression estimator and generalized regression-cum-ratio estimator. The sample data come from a very large survey on Information and Communication Technologies (ICT) usage in enterprises in 2009 with seat in Portugal (Smilhily and Storm, 2010). This survey intends to promote the development of the national statistical system in the information society and to contribute to a deeper knowledge about the usage of ICT by enterprises. The target population covers all industries with one and more persons employed in the sections of economic activity C (Manufacturing) to N (Administrative and support service activities) and S (Other service activities), from NACE¹ Rev. 2 (Eurostat, 2008). The data are collected mainly using Electronic Data Interchange, applying direct connection between information systems at the respondent and the National Statistics Institute. For some enterprises the paper questionnaire is still used. The questions in the structural business surveys mainly deal with characteristics that can be found in the organisations' annual reports and financial statements, such as employment, turnover and investment.

In our application the study variable Y is the purchase orders in 2010, collected by the ICT survey in that year. This is typically a confidential variable for enterprises, only known from business surveys. The auxiliary variable X is the turnover of each enterprise. This information can be easily obtained from enterprise records available in the public domain, as administrative information. In 2010, the population survey contained approximately 278000 enterprises and we know the value of X for all these enterprises. The purchase orders information was collected in the ICT survey and we have the values of Y for 5336 enterprises (which answered this question in the ICT survey in 2010). For this study, these 5336 enterprises are considered as our population so that its parameters are known. The scrambling variable S is taken to be a normal random variable with mean equal to zero and standard deviation equal to 10% of the standard deviation of X , that is $\sigma_S = 0.1\sigma_X$. The reported response is given by $Z = Y + S$ (the purchase order value plus a random quantity). The variables Y and X are strongly correlated so we can take advantage of this correlation by using the ratio and regression estimators, as well as a hybrid estimator that combines both. In Table 6.1 we present the results for the ratio estimator, the regression estimator and the generalized ratio-cum-regression estimator for different sample sizes.

¹NACE is derived from the French title "Nomenclature générale des Activités économiques dans les Communautés Européennes" (Statistical classification of economic activities in the European Communities).

Population Characteristics:

$N = 5336, \rho_{XY} = 0.9632$
$\mu_X = 22.99, \mu_Y = 30.19, \sigma_X = 172.09, \sigma_Y = 138.65$ (in millions of Euros)
and $\beta_{YX} = 0.7763$

We use the following samples sizes in our simulation study: $n = 100, 300, 500, 1000$ and 2000 .

Table 3.2 below presents the results for the empirical *MSE* estimates, the theoretical estimates, correct up to first degree of approximation, and the *PRE* of ratio, regression and generalized regression-cum-ratio estimators relative to the RRT mean estimator. We estimate the empirical *MSE* using 5000 samples of size n selected from the population.

Table 3.2: *MSE* correct up to 1st order approximation and *PRE* for the ratio estimator ($\hat{\mu}_R$), the regression estimator ($\hat{\mu}_{Reg}$) and the generalized regression-cum-ratio estimator ($\hat{\mu}_{GRR}$) relative to the RRT mean estimator.

Population		n	Estimator	MSE Estimation			Condition ¹
N	ρ_{XY}			Empirical	Theoretical	<i>PRE</i>	
5336	0.9636	100	$\hat{\mu}_R$	11.5741	16.4778	1162.46	0.3189
			$\hat{\mu}_{Reg}$	8.8601	16.4153	1166.88	
			$\hat{\mu}_{GRR}$	11.6905	15.3461	1248.18	
		300	$\hat{\mu}_R$	4.3423	5.2828	1162.46	0.1022
			$\hat{\mu}_{Reg}$	3.9360	5.2628	1166.88	
			$\hat{\mu}_{GRR}$	4.3858	5.0879	1206.99	
		500	$\hat{\mu}_R$	2.6995	3.0438	1162.46	0.0589
			$\hat{\mu}_{Reg}$	2.6596	3.0323	1166.88	
			$\hat{\mu}_{GRR}$	2.7166	2.9460	1201.03	
		1000	$\hat{\mu}_R$	1.4224	1.3645	1162.46	0.0264
			$\hat{\mu}_{Reg}$	1.4265	1.3594	1166.88	
			$\hat{\mu}_{GRR}$	1.4287	1.3253	1196.91	
2000	$\hat{\mu}_R$	0.5817	0.5249	1162.46	0.0102		
	$\hat{\mu}_{Reg}$	0.6100	0.5229	1166.88			
	$\hat{\mu}_{GRR}$	0.5869	0.5106	1194.95			

¹ *MSE* comparison base on 1st order approximation given in expression 3.23.

According to the results in Table 3.2, all of the percent relative efficiencies are greater than 100, so all the estimators perform better than the RRT mean estimator for the real data also. The *PRE* of the generalized regression-cum-ratio estimator is better than the other estimators, particularly when the sample size is small.

Note that the sample size does not play a role in the *PRE* calculation for the ratio and regression estimators, as can be seen from Equations (3.2), (3.5) and (3.10).

3.8 Conclusions

We can observe from this study that the estimation of the mean of a sensitive variable can be improved by using a non-sensitive auxiliary variable. Although both the ratio and regression estimators perform better than the ordinary RRT mean estimator, the improvement is much larger with the regression estimator. Our simulation study shows that this improvement can be quite substantial for large sample sizes, particularly if the correlation between the study variable and the auxiliary variable is high. Further gains, although modest, can be achieved by using a generalized regression-cum-ratio estimator.

References

- EICHHORN, B. H. & HAYRE, L. S. 1983. Scrambled randomized response methods for obtaining sensitive quantitative data. *Journal of Statistical Planning and Inference*, 7, 307-316.
- EUROSTAT. 2008. NACE Rev. 2 - Statistical classification of economic activities in the European Community. *Official Publications of the European Communities*, 112-285 and 306-311.
- GUPTA, S. N., GUPTA, B. C. & SINGH, S. 2002. Estimation of sensitivity level of personal interview survey questions. *Journal of Statistical Planning and Inference*, 100, 239-247.
- GUPTA, S. & SHABBIR, J. 2004. Sensitivity estimation for personal interview survey questions. *Statistica*, 64, 643-653.
- GUPTA, S., SHABBIR, J. & SEHRA, S. 2010. Mean and sensitivity estimation in optional randomized response models. *Journal of Statistical Planning and Inference*, 140(10), 2870-2874.
- KADILAR, C., CANDAN, M. & CINGI, H. 2007. Ratio estimators using robust regression. *Hacettepe Journal of Mathematics and Statistics*, 36(2), 81-188.
- KADILAR, C. & CINGI, H. 2004. Ratio estimators in simple random sampling. *Applied Mathematics and Computation*, 151, 893-902.
- KADILAR, C. & CINGI, H. 2005. A new estimator using two auxiliary variables. *Applied Mathematics and Computation*, 162, 901-908.
- KADILAR, C. & CINGI, H. 2006. Improvement in estimating the population mean in simple random sampling. *Applied Mathematics Letters*, 19(1), 75-79.

MUKHOPADHYAY, P. 1998. *Theory and Methods of Survey Sampling*, New Delhi, Prentice-Hall of India.

NANGSUE, N. 2009. Adjusted Ratio and Regression Type Estimators for Estimation of Population Mean when some Observations are missing. *World Academy of Science, Engineering and Technology*, 53, 781-784.

PERRI, P. F. 2004. On the efficient use of regression-in-ratio estimator in simple random sampling. *Proceeding of the XLII Meeting of the Italian Statistical Society*, Bari (Italy), 537-540.

PERRI, P. F. 2008. Modified randomized devices for Simmons' model. *Model Assisted Statistics and Applications*, 3, 233-239.

RAY, S. K. & SINGH, R. K. 1981. Difference-cum-ratio type estimators. *Journal of Indian Statistical Association*, 19, 147-151.

SAHA, A. 2008. A randomized response technique for quantitative data under unequal probability sampling. *Journal of statistical Theory and Practice*, 2(4), 589-596.

SHABBIR, J. & GUPTA, S. 2007. On improvement in variance estimation using auxiliary information. *Communication in Statistics-Theory and Methods*, 36(12), 2177-2185.

SHABBIR, J. & GUPTA, S. 2010. Estimation of the finite population mean in two-phase sampling when auxiliary variables are attribute. *Hacettepe Journal of Mathematics and Statistics*, 39(1), 121-129.

SMILHILY, M. & STORM, H. 2010. ICT usage in enterprises - 2009. *Eurostat Publications*, Issue 1.

SOUSA, R., SHABBIR, J. REAL, P. C. & GUPTA, S. 2010. Ratio estimation of the mean of a sensitive variable in the presence of auxiliary information. *Journal of Statistical Theory and Practice*, 4(3), 495-507.

WU, J-W, TIAN, G-L & TANG, M-L. 2008. Two new models for survey sampling with sensitive characteristic: design and analysis. *Metrika*, 67, 251-263.

Appendix B - R Routines

Listing 3.1: R Code for Simulation Study of Proposed Estimator in Chapter 3

```

1
2 proj2_2nd_estimator <- function(N, sigma, mu)
3 {
4
5   #Generation of a bivariate normal population
6   data_yx <- mvrnorm(N, mu, sigma)
7
8   #Study variable
9   Y <- data_yx[,1]
10  #Auxiliary variable, correlated with Y
11  X <- data_yx[,2]
12
13  #Coefficient of correlation between Y and X
14  Ro_YX <- cor(Y,X)
15
16  #Scrambling variable independent of Y and X, with mean=0
17  S <- rnorm(N, mean=0, sd=0.1*sd(X))
18  #Scrambled response
19  Z <- Y+S
20
21  #Coefficient of correlation between Z and X
22  Ro_ZX <- Ro_YX/sqrt(1+(var(S)/var(Y)))
23
24  #population
25  univ <- data.frame(cbind(Y=Y, S=S, Z=Z, X=X, NRAND=runif(N)))
26  univ <- univ[order(univ$NRAND),]
27
28  #Mean of Y
29  mz <- mean(univ$Z)
30  mx <- mean(univ$X)
31  my <- mean(univ$Y)
32
33  mu11 <- sum((univ$Z-mz)*(univ$X-mx))/(N-1)
34  mu12 <- sum((univ$Z-mz)*((univ$X-mx)^2))/(N-1)
35  mu02 <- sum((univ$X-mx)^2)/(N-1)
36  mu03 <- sum((univ$X-mx)^3)/(N-1)
37
38  beta_zx <- Ro_YX*(sd(univ$Y)/sd(univ$X))
39
40  #Samples dimension
41  dim_samp <- c(50,100,150,200,300)
42
43  #Initialize the variables...
44  for (i in 1:length(dim_samp))
45  {
46    #sample dimension

```

3. ESTIMATION OF THE MEAN OF A SENSITIVE VARIABLE IN THE PRESENCE OF AUXILIARY INFORMATION
Appendix B - R Routines

```

47  n <- dim_samp[i]
48  #sample
49  samp <- univ[1:n,]
50  #Sampling rate
51  f <- n/N
52
53  #Ratio
54  R <- mean(univ$X)/mean(samp$X)
55
56  #Ordinary meam
57  est1 <- mean(samp$Z)
58  #Ratio estimator
59  est2 <- mean(samp$Z) * (mx/mean(samp$X))
60  #Regression estimator
61  est3 <- mean(samp$Z) + beta_zx * (mx - mean(samp$X))
62
63  #Coefficient of variation
64  c_x <- sd(univ$X)/mx
65  c_y <- sd(univ$Y)/my
66  c2_x <- c_x^2
67  c2_y <- c_y^2
68  c2_z <- c2_y + (var(univ$S)/(my^2))
69  c_z <- sqrt(c2_z)
70
71  k1 <- (1 - ((1-f) * c2_x/n)) / (1 - ((1-f)/n) * (c2_x - c2_z * (1 - (Ro_ZX^2))))
72  k2 <- (mz/mx) * (1 + k1 * ((Ro_ZX * c_z/c_x) - 2))
73  #Generalized regression-cum-ratio estimator
74  est5 <- (k1 * mean(samp$Z) + k2 * (mx - mean(samp$X))) * (mx/mean(samp$X))
75
76  #Mean Square Error of 1st estimator (ordinal mean)
77  mse1 <- ((1-f)/n) * (var(univ$Y) + var(univ$S))
78
79  #Bias of ratio estimator - 1st degree approximation
80  bias2i <- ((1-f)/n) * my * (c2_x - Ro_ZX * c_z * c_x)
81  #Mean Square Error of ratio estimator - 1st degree approximation
82  mse2i <- ((1-f)/n) * (my^2) * (c2_z + c2_x - 2 * Ro_ZX * c_z * c_x)
83
84  #Bias of regression estimator - 1st degree approximation
85  bias3i <- -beta_zx * ((1-f)/n) * ((mu12/mu11) - (mu03/mu02))
86  #Mean Square Error of regression estimator - 1st degree approximation
87  mse3i <- ((1-f)/n) * (my^2) * c2_z * (1 - (Ro_ZX^2))
88
89  #Bias of genetalized regression-cum-ratio estimator
90  #1st degree approximation
91  bias5i <- (k1-1) * my + k1 * my * ((1-f)/n) * (c2_x - Ro_ZX * c_z * c_x)
92  + k2 * mx * ((1-f)/n) * c2_x
93  #Mean Square Error of generalized regression-cum-ratio estimator
94  #1st degree approximation
95  mse5i <- ((k1-1)^2) * (my^2) + (k1^2) * (my^2) * ((1-f)/n) * (c2_z
96  + 3 * c2_x - 4 * Ro_ZX * c_z * c_x) + (k2^2) * (mx^2) * ((1-f)/n) * c2_x

```

```

97         -2*k1*(my^2)*(1-f)/n*(c2_x-Ro_ZX*c_z*c_x)
98         -2*k2*my*mx*(1-f)/n*c2_x
99         -2*k1*k2*my*mx*(1-f)/n*(Ro_ZX*c_z*c_x-2*c2_x)
100
101     cond1 <- ((1-f)/n)*c2_x
102
103     #Empirical results
104     #Simulation of 5000 replicas of estimates
105     ...
106
107     #Results
108     res <- rbind(res,c(N,n,Ro_YX,Ro_ZX,R,
109                   c_x,c_y,c_z,mx,my,mz,
110                   med_est1,med_est2,med_est3,med_est5,
111                   bias2i,bias3i,bias5i,
112                   emp_mse1,mse1,emp_mse2,mse2i,
113                   emp_mse3,mse3i,emp_mse5,mse5i,cond1))
114 }
115 colnames(res) <- c("N","n","RhoXY","RhoZX","R",
116                  "Cx","Cy","Cz","mX","mY","mZ",
117                  "Est1","Est2","Est3","Est5",
118                  "BIAS2I","BIAS3I","BIAS5I",
119                  "EMP_MSE1","MSE1","EMP_MSE2","MSE2I",
120                  "EMP_MSE3","MSE3I","EMP_MSE5","MSE5I","COND1")
121 return(res)
122 }
123
124 #Package for generation
125 require(MASS)
126 N <- 1000
127
128 #Parameters
129 sigma1 <- matrix(c(9,1.9,1.9,4),2,2)
130 sigma2 <- matrix(c(10,3,3,2),2,2)
131 sigma3 <- matrix(c(6,3,3,2),2,2)
132 mu <- c(2,2)
133
134 res <- NULL
135 for (i in 1:length(N))
136 {
137     res <- rbind(res,proj2_2nd_estimator(N[i],sigma1,mu))
138     res <- rbind(res,proj2_2nd_estimator(N[i],sigma2,mu))
139     res <- rbind(res,proj2_2nd_estimator(N[i],sigma3,mu))
140 }
141 write.table(res,"chapter3_ss_results.txt",sep="\t",dec=",",row.names=FALSE)

```

Listing 3.2: R Code for Numerical Example of Proposed Estimator in Chapter 3

```
1
2 proj2_2nd_estimator_real <- function(Y,X,N)
3 {
4   #Coefficient of correlation between Y and X
5   Ro_YX <- cor(Y,X)
6
7   #Scrambling variable independent of Y and X, with mean=0
8   S <- rnorm(N,mean=0,sd=sd(X)*0.1)
9   #Scrambled response
10  Z <- Y+S
11
12  #Coefficient of correlation between Z and X
13  Ro_ZX <- Ro_YX/sqrt(1+(var(S)/var(Y)))
14
15  #population
16  univ <- data.frame(cbind(Y=Y,S=S,Z=Z,X=X,NRAND=runif(N)))
17  univ <- univ[order(univ$NRAND),]
18
19  #Mean of Y
20  mz <- mean(univ$Z)
21  mx <- mean(univ$X)
22  my <- mean(univ$Y)
23  ms <- mean(univ$S)
24
25  mu11 <- sum((univ$Z-mz)*(univ$X-mx))/(N-1)
26  mu12 <- sum((univ$Z-mz)*((univ$X-mx)^2))/(N-1)
27  mu02 <- sum((univ$X-mx)^2)/(N-1)
28  mu03 <- sum((univ$X-mx)^3)/(N-1)
29
30  beta_zx <- Ro_YX*(sd(univ$Y)/sd(univ$X))
31
32  #Samples dimension
33  dim_samp <- c(100,300,500,1000,2000)
34
35  #Initialize variables...
36  for (i in 1:length(dim_samp))
37  {
38    #sample dimension
39    n <- dim_samp[i]
40    #sample
41    samp <- univ[1:n,]
42    #Sampling rate
43    f <- n/N
44
45    #Ratio
46    R <- mean(univ$X)/mean(samp$X)
47
```



```

48  #Ordinary meam
49  est1 <- mean(samp$Z)
50  #Ratio estimator
51  est2 <- mean(samp$Z) * (mx/mean(samp$X))
52  #Regression estimator
53  est3 <- mean(samp$Z) + beta_zx * (mx - mean(samp$X))
54
55  #Coefficient of variation
56  c_x <- sd(univ$X) / mx
57  c_y <- sd(univ$Y) / my
58  c2_x <- c_x^2
59  c2_y <- c_y^2
60  c2_z <- c2_y + (var(univ$S) / (my^2))
61  c_z <- sqrt(c2_z)
62
63  k1 <- (1 - ((1-f) * c2_x / n)) / (1 - ((1-f) / n) * (c2_x - c2_z * (1 - (Ro_ZX^2))))
64  k2 <- (mz / mx) * (1 + k1 * ((Ro_ZX * c_z / c_x) - 2))
65  #Generalized regression-cum-ratio estimator
66  est5 <- (k1 * mean(samp$Z) + k2 * (mx - mean(samp$X))) * (mx / mean(samp$X))
67
68  #Mean Square Error of 1st estimator (ordinal mean)
69  mse1 <- ((1-f) / n) * (var(univ$Y) + var(univ$S))
70
71  #Bias of ratio estimator - 1st degree approximation
72  bias2i <- ((1-f) / n) * my * (c2_x - Ro_ZX * c_z * c_x)
73  #Mean Square Error of ratio estimator - 1st degree approximation
74  mse2i <- ((1-f) / n) * (my^2) * (c2_z + c2_x - 2 * Ro_ZX * c_z * c_x)
75
76  #Bias of regression estimator - 1st degree approximation
77  bias3i <- -beta_zx * ((1-f) / n) * ((mu12 / mu11) - (mu03 / mu02))
78  #Mean Square Error of regression estimator - 1st degree approximation
79  mse3i <- ((1-f) / n) * (my^2) * c2_z * (1 - (Ro_ZX^2))
80
81  #Bias of genetalized regression-cum-ratio estimator
82  #1st degree approximation
83  bias5i <- (k1 - 1) * my + k1 * my * ((1-f) / n) * (c2_x - Ro_ZX * c_z * c_x)
84  + k2 * mx * ((1-f) / n) * c2_x
85  #Mean Square Error of generalized regression-cum-ratio estimator
86  #1st degree approximation
87  mse5i <- ((k1 - 1)^2) * (my^2) + (k1^2) * (my^2) * ((1-f) / n)
88  * (c2_z + 3 * c2_x - 4 * Ro_ZX * c_z * c_x) + (k2^2) * (mx^2) * ((1-f) / n)
89  * c2_x - 2 * k1 * (my^2) * ((1-f) / n) * (c2_x - Ro_ZX * c_z * c_x)
90  - 2 * k2 * my * mx * ((1-f) / n) * c2_x - 2 * k1 * k2 * my * mx * ((1-f) / n) * (Ro_ZX * c_z * c_x - 2 * c2_x)
91
92  cond1 <- ((1-f) / n) * c2_x
93
94  #Empirical results
95  #Simulation of 5000 replicas of estimates
96  ...
97

```

3. ESTIMATION OF THE MEAN OF A SENSITIVE VARIABLE IN THE PRESENCE OF AUXILIARY INFORMATION

Appendix B - R Routines

```

98   #Results
99   res <- rbind(res, c(N, n, Ro_YX, Ro_ZX, R,
100                   c_x, c_y, c_z, mx, my, mz, ms,
101                   med_est1, med_est2, med_est3, med_est5,
102                   bias2i, bias3i, bias5i,
103                   emp_mse1, mse1, emp_mse2, mse2i,
104                   emp_mse3, mse3i, emp_mse5, mse5i, cond1))
105 }
106 colnames(res) <- c("N", "n", "RhoXY", "RhoZX", "R",
107                  "Cx", "Cy", "Cz", "mX", "mY", "mZ", "ms",
108                  "Est1", "Est2", "Est3", "Est5",
109                  "BIAS2I", "BIAS3I", "BIAS5I",
110                  "EMP_MSE1", "MSE1", "EMP_MSE2", "MSE2I",
111                  "EMP_MSE3", "MSE3I", "EMP_MSE5", "MSE5I", "COND1")
112 return(res)
113 }
114
115 #Package for generation
116 require(MASS)
117
118 #Import data
119 data_yx <- read.table("IUTICE10.txt", sep="\t", dec=",", header = T)
120 #Study variable (purchase, millions of euros)
121 Y <- data_yx[,3]
122 #Auxiliary variable, correlated with Y (turnover, millions of euros)
123 X <- data_yx[,2]
124
125 #Data application
126 N <- dim(data_yx)[1]
127 res <- proj2_2nd_estimator_real(Y, X, N)
128
129 #Export data
130 write.table(res_exp, "chapter3_ne_results.txt", sep="\t", dec=",", row.names=FALSE)

```



Exponential Type Estimators of the Mean of a Sensitive Variable in the Presence of Non Sensitive Auxiliary Information

Abstract

Sousa et al. (2010) and Gupta et al. (2012) suggested ratio and regression type estimators of the mean of a sensitive variable using non-sensitive auxiliary variable. This paper proposes exponential type estimators using one and two auxiliary variables to improve the efficiency of mean estimator based on a Randomized Response Technique (RRT). The expressions for the Mean Square Errors (MSE 's) and bias, up to first order approximation, have been obtained. It is shown that the proposed exponential type estimators are more efficient than the existing estimators. The gain in efficiency over the existing estimators has also been shown with a simulation study and by using real data.

Accepted as: KOYUNCU, N., GUPTA, S., SOUSA, R. 2013. Exponential type estimators of the mean of a sensitive variable in the presence of non-sensitive auxiliary information. *Communications in Statistics - Simulation and Computation*.

4.1 Introduction

Randomized Response Technique (RRT) is used to estimate the proportion of people in a community bearing a stigmatizing characteristic like habitual tax evasion, reckless driving, indiscriminate gambling, abortion etc. In such situations we cannot expect to get a truthful direct response to a sensitive question. Eichhorn and Hayre (1983), Gupta and Shabbir (2004), Gupta et al. (2002, 2010), Wu et al. (2008), Perri (2008), and many others have estimated the mean of a sensitive variable when the study variable is sensitive and there is no auxiliary variable. Sousa et al. (2010) and Gupta et al. (2012) suggested mean estimators based on RRT models using an auxiliary variable that can be directly observed. In sampling literature, Bahl and Tuteja (1991), Shabbir and Gupta (2007), Grover (2010) and Koyuncu (2012) have studied exponential type estimators to get more efficient estimates. In this study we have proposed exponential type estimators of the mean of a sensitive variable using non-sensitive auxiliary information. We have discussed the cases when one or two non-sensitive auxiliary variables are available.

4.2 Terminology

Let Y be the study variable, a sensitive variable which cannot be observed directly. Let X_1 and X_2 be non-sensitive auxiliary variables which have a positive correlation with Y . Let S , be a scrambling variable, independent of Y , X_1 and X_2 . The respondent is asked to report a scrambled response for Y given by $Z = Y + S$ but is asked to provide a true response for X_1 and X_2 . Let a random sample of size n be drawn without replacement from a finite population $U = (U_1, U_2, \dots, U_N)$. For the i^{th} unit ($i = 1, 2, \dots, N$), let y_i , x_{1i} and x_{2i} respectively be the values of the study variable Y and auxiliary variables X_1 and X_2 . Let $\bar{y} = \frac{\sum_{i=1}^n y_i}{n}$, $\bar{x}_1 = \frac{\sum_{i=1}^n x_{1i}}{n}$, $\bar{x}_2 = \frac{\sum_{i=1}^n x_{2i}}{n}$ and $\bar{z} = \frac{\sum_{i=1}^n z_i}{n}$ be the sample means and $\bar{Y} = E(Y)$, $\bar{X}_1 = E(X_1)$, $\bar{X}_2 = E(X_2)$ and $\bar{Z} = E(Z)$ be the population means for Y , X_1 , X_2 and Z respectively. We assume that \bar{X}_1 , \bar{X}_2 are known and $\bar{S} = E(S) = 0$. Thus $E(Z) = E(Y)$ and $C_z^2 = C_y^2 + (S_s^2/\bar{Y}^2)$, where C_z and C_y are the coefficients of the variation of z and y respectively.

To obtain the bias and MSE expressions, let us define

$$e_0 = \frac{\bar{z} - \bar{Z}}{\bar{Z}}, e_1 = \frac{\bar{x}_1 - \bar{X}_1}{\bar{X}_1}, e_2 = \frac{\bar{x}_2 - \bar{X}_2}{\bar{X}_2}, e_3 = \frac{s_{x1}^2 - S_{x1}^2}{S_{x1}^2} \text{ and } e_4 = \frac{s_{zx1}^2 - S_{zx1}^2}{S_{zx1}^2}.$$

Using these notations,

$$E(e_i) = 0, \quad i = 0, 1, 2, 3, 4.$$

$$E(e_0^2) = \lambda C_z^2, E(e_1^2) = \lambda C_{x1}^2, E(e_2^2) = \lambda C_{x2}^2, E(e_0 e_1) = \lambda C_{zx1},$$

$$E(e_0e_2) = \lambda C_{zx2}, E(e_1e_2) = \lambda C_{x1x2}, E(e_1e_3) = \lambda \frac{1}{\bar{X}_1} \frac{\mu_{03}}{\mu_{02}},$$

$$E(e_1e_4) = \lambda \frac{1}{\bar{X}_1} \frac{\mu_{12}}{\mu_{11}},$$

where $\lambda = \frac{1-f}{n}$ and $\mu_{rs} = \frac{1}{N-1} \sum_{i=1}^N (z_i - \bar{Z})^r (x_{1i} - \bar{X}_1)^s$.

4.3 Estimators Review

We describe below some existing mean estimators and their bias and *MSE* formulas.

(i) Ordinary sample mean (\bar{z}) of scrambled responses:

$$\hat{\mu}_Y = \bar{z}. \tag{4.1}$$

$$MSE(\hat{\mu}_y) = \lambda (S_y^2 + S_s^2), \tag{4.2}$$

where

$$S_y^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{Y})^2, S_{x1}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_{1i} - \bar{X}_1)^2, S_s^2 = \frac{1}{N-1} \sum_{i=1}^N (s_i - \bar{S})^2.$$

(ii) Sousa et al. (2010) ratio type estimator:

$$\hat{\mu}_R = \bar{z} \frac{\bar{X}_1}{\bar{x}_1}. \tag{4.3}$$

$$Bias(\hat{\mu}_R) \cong \bar{Y} \lambda (C_{x1}^2 - C_{x1z}), \tag{4.4}$$

where $C_z^2 = \left(C_y^2 + \frac{S_s^2}{\bar{Y}^2} \right), \rho_{zx1} = \frac{\rho_{yx1}}{\sqrt{1 + \frac{S_s^2}{S_y^2}}}, \bar{Z} = \bar{Y}$.

$$MSE(\hat{\mu}_R) \cong \lambda \bar{Y}^2 (C_z^2 + C_{x1}^2 - 2C_{x1z}). \tag{4.5}$$

(iii) Gupta et al. (2012) regression estimator:

$$\hat{\mu}_{Reg} = \bar{z} + \hat{\beta}_{zx1} (\bar{X}_1 - \bar{x}_1), \tag{4.6}$$

where $\hat{\beta}_{zx1} = \frac{S_{zx1}}{S_{x1}^2} = \frac{S_{yx1}}{S_{x1}^2}$, is the sample regression coefficient between *Z* and *X*₁.

$$Bias(\hat{\mu}_{Reg}) \cong -\beta_{zx1} \lambda \left(\frac{\mu_{12}}{\mu_{11}} - \frac{\mu_{03}}{\mu_{02}} \right), \tag{4.7}$$

where $\beta_{zx1} = \frac{S_{zx1}}{S_{x1}^2} = \frac{S_{yx1}}{S_{x1}^2} = \rho_{yx1} \frac{S_y}{S_{x1}} = \beta_{yx1}$ is the population regression coefficient between Z and X_1 .

Recognizing $\bar{Z} = \bar{Y}$

$$MSE(\hat{\mu}_{Reg}) \cong \lambda \bar{Y}^2 C_z^2 \left[1 - \frac{S_{zx1}^2}{S_{x1}^2 S_z^2} \right] = \lambda \bar{Y}^2 C_z^2 [1 - \rho_{zx1}^2]$$

or

$$MSE(\hat{\mu}_{Reg}) \cong \lambda S_y^2 \left[\left(1 + \frac{S_s^2}{S_y^2} \right) - \rho_{yx1}^2 \right]. \quad (4.8)$$

(iv) Gupta et al. (2012) generalized regression-cum-ratio estimator:

$$\hat{\mu}_{GRR} = [k_1 \bar{z} + k_2 (\bar{X} - \bar{x})] \left(\frac{\bar{X}}{\bar{x}} \right), \quad (4.9)$$

where k_1 and k_2 are constants.

$$Bias(\hat{\mu}_{GRR}) \cong (k_1 - 1)\bar{Y} + k_1 \bar{Y} \lambda \{C_x^2 - \rho_{zx} C_z C_x\} + k_2 \bar{X} \lambda C_x^2. \quad (4.10)$$

$$\begin{aligned} MSE(\hat{\mu}_{GRR}) \cong & (k_1 - 1)^2 \bar{Y}^2 + k_1^2 \bar{Y}^2 \lambda \{C_z^2 + 3C_x^2 - 4\rho_{zx} C_z C_x\} \\ & + k_2^2 \bar{X}^2 \lambda C_x^2 - 2k_1 \bar{Y}^2 \lambda \{C_x^2 - \rho_{zx} C_z C_x\} \\ & - 2k_2 \bar{Y} \bar{X} \lambda C_x^2 - 2k_1 k_2 \bar{Y} \bar{X} \lambda \{\rho_{zx} C_z C_x - 2C_x^2\}. \end{aligned} \quad (4.11)$$

From Equation (4.11), the optimum values of k_1 and k_2 are given by

$$k_{1(opt)} = \frac{1 - \lambda C_x^2}{1 - \lambda \{C_x^2 - C_z^2 (1 - \rho_{zx}^2)\}} \quad (4.12)$$

and

$$k_{2(opt)} = \frac{\bar{Y}}{\bar{X}} \left\{ 1 + k_{1(opt)} \left(\frac{\rho_{zx} C_z}{C_x} - 2 \right) \right\}, \quad (4.13)$$

the minimum MSE of $\hat{\mu}_{GRR}$ can be written as follows:

$$MSE(\hat{\mu}_{GRR})_{min} \cong \frac{\bar{Y}^2 C_z^2 (1 - \rho_{zx}^2) \lambda \{1 - \lambda C_x^2\}}{C_z^2 (1 - \rho_{zx}^2) \lambda + \{1 - \lambda C_x^2\}}. \quad (4.14)$$

4.4 Proposed Exponential Type Estimators

Our first proposed estimator, which we call "generalized regression-cum-exponential estimator" follows Grover (2010) and Shabbir and Gupta (2007), and is given by

$$\hat{\mu}_{exp1} = [w_1 \bar{z} + w_2 (\bar{X}_1 - \bar{x}_1)] \exp \left(\frac{\bar{X}_1 - \bar{x}_1}{\bar{X}_1 + \bar{x}_1} \right), \quad (4.15)$$

where w_1 and w_2 are suitable weights. Expressing (4.15) in terms of e 's (defined earlier) and retaining terms of e 's up to second-order we have

$$\hat{\mu}_{exp1} - \bar{Z} \cong \left[w_1 \bar{Z} + w_1 \bar{Z} e_0 - w_2 \bar{X}_1 e_1 - \frac{1}{2} w_1 \bar{Z} e_1 - \frac{1}{2} w_1 \bar{Z} e_0 e_1 + \frac{1}{2} w_2 \bar{X}_1 e_1^2 + \frac{3}{8} w_1 \bar{Z} e_1^2 - \bar{Z} \right]. \quad (4.16)$$

The *Bias* and *MSE* of $\hat{\mu}_{exp1}$, to the first order of approximation, are given by

$$Bias(\hat{\mu}_{exp1}) \cong (w_1 - 1) \bar{Y} + \lambda \left\{ \frac{1}{2} w_1 \bar{Y} \left(\frac{3}{4} C_{x1}^2 - C_{zx1} \right) + \frac{1}{2} w_2 \bar{X}_1 C_{x1}^2 \right\} \quad (4.17)$$

and

$$MSE(\hat{\mu}_{exp1}) \cong \{ \bar{Y}^2 + w_1^2 \bar{Y}^2 (1 + \lambda (C_z^2 + C_{x1}^2 - 2C_{zx1})) + w_2^2 \bar{X}_1^2 \lambda C_{x1}^2 + w_1 \bar{Y}^2 \left(\lambda \left(C_{zx1} - \frac{3}{4} C_{x1}^2 \right) - 2 \right) - w_2 \bar{Y} \bar{X}_1 \lambda C_{x1}^2 + 2w_1 w_2 \bar{Y} \bar{X}_1 \lambda (C_{x1}^2 - C_{zx1}) \}, \quad (4.18)$$

and optimum values of w_1 and w_2 respectively are found as

$$w_1^* = \frac{1 - \frac{1}{8} \lambda C_{x1}^2}{1 + \lambda C_z^2 (1 - \rho_{zx1}^2)} \quad (4.19)$$

and

$$w_2^* = \frac{\bar{Y}}{2\bar{X}_1} \frac{C_{x1}^2 - 2C_{x1}^2 + 2C_{zx1} + \lambda C_{x1}^2 \left(C_z^2 (1 - \rho_{zx1}^2) + \frac{1}{4} (C_{x1}^2 - C_{zx1}) \right)}{C_{x1}^2 [1 + \lambda C_z^2 (1 - \rho_{zx1}^2)]}. \quad (4.20)$$

Substituting these optimum values in (4.18), the minimum *MSE* of $\hat{\mu}_{exp1}$ can be written as follows:

$$MSE_{min}(\hat{\mu}_{exp1}) \cong \bar{Y}^2 \left[1 - \frac{\lambda^2 C_{x1}^2 \left(\frac{1}{16} C_{x1}^2 + C_z^2 (1 - \rho_{zx1}^2) \right) + 4}{4 + [1 + \lambda C_z^2 (1 - \rho_{zx1}^2)]} \right] \quad (4.21)$$

or

$$MSE_{min}(\hat{\mu}_{exp1}) \cong \left[\frac{MSE(\hat{\mu}_{Reg})}{\left[1 + \frac{MSE(\hat{\mu}_{Reg})}{\bar{Y}^2} \right]} - \frac{\lambda C_{x1}^2 \left(MSE(\hat{\mu}_{Reg}) + \lambda \frac{1}{16} C_{x1}^2 \bar{Y}^2 \right)}{4 \left[1 + \frac{MSE(\hat{\mu}_{Reg})}{\bar{Y}^2} \right]} \right]. \quad (4.22)$$

Note that the optimum choice of the constants w_1 and w_2 involve unknown parameters. These quantities can be guessed through a pilot sample survey or through experience gathered in due course of time, as mentioned Upadhyaya and Singh (2006), and Koyuncu and Kadilar (2009).

The estimator defined in (4.15) can be generalized to the case of multiple auxiliary variables. We consider below the case of two auxiliary non-sensitive variables. This estimator is given by

$$\hat{\mu}_{exp2} = [d_1\bar{z} + d_2(\bar{X}_1 - \bar{x}_1) + d_3(\bar{X}_2 - \bar{x}_2)] \exp\left(\frac{(\bar{X}_1 - \bar{x}_1) + (\bar{X}_2 - \bar{x}_2)}{(\bar{X}_1 + \bar{x}_1) + (\bar{X}_2 + \bar{x}_2)}\right). \quad (4.23)$$

Expressing (4.23) in terms of e 's and retaining up to second-order terms in e 's we have

$$\hat{\mu}_{exp2} \cong \left\{ d_1\bar{Z}(1 + e_0) - d_2\bar{X}_1e_1 - d_3\bar{X}_2e_2 \right\} \left\{ 1 - \frac{\bar{X}_1}{2(\bar{X}_1 + \bar{X}_2)}e_1 - \frac{\bar{X}_2}{2(\bar{X}_1 + \bar{X}_2)}e_2 + \frac{3\bar{X}_1^2}{8(\bar{X}_1 + \bar{X}_2)^2}e_1^2 + \frac{6\bar{X}_1\bar{X}_2}{8(\bar{X}_1 + \bar{X}_2)^2}e_1e_2 + \frac{3\bar{X}_2^2}{8(\bar{X}_1 + \bar{X}_2)^2}e_2^2 \right\}. \quad (4.24)$$

The Bias and MSE of $\hat{\mu}_{exp2}$, to the first order of approximation, are given by

$$\begin{aligned} Bias(\hat{\mu}_{exp2}) &\cong (d_1 - 1)\bar{Z} \\ &+ \frac{d_1\lambda\bar{Z}}{2(\bar{X}_1 + \bar{X}_2)} \left(-\bar{X}_1C_{zx1} - \bar{X}_2C_{zx2} + \frac{3\bar{X}_1^2}{4(\bar{X}_1 + \bar{X}_2)}C_{x1}^2 \right. \\ &\quad \left. + \frac{3\bar{X}_2^2}{4(\bar{X}_1 + \bar{X}_2)}C_{x2}^2 + \frac{3\bar{X}_1\bar{X}_2}{2(\bar{X}_1 + \bar{X}_2)}C_{x1x2} \right) \\ &+ \frac{d_2\lambda\bar{X}_1}{2(\bar{X}_1 + \bar{X}_2)}(\bar{X}_1C_{x1}^2 + \bar{X}_2C_{x1x2}) \\ &+ \frac{d_3\lambda\bar{X}_2}{2(\bar{X}_1 + \bar{X}_2)}\lambda(\bar{X}_1C_{x1x2} + \bar{X}_2C_{x2}^2) \end{aligned} \quad (4.25)$$

and

$$\begin{aligned} MSE(\hat{\mu}_{exp2}) &\cong \bar{Z}^2 + d_1A - d_2B - d_3C + d_1^2D + d_2^2\bar{X}_1^2\lambda C_{x1}^2 + d_3^2\bar{X}_2^2\lambda C_{x2}^2 \\ &+ 2d_1d_2F + 2d_1d_3G + 2d_2d_3\bar{X}_1\bar{X}_2\lambda C_{x1x2}, \end{aligned} \quad (4.26)$$

where

$$\begin{aligned} A &= \bar{Z}^2 \left(-2 + \lambda \left\{ \frac{\bar{X}_1C_{zx1}}{(\bar{X}_1 + \bar{X}_2)} + \frac{\bar{X}_2C_{zx2}}{(\bar{X}_1 + \bar{X}_2)} - \frac{3\bar{X}_1^2C_{x1}^2}{4(\bar{X}_1 + \bar{X}_2)^2} - \frac{6\bar{X}_1\bar{X}_2C_{x1x2}}{4(\bar{X}_1 + \bar{X}_2)^2} \right. \right. \\ &\quad \left. \left. - \frac{3\bar{X}_2^2C_{x2}^2}{4(\bar{X}_1 + \bar{X}_2)^2} \right\} \right), \\ B &= \lambda \frac{\bar{Z}}{(\bar{X}_1 + \bar{X}_2)} (\bar{X}_1^2C_{x1}^2 + \bar{X}_1\bar{X}_2C_{x1x2}), \\ C &= \lambda \frac{\bar{Z}}{(\bar{X}_1 + \bar{X}_2)} (\bar{X}_2^2C_{x2}^2 + \bar{X}_1\bar{X}_2C_{x1x2}), \\ D &= \bar{Z}^2 + \lambda \left\{ \bar{Z}^2C_z^2 + \frac{\bar{X}_1^2\bar{Z}^2C_{x1}^2}{(\bar{X}_1 + \bar{X}_2)^2} + \frac{\bar{X}_2^2\bar{Z}^2C_{x2}^2}{(\bar{X}_1 + \bar{X}_2)^2} - 2\frac{\bar{X}_1\bar{Z}^2C_{zx1}}{(\bar{X}_1 + \bar{X}_2)} - 2\frac{\bar{X}_2\bar{Z}^2C_{zx2}}{(\bar{X}_1 + \bar{X}_2)} \right. \\ &\quad \left. + \frac{2\bar{X}_1\bar{X}_2\bar{Z}^2C_{x1x2}}{(\bar{X}_1 + \bar{X}_2)^2} \right\}, \end{aligned}$$

$$F = \lambda \left(\frac{\bar{Z}\bar{X}_1^2}{(\bar{X}_1 + \bar{X}_2)} C_{x1}^2 + \frac{\bar{Z}\bar{X}_1\bar{X}_2}{(\bar{X}_1 + \bar{X}_2)} C_{x1x2} - \bar{Z}\bar{X}_1 C_{zx1} \right),$$

$$G = \lambda \left(\frac{\bar{Z}\bar{X}_2^2}{(\bar{X}_1 + \bar{X}_2)} C_{x2}^2 + \frac{\bar{Z}\bar{X}_1\bar{X}_2}{(\bar{X}_1 + \bar{X}_2)} C_{x1x2} - \bar{Z}\bar{X}_2 C_{zx2} \right),$$

and optimum values of d_1 , d_2 and d_3 are respectively found as

$$d_1^* = \frac{1}{2D} \frac{\left[\begin{array}{l} A(D\lambda S_{x1x2} - FG)^2 + (BDG + CDF + 2AFG)(D\lambda S_{x1x2} - FG) \\ -G(CD + AG)(D\lambda S_{x1}^2 - F^2) - F(AF + BD)(D\lambda S_{x2}^2 - G^2) \\ -A(D\lambda S_{x1}^2 - F^2)(D\lambda S_{x2}^2 - G^2) \end{array} \right]}{(D\lambda S_{x1}^2 - F^2)(D\lambda S_{x2}^2 - G^2) - (D\lambda S_{x1x2} - FG)^2}, \quad (4.27)$$

$$d_2^* = \frac{1}{2} \frac{(AF + BD)(D\lambda S_{x2}^2 - G^2) - (AG + CD)(D\lambda S_{x1x2} - FG)}{(D\lambda S_{x1}^2 - F^2)(D\lambda S_{x2}^2 - G^2) - (D\lambda S_{x1x2} - FG)^2}, \quad (4.28)$$

and

$$d_3^* = \frac{1}{2} \frac{(AG + CD)(D\lambda S_{x1}^2 - F^2) - (AF + BD)(D\lambda S_{x1x2} - FG)}{(D\lambda S_{x1}^2 - F^2)(D\lambda S_{x2}^2 - G^2) - (D\lambda S_{x1x2} - FG)^2}. \quad (4.29)$$

Substituting these optimum values in (4.26), the minimum MSE of $\hat{\mu}_{exp2}$ can be written as follows:

$$MSE(\hat{\mu}_{exp2})_{min} = \bar{Z}^2 - \frac{A^2}{4D} \frac{\left[\begin{array}{l} (AG + CD)^2(D\lambda S_{x1}^2 - F^2) \\ + (AF + BD)^2(D\lambda S_{x2}^2 - G^2) \\ - 2(AG + CD)(AF + BD)(D\lambda S_{x1x2} - FG) \end{array} \right]}{4D(D\lambda S_{x1}^2 - F^2)(D\lambda S_{x2}^2 - G^2) - (D\lambda S_{x1x2} - FG)^2}. \quad (4.30)$$

4.5 Comparison with Gupta et al. (2012) Estimators

First, we compare the proposed "generalized regression-cum-exponential estimator" with the Gupta et al. (2012) regression estimator. Note that

$MSE(\hat{\mu}_{exp1}) < MSE(\hat{\mu}_{Reg})$ if

$$MSE(\hat{\mu}_{Reg}) - \frac{MSE(\hat{\mu}_{Reg})}{\left[1 + \frac{MSE(\hat{\mu}_{Reg})}{\bar{Y}^2}\right]} + \frac{\lambda C_x^2 \left(MSE(\hat{\mu}_{Reg}) + \lambda \frac{1}{16} C_x^2 \bar{Z}^2 \right)}{4 \left[1 + \frac{MSE(\hat{\mu}_{Reg})}{\bar{Y}^2} \right]} > 0$$

or

$$\frac{\frac{(MSE(\hat{\mu}_{Reg}))^2}{\bar{Y}^2}}{\left[1 + \frac{MSE(\hat{\mu}_{Reg})}{\bar{Y}^2}\right]} + \frac{\lambda C_x^2 \left(MSE(\hat{\mu}_{Reg}) + \lambda \frac{1}{16} C_x^2 \bar{Z}^2 \right)}{4 \left[1 + \frac{MSE(\hat{\mu}_{Reg})}{\bar{Y}^2}\right]} > 0. \quad (4.31)$$

From (4.31), we can see easily that proposed "generalized regression-cum-exponential estimator" is always more efficient than regression estimator of Gupta et al. (2012).

Secondly, we compare the proposed "generalized regression-cum-exponential estimator" with Gupta et al. (2012) generalized regression-cum-ratio estimator

$$MSE(\hat{\mu}_{exp1})_{min} < MSE(\hat{\mu}_{GRR})_{min} \text{ if}$$

$$\frac{\lambda^2 C_x^2 \left(\frac{1}{16} C_x^2 + C_z^2 (1 - \rho_{zx}^2) \right) + 4}{4 [1 + \lambda C_z^2 (1 - \rho_{zx}^2)]} + \frac{C_z^2 (1 - \rho_{zx}^2) \lambda \{1 - \lambda C_x^2\}}{C_z^2 (1 - \rho_{zx}^2) \lambda + \{1 - \lambda C_x^2\}} > 1. \quad (4.32)$$

When the condition (4.32) is satisfied, we can infer that the suggested estimator is more efficient than Gupta et al. (2012) generalized regression-cum-ratio estimator.

4.6 Simulation Study

In this section, we investigate the efficiency of proposed exponential estimators to existing estimators. The simulation study is carried out to compare the *Bias* and *MSE* of the estimators both empirically and theoretically. In the simulation study, we consider two finite populations of size $N = 1000$ generated from a multivariate normal distribution with the same theoretical mean of $[Y, X_1, X_2]$ as $\mu = [5, 5, 5]$ and different covariance matrices as given below.

Population 1

$$\sigma^2 = \begin{bmatrix} 10 & 3 & 2.9 \\ 3 & 2 & 1.1 \\ 2.9 & 1.1 & 2 \end{bmatrix}, \rho_{X_1Y} = 0.6817, \rho_{X_2Y} = 0.6705.$$

Population 2

$$\sigma^2 = \begin{bmatrix} 6 & 3 & 2.9 \\ 3 & 2 & 1.1 \\ 2.9 & 1.1 & 2 \end{bmatrix}, \rho_{X_1Y} = 0.8706, \rho_{X_2Y} = 0.8706.$$

The scrambling variable S is taken to be a normal variable with mean equal to zero and standard deviation equal to 10% of the standard deviation of X_1 . The reported response is given by $Z = Y + S$. For each population we considered four sample sizes:

$n = 50, 100, 200$ and 300 . The percent relative efficiency (PRE) is calculated from following equations

$$PRE = \frac{MSE(\hat{\mu}_Y)}{MSE(\hat{\mu}_i)} \times 100,$$

where $i = R, Reg, GRR, exp1, exp2$.

The empirical MSE , theoretical MSE and Percent Relative Efficiency (PRE) values for all estimators are given in Table 4.1 and Table 4.2.

Table 4.1: Empirical MSE , theoretical MSE correct up to 1^{st} order approximation and PRE of all estimators.

Population			MSE Estimation			
N	ρ_{X1Y} ρ_{X2Y}	n	Estimator	Empirical	Theoretical	PRE
1000	0.6817 0.6705	50	$\hat{\mu}_Y$	0.1193	0.1953	100.00
			$\hat{\mu}_R$	0.1193	0.1145	170.64
			$\hat{\mu}_{Reg}$	0.1083	0.1047	186.50
			$\hat{\mu}_{GRR}$	0.1094	0.1043	187.26
			$\hat{\mu}_{exp1}$	0.1089	0.1042	187.34
			$\hat{\mu}_{exp2}$	0.0857	0.0827	236.13
		100	$\hat{\mu}_Y$	0.0900	0.0925	100.00
			$\hat{\mu}_R$	0.0544	0.0542	170.64
			$\hat{\mu}_{Reg}$	0.0499	0.0496	186.50
			$\hat{\mu}_{GRR}$	0.0503	0.0495	186.86
			$\hat{\mu}_{exp1}$	0.0501	0.0495	186.90
			$\hat{\mu}_{exp2}$	0.0390	0.0393	235.69
		200	$\hat{\mu}_Y$	0.0404	0.0411	100.00
			$\hat{\mu}_R$	0.0240	0.0241	170.64
			$\hat{\mu}_{Reg}$	0.0220	0.0220	186.50
			$\hat{\mu}_{GRR}$	0.0221	0.0220	186.66
			$\hat{\mu}_{exp1}$	0.0220	0.0220	186.67
			$\hat{\mu}_{exp2}$	0.0172	0.0175	235.47
		300	$\hat{\mu}_Y$	0.0236	0.0240	100.00
			$\hat{\mu}_R$	0.0141	0.0141	170.64
			$\hat{\mu}_{Reg}$	0.0129	0.0129	186.50
			$\hat{\mu}_{GRR}$	0.0130	0.0129	186.59
			$\hat{\mu}_{exp1}$	0.0130	0.0129	186.60
			$\hat{\mu}_{exp2}$	0.0103	0.0102	235.40

From Table 4.1 and Table 4.2 we can confirm that suggested generalized regression-cum-exponential estimator is always more efficient than Gupta et al. (2012) regression estimator. Generalized regression-cum-exponential estimator is the most efficient estimator for using one auxiliary variable. Suggested exponential estimator with two auxiliary variables performs better than the estimator with one auxiliary variable, as expected.

Table 4.2: Table 4.1 Continued.

Population			MSE Estimation			
N	ρ_{X1Y} ρ_{X2Y}	n	Estimator	Empirical	Theoretical	PRE
1000	0.8706 0.8428	50	$\hat{\mu}_Y$	0.1198	0.1181	100.00
			$\hat{\mu}_R$	0.0400	0.0395	299.42
			$\hat{\mu}_{Reg}$	0.0287	0.0289	409.40
			$\hat{\mu}_{GRR}$	0.0291	0.0288	409.86
			$\hat{\mu}_{exp1}$	0.0289	0.0288	410.03
			$\hat{\mu}_{exp2}$	0.0078	0.0073	1626.30
		100	$\hat{\mu}_Y$	0.0547	0.0560	100.00
			$\hat{\mu}_R$	0.0188	0.0187	299.42
			$\hat{\mu}_{Reg}$	0.0138	0.0137	409.40
			$\hat{\mu}_{GRR}$	0.0140	0.0137	409.62
			$\hat{\mu}_{exp1}$	0.0139	0.0137	409.70
			$\hat{\mu}_{exp2}$	0.0037	0.0034	1625.73
		200	$\hat{\mu}_Y$	0.0246	0.0249	100.00
			$\hat{\mu}_R$	0.0086	0.0083	299.42
			$\hat{\mu}_{Reg}$	0.0063	0.0061	409.40
			$\hat{\mu}_{GRR}$	0.0063	0.0061	409.49
			$\hat{\mu}_{exp1}$	0.0063	0.0061	409.53
			$\hat{\mu}_{exp2}$	0.0017	0.0015	1625.44
300	$\hat{\mu}_Y$	0.0143	0.0145	100.00		
	$\hat{\mu}_R$	0.0049	0.0048	299.42		
	$\hat{\mu}_{Reg}$	0.0037	0.0035	409.40		
	$\hat{\mu}_{GRR}$	0.0037	0.0035	409.45		
	$\hat{\mu}_{exp1}$	0.0037	0.0035	409.47		
	$\hat{\mu}_{exp2}$	0.0011	0.0009	1625.35		

4.7 Numerical Example

We consider the real population used in Sousa et al (2010) and in Gupta et al. (2012). It is based on the survey on Information and Communication Technologies (ICT) usage in enterprises in 2009 with seat in Portugal (Smilhily and Storm, 2010). This survey intends to promote the development of the national statistical system in the information society and to contribute to a deeper knowledge about the usage of ICT by enterprises. The target population covers all industries with one and more persons employed in the sections of economic activity C (Manufacturing) to N (Administrative and support service activities) and S (Other service activities), from NACE¹ Rev. 2 (Eurostat, 2008). The data are essentially collected using Electronic Data Interchange, applying direct connection between information systems at the respondent and the National Statistics Institute. For some enterprises the paper questionnaire is still used. The questions in the structural business surveys mainly deal with characteristics that can be found in the organisations' annual reports and financial statements, such as employment, turnover and investment.

¹NACE is derived from the French title "Nomenclature générale des Activités économiques dans les Communautés Européennes" (Statistical classification of economic activities in the European Communities).

In our application the study variable Y is the purchase orders in 2010, collected by the ICT survey in that year. This is typically a confidential variable for enterprises, only known from business surveys. The auxiliary variable X is the turnover of each enterprise. This information can be easily obtained from enterprise records available in the public domain, as administrative information. In 2010 the population survey contained approximately 278000 enterprises and we know the value of X for all these enterprises. The purchase orders information was collected in the ICT survey and we have the values of Y for 5336 enterprises (which answered this question in the ICT survey in 2010). For this study, these 5336 enterprises are considered as our population. The scrambling variable S is taken to be a normal random variable with mean equal to zero and standard deviation equal to 10% of the standard deviation of X , that is $\sigma_S = 0.1\sigma_X$. The reported response is given by $Z = Y + S$ (the purchase order value plus a random quantity). The variables Y and X are strongly correlated so we can take advantage of this correlation by using the ratio and regression estimators.

Population Characteristics:

$N = 5336, \rho_{XY} = 0.9632$
$\mu_X = 22.99, \mu_Y = 30.19, \sigma_X = 172.09, \sigma_Y = 138.65$ (in millions of Euros)
and $\beta_{YX} = 0.7763$

We use the following samples sizes in our simulation study: $n = 100, 500, 1000$ and 2000 .

The empirical, theoretical MSE and Percent Relative Efficiency (PRE) values for all estimators are given in Table 4.3.

Table 4.3: *MSE* and *PRE* for the ratio estimator ($\hat{\mu}_R$), the regression estimator ($\hat{\mu}_{Reg}$), the generalized regression-cum-ratio estimator ($\hat{\mu}_{GRR}$) and the exponential estimator ($\hat{\mu}_{exp1}$) relative to the RRT mean estimator.

Population			MSE Estimation			
<i>N</i>	ρ_{XY}	<i>n</i>	Estimator	Empirical	Theoretical	<i>PRE</i>
5336	0.9636	100	$\hat{\mu}_Y$	196.8002	191.5088	100.00
			$\hat{\mu}_R$	11.3683	16.4393	1164.94
			$\hat{\mu}_{Reg}$	16.7963	16.3768	1169.39
			$\hat{\mu}_{GRR}$	11.3909	15.6644	1222.58
			$\hat{\mu}_{exp1}$	12.3339	13.1849	1452.49
		500	$\hat{\mu}_Y$	34.6507	35.3757	100.00
			$\hat{\mu}_R$	2.7259	3.0367	1164.94
			$\hat{\mu}_{Reg}$	3.0170	3.0252	1169.39
			$\hat{\mu}_{GRR}$	2.7509	3.0069	1176.50
			$\hat{\mu}_{exp1}$	2.8631	2.9173	1212.61
		1000	$\hat{\mu}_Y$	15.7543	15.8591	100.00
			$\hat{\mu}_R$	1.3092	1.3614	1164.94
			$\hat{\mu}_{Reg}$	1.3592	1.3562	1169.39
			$\hat{\mu}_{GRR}$	1.3175	1.3526	1172.47
			$\hat{\mu}_{exp1}$	1.3381	1.3345	1188.38
2000	$\hat{\mu}_Y$	6.2451	6.1008	100.00		
	$\hat{\mu}_R$	0.5691	0.5237	1164.94		
	$\hat{\mu}_{Reg}$	0.5573	0.5217	1169.39		
	$\hat{\mu}_{GRR}$	0.5718	0.5212	1170.55		
	$\hat{\mu}_{exp1}$	0.5673	0.5185	1176.62		

From Table 4.3 we can say that generalized regression-cum-exponential estimator has the largest *PRE*.

4.8 Conclusions

This paper proposed type estimators using non-sensitive one or two auxiliary variables to improve the efficiency of RRT estimators of mean. The expression for *Bias* and *MSE* are derived. We found that the proposed exponential type estimators are more efficient than the existing estimators in literature. These results are also supported with a simulation study and using a real data.

References

- BAHL, S. & TUTEJA, R. K. 1991. Ratio and product type exponential estimators. *Information and Optimization Sciences*, 12(1), 159-163.
- EICHHORN, B. H. & HAYRE, L. S. 1983. Scrambled randomized response methods for obtaining sensitive quantitative data. *Journal of Statistical Planning and Inference*, 7, 307-316.

- EUROSTAT. 2008. NACE Rev. 2 - Statistical classification of economic activities in the European Community. *Official Publications of the European Communities*, 112-285 and 306-311.
- GROVER, L.K. 2010. A correction note on improvement in variance estimation using auxiliary information. *Communications in Statistics - Theory and Methods*, 39, 753-764.
- GUPTA, S. N., GUPTA, B. C. & SINGH, S. 2002. Estimation of sensitivity level of personal interview survey questions. *Journal of Statistical Planning and Inference*, 100, 239-247.
- GUPTA, S. & SHABBIR, J. 2004. Sensitivity estimation for personal interview survey questions. *Statistica*, 64, 643-653.
- GUPTA, S., SHABBIR, J. & SEHRA, S. 2010. Mean and sensitivity estimation in optional randomized response models. *Journal of Statistical Planning and Inference*, 140(10), 2870-2874.
- GUPTA, S., SHABBIR, J., SOUSA, R. & REAL, P.C. 2012. Estimation of the Mean of a Sensitive Variable in the Presence of Auxiliary Information. *Communications in Statistics - Theory and Methods*, 41(13-14), 2394-2404.
- KOYUNCU, N. 2012. Efficient estimators of population mean using auxiliary attributes. *Applied Mathematics and Computation*, 218, 10900-10905.
- KOYUNCU, N. & KADILAR, C. 2009. Family of estimators of population mean using two auxiliary variables in stratified random sampling. *Communications in Statistics: Theory and Methods*, 38(14), 2398-2417.
- PERRI, P. F. 2008. Modified randomized devices for Simmons' model. *Model Assisted Statistics and Applications*, 3, 233-239.
- SAHA, A. 2008. A randomized response technique for quantitative data under unequal probability sampling. *Journal of statistical Theory and Practice*, 2(4), 589-596.
- SHABBIR, J. & GUPTA, S. 2007. On improvement in variance estimation using auxiliary information. *Communication in Statistics-Theory and Methods*, 36(12), 2177-2185.
- SMILHILY, M. & STORM, H. 2010. ICT usage in enterprises - 2009. *Eurostat Publications*, Issue 1.
- SOUSA, R., SHABBIR, J. REAL, P. C. & GUPTA, S. 2010. Ratio estimation of the mean of a sensitive variable in the presence of auxiliary information. *Journal of Statistical Theory and Practice*, 4(3), 495-507.
- UPADHYAYA, L.N. & SINGH, H.P. 2006. Almost unbiased ratio and product-type estimators of finite population variance in sample surveys. *Statistics in Transition*, 7(5), 1087-1096.

WU, J-W, TIAN, G-L & TANG, M-L. 2008. Two new models for survey sampling with sensitive characteristic: design and analysis. *Metrika*, 67, 251-263.

Appendix C - R Routines

Listing 4.1: R Code for Simulation Study of Proposed Estimator in Chapter 4

```

1
2 proj_exponential <- function (N, sigma, mu)
3 {
4
5   #Generation of a bivariate normal population
6   data_yx <- mvrnorm(N, mu, sigma)
7
8   #Study variable
9   Y <- data_yx[,1]
10  #Auxiliary variable, correlated with Y
11  X <- data_yx[,2]
12
13  #Coefficient of correlation between Y and X
14  Ro_YX <- cor(Y,X)
15
16  #Scrambling variable independent of Y and X, with mean=0
17  S <- rnorm(N,mean=0,sd=0.1*sd(X))
18
19  #Scrambled response
20  Z <- Y+S
21
22  #Coefficient of correlation between Z and X
23  Ro_ZX <- Ro_YX/sqrt(1+(var(S)/var(Y)))
24
25  #population
26  univ <- data.frame(cbind(Y=Y,S=S,Z=Z,X=X,NRAND=runif(N)))
27  univ <- univ[order(univ$NRAND),]
28
29  #Mean of Y
30  mz <- mean(univ$Z)
31  mx <- mean(univ$X)
32  my <- mean(univ$Y)
33
34  mu11 <- sum((univ$Z-mz) * (univ$X-mx)) / (N-1)
35  mu12 <- sum((univ$Z-mz) * ((univ$X-mx)^2)) / (N-1)
36  mu02 <- sum((univ$X-mx)^2) / (N-1)
37  mu03 <- sum((univ$X-mx)^3) / (N-1)
38
39  beta_zx <- Ro_YX*(sd(univ$Y)/sd(univ$X))
40
41  #Samples dimension
42  dim_samp <- c(50,100,200,300)
43
44  #Initialize the variables...
45
46  for (i in 1:length(dim_samp))

```

```

47 {
48   #sample dimension
49   n <- dim_samp[i]
50   #sample
51   samp <- univ[1:n,]
52   #Sampling rate
53   f <- n/N
54
55   #Ratio
56   R <- mean(univ$X)/mean(samp$X)
57
58   #Ordinary mean
59   est1 <- mean(samp$Z)
60   #Ratio estimator
61   est2 <- mean(samp$Z)*(mx/mean(samp$X))
62   #Regression estimator
63   est3 <- mean(samp$Z)+beta_zx*(mx-mean(samp$X))
64   #Regression-cum-ratio estimator
65   est4 <- (mean(samp$Z)+beta_zx*(mx-mean(samp$X)))*(mx/mean(samp$X))
66
67   #Coefficient of variation
68   c_x <- sd(univ$X)/mx
69   c_y <- sd(univ$Y)/my
70   c2_x <- c_x^2
71   c2_y <- c_y^2
72   c2_z <- c2_y+(var(univ$S)/(my^2))
73   c_z <- sqrt(c2_z)
74
75   A <- (1+((1-f)/n)*(c2_z+c2_x-2*Ro_ZX*c_z*c_x))
76   B <- (((1-f)/n)*(Ro_ZX*c_z*c_x-0.75*c2_x)-2)
77   w1 <- (c2_x/2)*(-B-((1-f)/n)*(c2_x-Ro_ZX*c_z*c_x))
78   / (A*c2_x-((1-f)/n)*((c2_x-Ro_ZX*c_z*c_x)^2))
79   w2 <- (my*c2_x-2*w1*my*(c2_x-Ro_ZX*c_z*c_x))/(2*mx*c2_x)
80
81   #Auxiliar coefficients
82   k1 <- (1-((1-f)*c2_x/n))/(1-((1-f)/n)*(c2_x-c2_z*(1-(Ro_ZX^2))))
83   k2 <- (my/mx)*(1+k1*((Ro_ZX*c_z/c_x)-2))
84
85   #Generalized Regression-cum-ratio Estimator
86   est5 <- (k1*mean(samp$Z)+k2*(mx-mean(samp$X)))*(mx/mean(samp$X))
87
88   #Exponential Type Estimator
89   est6 <- (mean(samp$Z)+beta_zx*(mx-mean(samp$X)))
90   *exp((mx-mean(samp$X))/(mx+mean(samp$X)))
91   #Generalized Exponential Type Estimator
92   est7 <- (w1*mean(samp$Z)+w2*(mx-mean(samp$X)))
93   *exp((mx-mean(samp$X))/(mx+mean(samp$X)))
94
95   #Mean Square Error of 1st estimator (ordinal mean)
96   msel <- ((1-f)/n)*(var(univ$Y)+var(univ$S))

```

```

97
98 #Bias of ratio estimator - 1st degree approximation
99 bias2i <- ((1-f)/n)*my*(c2_x-Ro_ZX*c_z*c_x)
100 #Mean Square Error of ratio estimator - 1st degree approximation
101 mse2i <- ((1-f)/n)*(my^2)*(c2_z+c2_x-2*Ro_ZX*c_z*c_x)
102
103 #Bias of regression estimator - 1st degree approximation
104 bias3i <- -beta_zx*((1-f)/n)*((mu12/mu11)-(mu03/mu02))
105 #Mean Square Error of regression estimator - 1st degree approximation
106 mse3i <- ((1-f)/n)*(my^2)*c2_z*(1-(Ro_ZX^2))
107
108 #Bias of regression-cum-ratio estimator - 1st degree approximation
109 bias4i <- ((1-f)/n)*(my*c2_x-beta_zx*((mu12/mu11)-(mu03/mu02)))
110 #Mean Square Error of regression-cum-ratio estimator
111 #1st degree approximation
112 mse4i <- ((1-f)/n)*(my^2)*(c2_z*(1-(Ro_ZX^2))+c2_x)
113
114 #Bias of genetalized regression-cum-ratio estimator
115 #1st degree approximation
116 bias5i <- (k1-1)*my+k1*my*((1-f)/n)*(c2_x-Ro_ZX*c_z*c_x)
117 +k2*mx*((1-f)/n)*c2_x
118 #Mean Square Error of generalized regression-cum-ratio estimator
119 #1st degree approximation
120 mse5i <- ((k1-1)^2)*(my^2)+(k1^2)*(my^2)*((1-f)/n)
121 *(c2_z+3*c2_x-4*Ro_ZX*c_z*c_x)+(k2^2)*(mx^2)*((1-f)/n)
122 *c2_x-2*k1*(my^2)*((1-f)/n)*(c2_x-Ro_ZX*c_z*c_x)
123 -2*k2*my*mx*((1-f)/n)*c2_x-2*k1*k2*my*mx*((1-f)/n)
124 *(Ro_ZX*c_z*c_x-2*c2_x)
125
126 #Bias of exponential type estimator - 1st degree approximation
127 bias6i <- ((1-f)/n)*(beta_zx*((mu03/mu02)-(mu12/mu11))+(3/8)*my*c2_x)
128 #Mean Square Error of exponential type estimator - 1st degree approximation
129 mse6i <- ((1-f)/n)*(my^2)*(c2_z*(1-(Ro_ZX^2))+0.25*c2_x)
130
131 #Bias of generalized exponential type estimator
132 #1st degree approximation
133 bias7i <- (w1-1)*my+((1-f)/n)
134 *(0.5*w1*my*(0.75*c2_x-Ro_ZX*c_z*c_x)+0.5*w2*mx*c2_x)
135 #Mean Square Error of generalized exponential type estimator
136 #1st degree approximation
137 mse7i <- (mse3i/(1+(mse3i/(my^2))))-(((1-f)/n)*c2_x*(mse3i
138 +((1-f)/n)*(1/16)*c2_x*(my^2)))/(4*(1+mse3i/(my^2)))
139
140 #Empirical results
141 #Simulation of 5000 replicas of estimates
142 ...
143
144 #Results
145 res <- rbind(res,c(N,n,Ro_YX,Ro_ZX,R,
146 c_x,c_y,c_z,k1,k2,w1,w2,

```

```

147         mx,my,mz,
148         med_est1,med_est2,med_est3,med_est4,
149         med_est5,med_est6,med_est7,
150         bias2i,bias3i,bias4i,bias5i,
151         bias6i,bias7i,
152         emp_mse1,mse1,emp_mse2,mse2i,
153         emp_mse3,mse3i,emp_mse4,mse4i,
154         emp_mse5,mse5i,emp_mse6,mse6i,
155         emp_mse7,mse7i)
156     }
157     colnames(res) <- c("N", "n", "RhoXY", "RhoZX", "R",
158                      "Cx", "Cy", "Cz", "k1", "k2", "w1", "w2",
159                      "mX", "mY", "mZ",
160                      "Est1", "Est2", "Est3", "Est4",
161                      "Est5", "Est6", "Est7",
162                      "BIAS2I", "BIAS3I", "BIAS4I", "BIAS5I",
163                      "BIAS6I", "BIAS7I",
164                      "EMP_MSE1", "MSE1", "EMP_MSE2", "MSE2I",
165                      "EMP_MSE3", "MSE3I", "EMP_MSE4", "MSE4I",
166                      "EMP_MSE5", "MSE5I", "EMP_MSE6", "MSE6I",
167                      "EMP_MSE7", "MSE7I")
168     return(res)
169 }
170
171 #Package for generation
172 require(MASS)
173 N <- 1000
174
175 #Parameters
176 sigma1 <- matrix(c(9,1.9,1.9,4),2,2)
177 sigma2 <- matrix(c(10,3,3,2),2,2)
178 sigma3 <- matrix(c(6,3,3,2),2,2)
179 mu <- c(2,2)
180
181 res <- NULL
182 for (i in 1:length(N))
183 {
184     res <- rbind(res,proj_exponential(N[i],sigma1,mu))
185     res <- rbind(res,proj_exponential(N[i],sigma2,mu))
186     res <- rbind(res,proj_exponential(N[i],sigma3,mu))
187 }
188 write.table(res,"chapter4_ss_results.txt",sep="\t",dec="," , row.names=FALSE)

```

Listing 4.2: R Code for Numerical Example of Proposed Estimator in Chapter 4

```

1
2 proj_exponential_real <- function(Y,X,N)
3 {
4   #Coefficient of correlation between Y and X
5   Ro_YX <- cor(Y,X)
6
7   #Scrambling variable independent of Y and X, with mean=0
8   S <- rnorm(N,mean=0,sd=sd(X)*0.1)
9   #Scrambled response
10  Z <- Y+S
11
12  #Coefficient of correlation between Z and X
13  Ro_ZX <- Ro_YX/sqrt(1+(var(S)/var(Y)))
14
15  #population
16  univ <- data.frame(cbind(Y=Y,S=S,Z=Z,X=X,NRAND=runif(N)))
17  univ <- univ[order(univ$NRAND),]
18
19  #Mean of Y
20  mz <- mean(univ$Z)
21  mx <- mean(univ$X)
22  my <- mean(univ$Y)
23  ms <- mean(univ$S)
24
25  mu11 <- sum((univ$Z-mz)*(univ$X-mx))/(N-1)
26  mu12 <- sum((univ$Z-mz)*((univ$X-mx)^2))/(N-1)
27  mu02 <- sum((univ$X-mx)^2)/(N-1)
28  mu03 <- sum((univ$X-mx)^3)/(N-1)
29
30  beta_zx <- Ro_YX*(sd(univ$Y)/sd(univ$X))
31
32  #Samples dimension
33  dim_samp <- c(100,500,1000,2000)
34
35  #Initialize the variables...
36
37  for (i in 1:length(dim_samp))
38  {
39    #sample dimension
40    n <- dim_samp[i]
41    #sample
42    samp <- univ[1:n,]
43    #Sampling rate
44    f <- n/N
45
46    #Ratio
47    R <- mean(univ$X)/mean(samp$X)

```

```

48
49 #Ordinary meam
50 est1 <- mean(samp$Z)
51 #Ratio estimator
52 est2 <- mean(samp$Z) * (mx/mean(samp$X))
53 #Regression estimator
54 est3 <- mean(samp$Z) + beta_zx * (mx - mean(samp$X))
55 #Regression-cum-ratio estimator
56 est4 <- (mean(samp$Z) + beta_zx * (mx - mean(samp$X))) * (mx/mean(samp$X))
57
58
59 #Coefficient of variation
60 c_x <- sd(univ$X) / mx
61 c_y <- sd(univ$Y) / my
62 c2_x <- c_x^2
63 c2_y <- c_y^2
64 c2_z <- c2_y + (var(univ$S) / (my^2))
65 c_z <- sqrt(c2_z)
66
67 A <- (1 + ((1-f)/n) * (c2_z + c2_x - 2 * Ro_ZX * c_z * c_x))
68 B <- (((1-f)/n) * (Ro_ZX * c_z * c_x - 0.75 * c2_x) - 2)
69 w1 <- (c2_x / 2) * (-B - ((1-f)/n) * (c2_x - Ro_ZX * c_z * c_x))
70 / (A * c2_x - ((1-f)/n) * ((c2_x - Ro_ZX * c_z * c_x)^2))
71 w2 <- (my * c2_x - 2 * w1 * my * (c2_x - Ro_ZX * c_z * c_x)) / (2 * mx * c2_x)
72
73 #Auxiliar coefficients
74 k1 <- (1 - ((1-f) * c2_x / n)) / (1 - ((1-f)/n) * (c2_x - c2_z * (1 - (Ro_ZX^2))))
75 k2 <- (my / mx) * (1 + k1 * ((Ro_ZX * c_z / c_x) - 2))
76
77 #Generalized Regression-cum-ratio Estimator
78 est5 <- (k1 * mean(samp$Z) + k2 * (mx - mean(samp$X))) * (mx/mean(samp$X))
79
80 #Exponential Type Estimator
81 est6 <- (mean(samp$Z) + beta_zx * (mx - mean(samp$X)))
82 * exp((mx - mean(samp$X)) / (mx + mean(samp$X)))
83 #Generalized Exponential Type Estimator
84 est7 <- (w1 * mean(samp$Z) + w2 * (mx - mean(samp$X)))
85 * exp((mx - mean(samp$X)) / (mx + mean(samp$X)))
86
87 #Mean Square Error of 1st estimator (ordinal mean)
88 mse1 <- ((1-f)/n) * (var(univ$Y) + var(univ$S))
89
90 #Bias of ratio estimator - 1st degree approximation
91 bias2i <- ((1-f)/n) * my * (c2_x - Ro_ZX * c_z * c_x)
92 #Mean Square Error of ratio estimator - 1st degree approximation
93 mse2i <- ((1-f)/n) * (my^2) * (c2_z + c2_x - 2 * Ro_ZX * c_z * c_x)
94
95 #Bias of regression estimator - 1st degree approximation
96 bias3i <- -beta_zx * ((1-f)/n) * ((mu12/mu11) - (mu03/mu02))
97 #Mean Square Error of regression estimator - 1st degree approximation

```

```

98 #mse3i <- ((1-f)/n)*var(univ$Y)*((1+(var(univ$S)/var(univ$Y)))-(Ro_YX^2))
99 mse3i <- ((1-f)/n)*(my^2)*c2_z*(1-(Ro_ZX^2))
100
101 #Bias of regression-cum-ratio estimator - 1st degree approximation
102 bias4i <- ((1-f)/n)*(my*c2_x-beta_zx*((mu12/mu11)-(mu03/mu02)))
103 #Mean Square Error of regression-cum-ratio estimator
104 #1st degree approximation
105 mse4i <- ((1-f)/n)*(my^2)*(c2_z*(1-(Ro_ZX^2))+c2_x)
106
107 #Bias of genetalized regression-cum-ratio estimator
108 #1st degree approximation
109 bias5i <- (k1-1)*my+k1*my*((1-f)/n)*(c2_x-Ro_ZX*c_z*c_x)
110 +k2*mx*((1-f)/n)*c2_x
111 #Mean Square Error of generalized regression-cum-ratio estimator
112 #1st degree approximation
113 mse5i <- ((k1-1)^2)*(my^2)+(k1^2)*(my^2)*((1-f)/n)
114 *(c2_z+3*c2_x-4*Ro_ZX*c_z*c_x)+(k2^2)*(mx^2)
115 *((1-f)/n)*c2_x-2*k1*(my^2)*((1-f)/n)
116 *(c2_x-Ro_ZX*c_z*c_x)-2*k2*my*mx*((1-f)/n)
117 *c2_x-2*k1*k2*my*mx*((1-f)/n)*(Ro_ZX*c_z*c_x-2*c2_x)
118
119 #Bias of exponential type estimator - 1st degree approximation
120 bias6i <- ((1-f)/n)*(beta_zx*((mu03/mu02)-(mu12/mu11)))+(3/8)*my*c2_x)
121 #Mean Square Error of exponential type estimator - 1st degree approximation
122 mse6i <- ((1-f)/n)*(my^2)*(c2_z*(1-(Ro_ZX^2))+0.25*c2_x)
123
124 #Bias of generalized exponential type estimator
125 #1st degree approximation
126 bias7i <- (w1-1)*my+((1-f)/n)*(0.5*w1*my*(0.75*c2_x-Ro_ZX*c_z*c_x)
127 +0.5*w2*mx*c2_x)
128 #Mean Square Error of generalized exponential type estimator
129 #1st degree approximation
130 mse7i <- (mse5i/(1+(mse5i/(my^2))))
131 -(((1-f)/n)*c2_x*(mse5i+((1-f)/n)*(1/16)
132 *c2_x*(my^2)))/(4*(1+mse5i/(my^2)))
133
134 #Empirical results
135 #Simulation of 5000 replicas of estimates
136 ...
137
138 #Results
139 res <- rbind(res,c(N,n,Ro_YX,Ro_ZX,R,
140 c_x,c_y,c_z,k1,k2,w1,w2,
141 mx,my,mz,
142 med_est1,med_est2,med_est3,med_est4,
143 med_est5,med_est6,med_est7,
144 bias2i,bias3i,bias4i,bias5i,
145 bias6i,bias7i,
146 emp_mse1,mse1,emp_mse2,mse2i,
147 emp_mse3,mse3i,emp_mse4,mse4i,

```

```
148         emp_mse5,mse5i,emp_mse6,mse6i,
149         emp_mse7,mse7i))
150     }
151     colnames(res) <- c("N","n","RhoXY","RhoZX","R",
152         "Cx","Cy","Cz","k1","k2","w1","w2",
153         "mX","mY","mZ",
154         "Est1","Est2","Est3","Est4",
155         "Est5","Est6","Est7",
156         "BIAS2I","BIAS3I","BIAS4I","BIAS5I",
157         "BIAS6I","BIAS7I",
158         "EMP_MSE1","MSE1","EMP_MSE2","MSE2I",
159         "EMP_MSE3","MSE3I","EMP_MSE4","MSE4I",
160         "EMP_MSE5","MSE5I","EMP_MSE6","MSE6I",
161         "EMP_MSE7","MSE7I")
162     return(res)
163 }
164
165 #Package for generation
166 require(MASS)
167
168 #Import data
169 data_yx <- read.table("IUTICE10.txt",sep="\t",dec="," ,header = T)
170 #Study variable (purchase, millions of euros)
171 Y <- data_yx[,3]
172 #Auxiliary variable, correlated with Y (turnover, millions of euros)
173 X <- data_yx[,2]
174
175 #Data application
176 N <- dim(data_yx)[1]
177 res <- proj_exponential_real(Y,X,N)
178
179 #Export data
180 write.table(res,"chapter4_ne_results.txt",sep="\t",dec="," ,row.names=FALSE)
```




Improved Exponential Type Estimators of the Mean of a Sensitive Variable in the Presence of Non-Sensitive Auxiliary Information

Abstract

Recently Koyuncu et al. (2013) proposed an exponential type estimator to improve the efficiency of mean estimator based on Randomized Response Technique (RRT). In this paper, we propose an improved exponential type estimator which is more efficient than the Koyuncu et al. (2013) estimator, which in turn was shown to be more efficient than the usual mean estimator, ratio estimator, regression estimator, and the Gupta et al. (2012) estimator. Under simple random sampling without replacement (SRSWOR) scheme, *Bias* and Mean Square Error (*MSE*) expressions for the proposed estimator are obtained up to first order of approximation and comparisons are made with the Koyuncu et al. (2013) estimator. A simulation study is used to observe the performances of these two estimators. Theoretical findings are also supported by a numerical example with real data.

Submitted as: GUPTA, S., SHABBIR, J., SOUSA, R. & REAL, P. C. 2013. Improved exponential type estimators of the mean of a sensitive variable in the presence of non-sensitive auxiliary information.

5.1 Introduction

This study proposes an improved exponential type estimator for estimating the population mean of a sensitive variable when information about a non-sensitive auxiliary variable is available. A common problem in conducting a statistical sample survey is that of response bias in the face of sensitive questions. Warner (1965) introduced the Randomized Response Technique (RRT) in order to solve this problem. Our main purpose in this study is to improve the mean estimation of a sensitive variable based on a RRT when some non-sensitive auxiliary information is available.

Many authors such as Kadilar and Cingi (2004), Kadilar et al. (2007), Shabbir and Gupta (2007, 2010) and Nangsue (2009) have presented ratio and regression estimators when both the study variable and the auxiliary variable are directly observable.

In this study we propose an exponential type estimator for the mean of a sensitive variable using known information on a correlated but non-sensitive auxiliary variable. The proposed estimator performs better than the recently introduced estimator by Koyuncu et al. (2013) which was shown to outperform many existing estimators of this type.

5.2 Terminology

Consider a finite population with N units $U = (U_1, U_2, \dots, U_N)$ from which a sample of size n is drawn using simple random sampling without replacement (SRSWOR). Let Y be the study variable, a sensitive variable which cannot be observed directly due to respondent bias. Let X be the non-sensitive auxiliary variable which is correlated with Y . Let S be a scrambling variable independent of Y and X . The respondent is asked to report a scrambled response for Y given by $Z = Y + S$ but is asked to provide a true response for X . Let (\bar{y}, \bar{x}) be the sample means corresponding to (\bar{Y}, \bar{X}) , the population means of Y and X , respectively. Consider \bar{Z} to be the population mean of the scrambled variable Z .

Let S_x^2 and s_x^2 respectively be the population variance and the sample variance of X . On the other hand, S_{zx}^2 and s_{zx}^2 are the population covariance and the sample covariance between Z and X , respectively.

To obtain the *Bias* and *MSE* expressions, let us define $e_0 = \frac{\bar{z} - \bar{Z}}{\bar{Z}}$, $e_1 = \frac{\bar{x} - \bar{X}}{\bar{X}}$, $e_2 = \frac{s_x^2 - S_x^2}{s_x^2}$ and $e_3 = \frac{s_{zx} - S_{zx}}{s_{zx}}$ such that $E(e_i) = 0$, $i = 0, 1, 2, 3$. To first degree of approximations, we have:

$$E(e_0^2) = \lambda C_z^2 = v_{20}, E(e_1^2) = \lambda C_x^2 = v_{02}, E(e_0 e_1) = \lambda C_{zx} = \lambda \rho_{zx} C_z C_x = v_{11},$$

$$E(e_1 e_2) = \lambda \frac{\mu_{03}}{\bar{X} \mu_{02}}, E(e_1 e_3) = \lambda \frac{\mu_{12}}{\bar{X} \mu_{11}},$$

where $\lambda = \frac{1-f}{n}$, $f = n/N$, $C_{zx} = \rho_{zx}C_zC_x$ and $\mu_{rs} = \frac{1}{N-1} \sum_{i=1}^N (z_i - \bar{Z})^r (x_i - \bar{X})^s$.

5.3 Difference-cum-exponential Estimator (Koyuncu et al., 2013)

Recently Koyuncu et al. (2013) have suggested a combination of the difference estimator and the exponential estimator with some gain in the efficiency. This estimator is given by

$$\hat{\mu}_{DE} = [w_1\bar{z} + w_2(\bar{X} - \bar{x})] \exp\left(\frac{\bar{X} - \bar{x}}{\bar{X} + \bar{x}}\right), \quad (5.1)$$

where w_1 and w_2 are constants.

The *Bias* and *MSE* of $\hat{\mu}_{DE}$, up to first degree of approximation, at optimum values

$$w_{1(opt)} = \frac{1 - (\lambda C_x^2/8)}{1 + \lambda C_z^2(1 - \rho_{zx}^2)} \quad \text{and} \quad w_{2(opt)} = \frac{\bar{Y}}{\bar{X}} \left\{ \frac{1}{2} - w_{1(opt)} \left(1 - \frac{\rho_{zx}C_z}{C_x}\right) \right\}$$

are given by

$$Bias(\hat{\mu}_{DE}) \cong (w_{1(opt)} - 1)\bar{Y} + w_{1(opt)}\bar{Y}\lambda \left\{ \frac{3}{8}C_x^2 - \frac{1}{2}\rho_{zx}C_zC_x \right\} + w_{2(opt)}\bar{X}\lambda C_x^2, \quad (5.2)$$

and

$$MSE(\hat{\mu}_{DE})_{opt} \cong \bar{Y}^2 \left\{ \left(1 - \frac{1}{4}\lambda C_x^2\right) - \frac{\left(1 - \frac{1}{8}\lambda C_x^2\right)^2}{1 + \lambda C_z^2(1 - \rho_{zx}^2)} \right\}$$

or

$$MSE(\hat{\mu}_{DE})_{opt} \cong \bar{Y}^2 \left\{ \left(1 - \frac{1}{4}v_{02}\right) - \frac{v_{02}(8 - v_{02})^2}{64(v_{02} + v_{20}v_{02} - v_{11}^2)} \right\}. \quad (5.3)$$

It is shown in Koyuncu et al. (2013) that this estimator is better than all the other similar estimators such as Sousa et al. (2010) and Gupta et al. (2012).

5.4 Proposed estimator

The combined product estimators have shown advantages in efficiency spite of being more biased than the traditional ratio or regression estimators (Koyuncu et al., 2013). Motivated by this fact, we propose a change in the difference-cum-exponential estimator in (5.1) so that when the sample mean (\bar{x}) of the auxiliary variable X is close to population mean (\bar{X}), the expected value of the proposed estimator is closer to the variable of interest Y . So, our proposed estimator is an improved exponential estimator as an modified version of the difference-cum-exponential estimator in (5.1) and is given by the following

expression

$$\hat{\mu}_{IE} = [d_1\bar{z} + d_2] \exp\left(\frac{\bar{X} - \bar{x}}{\bar{X} + \bar{x}}\right), \quad (5.4)$$

where d_1 and d_2 are constants.

Using Taylor's approximation and retaining terms of order up to 2, (5.4) can be rewritten as

$$\hat{\mu}_{IE} - \bar{Z} \cong [(d_1 - 1)\bar{Z} + d_1\bar{Z}e_0 + d_2] \left\{ 1 - \frac{1}{2}e_1 + \frac{3}{8}e_1^2 \right\}. \quad (5.5)$$

Recognizing that $\bar{Z} = \bar{Y}$, the optimum *Bias* and *MSE* of $\hat{\mu}_{IE}$, to first degree of approximation, are given by

$$Bias(\hat{\mu}_{IE}) \cong (d_1 - 1)\bar{Y} + d_1\bar{Y} \left(\frac{3}{8}v_{02} - \frac{1}{2}v_{11} \right) + d_2 \left(1 + \frac{3}{8}v_{02} \right), \quad (5.6)$$

and

$$MSE(\hat{\mu}_{IE}) \cong d_1^2\bar{Y}^2A + d_2^2B - 2d_1\bar{Y}^2C - 2d_2\bar{Y}D + 2d_1d_2\bar{Y}E + \bar{Y}^2, \quad (5.7)$$

where $A = 1 + v_{20} + v_{02} - 2v_{11}$, $B = 1 + v_{02}$, $C = 1 + \frac{3}{8}v_{02} - \frac{1}{2}v_{11}$, $D = 1 + \frac{3}{8}v_{02}$, $E = 1 + v_{02} - v_{11}$.

Using (5.7), the optimum values are

$$d_{1(opt)} = \frac{BC - DE}{AB - E^2},$$

and

$$d_{2(opt)} = \frac{\bar{Y}(AD - CE)}{AB - E^2}.$$

Considering the *MSE* at optimum values we get

$$MSE(\hat{\mu}_{IE})_{opt} \cong \bar{Y}^2 \left[1 - \frac{BC^2 + AD^2 - 2CDE}{AB - E^2} \right]$$

or

$$MSE(\hat{\mu}_{IE})_{opt} \cong \bar{Y}^2 \left[1 - \frac{v_{20} + \frac{3}{4}v_{20}v_{02} (1 - \rho_{zx}^2) \left(1 + \frac{3}{16}v_{02} \right) + \frac{1}{64}v_{02}v_{11}^2}{v_{20} + v_{02}v_{20} (1 - \rho_{zx}^2)} \right], \quad (5.8)$$

where $\rho_{zx} = \frac{v_{11}}{\sqrt{v_{20}}\sqrt{v_{02}}}$.

Comparing the *MSE* of this estimator to the *MSE* of difference-cum-exponential estimator given in (5.3), we note that the proposed estimator will be more efficient if

$$MSE(\hat{\mu}_{IE})_{opt} < MSE(\hat{\mu}_{DE})_{opt}.$$

This will be so if

$$\frac{64v_{20} - 48v_{11}^2 - 8v_{02}v_{11}^2 + 48v_{20}v_{02} + 9v_{20}v_{02}^2}{64(v_{20} + v_{02}v_{20} - v_{11}^2)} - \frac{v_{02}(8 - v_{02})^2}{64(v_{02} + v_{02}v_{20} - v_{11}^2)} - \frac{1}{4}v_{02} > 0$$

or if

$$\frac{\left[\begin{array}{l} 64v_{20}^2v_{02} + 32v_{20}^2v_{02}^2 + 8v_{20}v_{02}^3 - 7v_{20}^2v_{02}^3 - v_{20}v_{02}^4 + 15v_{02}^2v_{20}v_{11}^2 - 80v_{20}v_{02}v_{11}^2 \\ -16v_{20}v_{02}^2 - 64v_{20}v_{11}^2 - 8v_{02}^2v_{11}^2 - 8v_{02}v_{11}^4 + v_{02}^3v_{11}^2 + 48v_{11}^4 + 16v_{02}v_{11}^2 \end{array} \right]}{M} > 0,$$

where $M = 64 \left[\{v_{20} + v_{20}v_{02}(1 - \rho_{zx}^2)\} \{v_{02} + v_{20}v_{02}(1 - \rho_{zx}^2)\} \right]$,

or if

$$\frac{v_{20}v_{02}(1 - \rho_{zx}^2) \left[8(8v_{20} + v_{02}v_{11}^2 - 6v_{11}^2) - v_{02}(4 - v_{02})^2 + v_{20}v_{02}(32 - 7v_{02}) \right]}{64 \left[\{v_{20} + v_{20}v_{02}(1 - \rho_{zx}^2)\} \{v_{02} + v_{20}v_{02}(1 - \rho_{zx}^2)\} \right]} > 0. \tag{5.9}$$

The above condition is likely to be true if numerator is positive.

5.5 Simulation Study

In this section, we conduct a simulation study with particular focus on comparing the performance of the proposed combined estimator $\hat{\mu}_P$ to the estimator $\hat{\mu}_{DE}$ suggested by Koyuncu et al. (2013), using the *Bias* and *MSE* results, correct up to first order of approximation.

We consider 2 different bivariate normal distributions for (Y, X) . The scrambling variable S is taken to be a normal variable with mean equal to zero and standard deviation equal to 10% of the standard deviation of X . The reported response is given by $Z = Y + S$. The summary statistics about the bivariate normal populations are given below.

Population Statistics:

I	$N = 1000, \mu_Y = 2, \sigma_Y = \sqrt{10}, \mu_X = 2, \sigma_X = \sqrt{2}, \sigma_{XY} = 3$ and $\rho_{XY} = 0.6708$
II	$N = 1000, \mu_Y = 2, \sigma_Y = \sqrt{6}, \mu_X = 2, \sigma_X = \sqrt{2}, \sigma_{XY} = 3$ and $\rho_{XY} = 0.8660$

We take samples of size $n = 50, 100, 200$ and 300 from each population to compare the results. We estimate the empirical *Bias* and *MSE* using 5000 samples of various sizes from the study populations. The absolute relative bias (*ARB*) is given by

$$\left| \frac{Bias(\mu_\alpha)}{\bar{Y}} \right|,$$

where $\alpha = DE$ and IE .

The empirical and the theoretical results for the two estimators under study are presented in Table 5.1 and Table 5.2, respectively. From these tables we can observe that the proposed estimator shows reduced *Bias* when compared to other estimator.

Table 5.1: Empirical *ARB* for the difference-cum-exponential estimator ($\hat{\mu}_{DE}$) and for the improved exponential estimator ($\hat{\mu}_{IE}$).

Population			Empirical <i>ARB</i>			
<i>N</i>	ρ_{XY}	Estimator	<i>n</i> = 50	<i>n</i> = 100	<i>n</i> = 200	<i>n</i> = 300
1000	0.6867	$\hat{\mu}_{DE}$	0.0267	0.0122	0.0042	0.0012
		$\hat{\mu}_{IE}$	0.0009	0.0007	0.0000	0.0003
	0.8713	$\hat{\mu}_{DE}$	0.0064	0.0027	0.0001	0.0009
		$\hat{\mu}_{IE}$	0.0001	0.0001	0.0002	0.0004

Table 5.2: Theoretical *ARB* for the difference-cum-exponential estimator ($\hat{\mu}_{DE}$) and for the improved exponential estimator ($\hat{\mu}_{IE}$).

Population			Theoretical <i>ARB</i>			
<i>N</i>	ρ_{XY}	Estimator	<i>n</i> = 50	<i>n</i> = 100	<i>n</i> = 200	<i>n</i> = 300
1000	0.6867	$\hat{\mu}_{DE}$	0.0214	0.0103	0.0046	0.0027
		$\hat{\mu}_{IE}$	0.0013	0.0006	0.0003	0.0002
	0.8713	$\hat{\mu}_{DE}$	0.0023	0.0011	0.0005	0.0003
		$\hat{\mu}_{IE}$	0.0006	0.0003	0.0001	0.0001

As expected, the absolute relative bias generally decreases as the sample size increases, however this effect becomes less pronounced when the correlation between *X* and *Y* is higher. Although the proposed estimator is not unbiased, the bias results show a very good performance for this estimator.

Table 5.3 above gives the empirical and theoretical *MSE*'s for the two competing estimators.

The *MSE* values for the proposed estimator are all less than the *MSE* values for the Koyuncu et al. (2013) estimator. The estimators under study get more and more efficient as ρ_{XY} increases. These results were expected from the condition in (5.9).

Table 5.3: Empirical and theoretical MSE for the difference-cum-exponential estimator ($\hat{\mu}_{DE}$) and for the improved exponential estimator ($\hat{\mu}_{IE}$).

Population			MSE Estimation			
N	ρ_{XY}	n	Estimator	Empirical	Theoretical	MSE Condition ¹
1000	0.6867	50	$\hat{\mu}_{DE}$	0.1025	0.1007	0.0253
			$\hat{\mu}_{IE}$	0.0052	0.0050	
		100	$\hat{\mu}_{DE}$	0.0483	0.0484	0.0122
			$\hat{\mu}_{IE}$	0.0024	0.0024	
		200	$\hat{\mu}_{DE}$	0.0217	0.0217	0.0055
			$\hat{\mu}_{IE}$	0.0011	0.0011	
	300	$\hat{\mu}_{DE}$	0.0127	0.0127	0.0032	
		$\hat{\mu}_{IE}$	0.0006	0.0006		
	0.8713	50	$\hat{\mu}_{DE}$	0.0285	0.0283	0.0068
			$\hat{\mu}_{IE}$	0.0024	0.0024	
		100	$\hat{\mu}_{DE}$	0.0132	0.0135	0.0032
			$\hat{\mu}_{IE}$	0.0011	0.0011	
200		$\hat{\mu}_{DE}$	0.0060	0.0060	0.0014	
		$\hat{\mu}_{IE}$	0.0005	0.0005		
300	$\hat{\mu}_{DE}$	0.0035	0.0035	0.0008		
	$\hat{\mu}_{IE}$	0.0003	0.0003			

¹ MSE comparison based on expression (5.9).

5.6 Numerical Example

In this section, we use real data concerning enterprises for the Monthly Economic Survey (MES) in Portugal. The survey is conducted to provide an accurate picture of business trends of enterprises. It provides short-term indicators on a monthly basis compiled for four sectors: industry, retail trade, construction and service sector. The survey results are broken down by branches according to the NACE¹ Rev. 2 (Eurostat, 2008). In this survey the main questions refer to an assessment of recent trends in production, of the current levels of order books and stocks, as well as expectations about production, selling prices and employment. We consider as population the enterprises collected in the 2010 sample which provided results for the industry sector, taking the monthly salaries as study variable and number of employees as auxiliary variable in each enterprise.

Let Y be the monthly salaries amount in 2010 collected by the MES in that year. This is typically a confidential variable for enterprises, only known from business surveys. The auxiliary variable X is the number of employees available from business data registers. The variables Y and X are strongly correlated so we can take advantage of this correlation by using the estimators under study. The MES provided 26980 monthly salary values in 2010, collected for about 2316 enterprises which answered this survey in that

¹NACE is derived from the French title "Nomenclature générale des Activités économiques dans les Communautés Européennes" (Statistical classification of economic activities in the European Communities).

same year. We take these 26980 values as our population. For the RRT part, let S be a normal random variable with mean equal to zero and standard deviation equal to 10% of the standard deviation of X . The reported response is given by $Z = Y + S$ (the salary amount plus a random quantity). The summary statistics about the populations are given below.

Population Characteristics:

$N = 26980, \rho_{XY} = 0.8599$
$\mu_X = 113.91, \mu_Y = 167.18$ (in thousands of Euros)
$\sigma_X = 215.8, \sigma_Y = 501.4$ and $\sigma_{XY} = 93043$

We use the following samples sizes in our simulation study: $n = 1000, 2500, 5000$ and 10000 .

In Tables 5.4 and 5.5 below we present the empirical and the theoretical *ARB* results, respectively, for the difference-cum-exponential estimator ($\hat{\mu}_{DE}$) and for the proposed estimator ($\hat{\mu}_{IE}$).

Table 5.4: Empirical *ARB* for the difference-cum-exponential estimator ($\hat{\mu}_{DE}$) and for the improved exponential estimator ($\hat{\mu}_{IE}$).

Population			Empirical <i>ARB</i>			
N	ρ_{XY}	Estimator	$n = 1000$	$n = 2500$	$n = 5000$	$n = 10000$
26980	0.8599	$\hat{\mu}_{DE}$	0.0025	0.0022	0.0016	0.0012
		$\hat{\mu}_{IE}$	0.0003	0.0004	0.0004	0.0003

Table 5.5: Theoretical *ARB* for the difference-cum-exponential estimator ($\hat{\mu}_{DE}$) and for the improved exponential estimator ($\hat{\mu}_{IE}$).

Population			Theoretical <i>ARB</i>			
N	ρ_{XY}	Estimator	$n = 1000$	$n = 2500$	$n = 5000$	$n = 10000$
26980	0.8599	$\hat{\mu}_{DE}$	0.0008	0.0003	0.0001	0.0001
		$\hat{\mu}_{IE}$	0.0002	0.0001	0.0000	0.0000

The *ARB* results show the good performance for the improved difference-cum-exponential estimator.

The theoretical *MSE* values for both estimators have been obtained using (5.3) and (5.8). These values are given in Table 5.6.

According to the *MSE* results in Table 5.6, the proposed estimator is considerably better than the difference-cum-exponential estimator ($\hat{\mu}_{DE}$). These results are in line with the theoretical findings and the simulation results.

Table 5.6: Empirical and theoretical MSE for the difference-cum-exponential estimator ($\hat{\mu}_{DE}$) and for the improved exponential estimator ($\hat{\mu}_{IE}$).

Population			MSE Estimation			
N	ρ_{XY}	n	Estimator	Empirical	Theoretical	MSE Condition ²
26980	0.8599	1000	$\hat{\mu}_{DE}$	62.58	63.34	0.0205
			$\hat{\mu}_{IE}$	6.28	6.3	
		2500	$\hat{\mu}_{DE}$	24.50	23.92	0.0083
			$\hat{\mu}_{IE}$	2.40	2.38	
		5000	$\hat{\mu}_{DE}$	10.88	10.75	0.0041
			$\hat{\mu}_{IE}$	1.09	1.07	
		10000	$\hat{\mu}_{DE}$	4.19	4.15	0.0020
			$\hat{\mu}_{IE}$	0.41	0.41	

² MSE comparison based on expression (5.9).

5.7 Conclusions

We can conclude from this study that the estimation of the mean of a sensitive variable can be improved by using a correlated non-sensitive auxiliary variable. Our simulation study and the numerical example show that improved difference-cum-exponential estimator can produce further improvement.

In this paper we show that the proposed estimator is more efficient than the difference-cum-exponential estimator recently proposed by Koyuncu et al. (2013), which in turn was better than most of the existing estimators of finite population mean.

References

- EUROSTAT. 2008. NACE Rev. 2 - Statistical classification of economic activities in the European Community. *Official Publications of the European Communities*, 112-285 and 306-311.
- GUPTA, S., SHABBIR, J., SOUSA, R. & REAL, P.C. 2012. Estimation of the Mean of a Sensitive Variable in the Presence of Auxiliary Information. *Communications in Statistics - Theory and Methods*, 41(13-14), 2394-2404.
- KADILAR, C., CANDAN, M. & CINGI, H. 2007. Ratio estimators using robust regression. *Hacetatepe Journal of Mathematics and Statistics*, 36(2), 81-188.
- KADILAR, C. & CINGI, H. 2004. Ratio estimators in simple random sampling. *Applied Mathematics and Computation*, 151, 893-902.
- KOYUNCU, N., GUPTA, S. & SOUSA, R. 2013. Exponential type estimators of the mean

of a sensitive variable in the presence of non-sensitive auxiliary information. *Communications in Statistics - Simulation and Computation*. (accepted).

NANGSUE, N. 2009. Adjusted Ratio and Regression Type Estimators for Estimation of Population Mean when some Observations are missing. *World Academy of Science, Engineering and Technology*, 53, 781-784.

SHABBIR, J. & GUPTA, S. 2007. On improvement in variance estimation using auxiliary information. *Communication in Statistics-Theory and Methods*, 36(12), 2177-2185.

SHABBIR, J. & GUPTA, S. 2010. Estimation of the finite population mean in two-phase sampling when auxiliary variables are attribute. *Hacettepe Journal of Mathematics and Statistics*, 39(1), 121-129.

SOUSA, R., SHABBIR, J. REAL, P. C. & GUPTA, S. 2010. Ratio estimation of the mean of a sensitive variable in the presence of auxiliary information. *Journal of Statistical Theory and Practice*, 4(3), 495-507.

WARNER, S. L. 1965. Randomized response: a survey technique for elimination evasive answer bias. *Journal of American Statistical Association*, 60, 63-69.

Appendix D - R Routines

Listing 5.1: R Code for Simulation Study of Proposed Estimator in Chapter 5

```

1
2 proj_improved_exp <- function(N, sigma, mu)
3 {
4
5   #Generation of a bivariate normal population
6   data_yx <- mvrnorm(N, mu, sigma)
7
8   #Study variable
9   Y <- data_yx[,1]
10  #Auxiliary variable, correlated with Y
11  X <- data_yx[,2]
12
13  #Coefficient of correlation between Y and X
14  Ro_YX <- cor(Y,X)
15
16  #Scrambling variable independent of Y and X, with mean=0
17  S <- rnorm(N,mean=0,sd=0.1*sd(X))
18  #Scrambled response
19  Z <- Y+S
20
21  #Coefficient of correlation between Z and X
22  Ro_ZX <- Ro_YX/sqrt(1+(var(S)/var(Y)))
23
24  #population
25  univ <- data.frame(cbind(Y=Y,S=S,Z=Z,X=X,NRAND=runif(N)))
26  univ <- univ[order(univ$NRAND),]
27
28  #Mean of Y
29  mz <- mean(univ$Z)
30  mx <- mean(univ$X)
31  my <- mean(univ$Y)
32
33  mu11 <- sum((univ$Z-mz) * (univ$X-mx)) / (N-1)
34  mu12 <- sum((univ$Z-mz) * ((univ$X-mx)^2)) / (N-1)
35  mu02 <- sum((univ$X-mx)^2) / (N-1)
36  mu03 <- sum((univ$X-mx)^3) / (N-1)
37
38  beta_zx <- Ro_YX*(sd(univ$Y)/sd(univ$X))
39
40  #Samples dimension
41  dim_samp <- c(50,100,200,300)
42
43  #Initialize the variables...
44
45  for (i in 1:length(dim_samp))
46  {

```

```

47  #sample dimension
48  n <- dim_samp[i]
49  #sample
50  samp <- univ[1:n,]
51  #Sampling rate
52  f <- n/N
53
54  #Coefficient of variation
55  c_x <- sd(univ$X)/mx
56  c_y <- sd(univ$Y)/my
57  c2_x <- c_x^2
58  c2_y <- c_y^2
59  c2_z <- c2_y+(var(univ$S)/(my^2))
60  c_z <- sqrt(c2_z)
61
62  l <- (1-f)/n
63
64  #Difference-cum-exponential type Estimator
65  w1 <- (1-(l*c2_x/8))/(1+l*c2_z*(1-(Ro_ZX^2)))
66  w2 <- (my/mx)*(0.5-w1*(1-(Ro_ZX*c_z/c_x)))
67  est5 <- (w1*mean(samp$Z)+w2*(mx-mean(samp$X)))
68  *exp((mx-mean(samp$X))/(mx+mean(samp$X)))
69
70  #2nd Improved Exponential Estimator
71  A <- 1+l*c2_z+l*c2_x-2*l*Ro_ZX*c_z*c_x
72  B <- 1+l*c2_x
73  C <- 1+(3/8)*l*c2_x-0.5*l*Ro_ZX*c_z*c_x
74  D <- 1+(3/8)*l*c2_x
75  E <- 1+l*c2_x-l*Ro_ZX*c_z*c_x
76  z1 <- (B*C-D*E)/(A*B-(E^2))
77  z2 <- my*(A*D-C*E)/(A*B-(E^2))
78  est6 <- (z1*mean(samp$Z)+z2)*exp((mx-mean(samp$X))/(mx+mean(samp$X)))
79
80  #Bias of Difference-cum-exponential estimator - 1st degree approximation
81  bias5i <- (w1-1)*my+w1*my*l*((3/8)*c2_x-0.5*Ro_ZX*c_z*c_x)
82  +w2*mx*l*c2_x
83  mse5i <- (my^2)*((1-0.25*l*c2_x)-(((1-(1/8)*l*c2_x)^2)
84  /(1+l*c2_z*(1-(Ro_ZX^2)))))
85
86  #Bias of Improved Exponential - 1st degree approximation
87  bias6i <- (z1-1)*my+z1*my*((3/8)*l*c2_x
88  -0.5*l*Ro_ZX*c_z*c_x)+z2*(1+(3/8)*l*c2_x)
89  #Mean Square Error of improved exponential estimator 2
90  #1st degree approximation
91  mse6i <- (my^2)*(1-((B*(C^2)+A*(D^2)-2*C*D*E)/(A*B-(E^2))))
92
93  #Condition to compare Est6(P) with Est8(DE)
94  cond <- (v20*v02*(1-(Ro_ZX^2))*(8*(8*v20+v02*(v11^2)
95  -6*(v11^2))-v02*((4-v02)^2)+v20*v02*(32-7*v02))
96  /(64*(v20+v20*v02*(1-(Ro_ZX^2)))*(v02+v20*v02*(1-(Ro_ZX^2))))

```

```

97
98   #Empirical results
99   #Simulation of 5000 replicas of estimates
100   ...
101
102   #Results
103   res <- rbind(res, c(N, n, Ro_YX, Ro_ZX,
104                     c_x, c_y, c_z, k1, k2, w1, w2,
105                     z1, z2, mx, my, mz,
106                     med_est1, med_est2, med_est3, med_est4,
107                     med_est5, med_est6, med_est7,
108                     bias2i, bias3i, bias4i, bias5i,
109                     bias6i, bias6i,
110                     emp_mse1, mse1, emp_mse2, mse2i,
111                     emp_mse3, mse3i, emp_mse4, mse4i,
112                     emp_mse5, mse5i, emp_mse6, mse6i,
113                     emp_mse7, mse6i,
114                     cond))
115 }
116 colnames(res) <- c("N", "n", "RhoXY", "RhoZX",
117                  "Cx", "Cy", "Cz", "k1", "k2", "w1", "w2",
118                  "z1", "z2", "mX", "mY", "mZ",
119                  "Est1", "Est2", "Est3", "Est4",
120                  "Est5", "Est6", "Est7",
121                  "BIAS2I", "BIAS3I", "BIAS4I", "BIAS5I",
122                  "BIAS6I", "BIAS7I",
123                  "EMP_MSE1", "MSE1", "EMP_MSE2", "MSE2I",
124                  "EMP_MSE3", "MSE3I", "EMP_MSE4", "MSE4I",
125                  "EMP_MSE5", "MSE5I", "EMP_MSE6", "MSE6I",
126                  "EMP_MSE7", "MSE7I",
127                  "COND")
128   return(res)
129 }
130
131 #Package for generation
132 require(MASS)
133
134 N <- 1000
135
136 #Parameters
137 sigma1 <- matrix(c(9, 1.9, 1.9, 4), 2, 2)
138 sigma2 <- matrix(c(10, 3, 3, 2), 2, 2)
139 sigma3 <- matrix(c(6, 3, 3, 2), 2, 2)
140
141 mu <- c(2, 2)
142
143 res <- NULL
144 for (i in 1:length(N))
145 {
146   res <- rbind(res, proj_improved_exp(N[i], sigma1, mu))

```

```
147   res <- rbind(res,proj_improved_exp(N[i],sigma2,mu))
148   res <- rbind(res,proj_improved_exp(N[i],sigma3,mu))
149 }
150 write.table(res,"chapter5_ss_results.txt",sep="\t",dec=" ",row.names=FALSE)
```

Listing 5.2: R Code for Numerical Example of Proposed Estimator in Chapter 5

```

1
2 proj_improved_exp_real <- function(Y,X,N)
3 {
4
5   #Coefficient of correlation between Y and X
6   Ro_YX <- cor(Y,X)
7
8   #Scrambling variable independent of Y and X, with mean=0
9   S <- rnorm(N,mean=0,sd=0.1*sd(X))
10  #Scrambled response
11  Z <- Y+S
12
13  #Coefficient of correlation between Z and X
14  Ro_ZX <- Ro_YX/sqrt(1+(var(S)/var(Y)))
15
16  #population
17  univ <- data.frame(cbind(Y=Y,S=S,Z=Z,X=X,NRAND=runif(N)))
18  univ <- univ[order(univ$NRAND),]
19
20  #Mean of Y
21  mz <- mean(univ$Z)
22  mx <- mean(univ$X)
23  my <- mean(univ$Y)
24
25  mu11 <- sum((univ$Z-mz)*(univ$X-mx))/(N-1)
26  mu12 <- sum((univ$Z-mz)*((univ$X-mx)^2))/(N-1)
27  mu02 <- sum((univ$X-mx)^2)/(N-1)
28  mu03 <- sum((univ$X-mx)^3)/(N-1)
29
30  beta_zx <- Ro_YX*(sd(univ$Y)/sd(univ$X))
31
32  #Samples dimension
33  dim_samp <- c(1000,2500,5000,10000)
34
35  #Initialize the variables...
36
37  for (i in 1:length(dim_samp))
38  {
39    #sample dimension
40    n <- dim_samp[i]
41    #sample
42    samp <- univ[1:n,]
43    #Sampling rate
44    f <- n/N
45
46    #Coefficient of variation
47    c_x <- sd(univ$X)/mx

```

```

48  c_y <- sd(univ$Y)/my
49  c2_x <- c_x^2
50  c2_y <- c_y^2
51  c2_z <- c2_y+(var(univ$S)/(my^2))
52  c_z <- sqrt(c2_z)
53
54  l <- (1-f)/n
55
56  #Difference-cum-exponential type Estimator
57  w1 <- (1-(1*c2_x/8))/(1+1*c2_z*(1-(Ro_ZX^2)))
58  w2 <- (my/mx)*(0.5-w1*(1-(Ro_ZX*c_z/c_x)))
59  est7 <- (w1*mean(samp$Z)+w2*(mx-mean(samp$X)))
60  *exp((mx-mean(samp$X))/(mx+mean(samp$X)))
61
62  #2nd Improved Exponential Estimator
63  A <- 1+1*c2_z+1*c2_x-2*1*Ro_ZX*c_z*c_x
64  B <- 1+1*c2_x
65  C <- 1+(3/8)*1*c2_x-0.5*1*Ro_ZX*c_z*c_x
66  D <- 1+(3/8)*1*c2_x
67  E <- 1+1*c2_x-1*Ro_ZX*c_z*c_x
68  z1 <- (B*C-D*E)/(A*B-(E^2))
69  z2 <- my*(A*D-C*E)/(A*B-(E^2))
70  est8 <- (z1*mean(samp$Z)+z2)*exp((mx-mean(samp$X))/(mx+mean(samp$X)))
71
72  #Bias of generalized exponential type estimator
73  #1st degree approximation
74  bias7i <- (w1-1)*my+w1*my*1*((3/8)*c2_x-0.5*Ro_ZX*c_z*c_x)
75  +w2*mx*1*c2_x
76  mse7i <- (my^2)*((1-0.25*1*c2_x)-(((1-(1/8)*1*c2_x)^2)
77  /((1+1*c2_z*(1-(Ro_ZX^2))))))
78
79  #Bias of improved exponential estimator 2 - 1st degree approximation
80  bias8i <- (z1-1)*my+z1*my*((3/8)*1*c2_x-0.5*1*Ro_ZX*c_z*c_x)
81  +z2*(1+(3/8)*1*c2_x)
82  #Mean Square Error of improved exponential estimator 2
83  #1st degree approximation
84  mse8i <- (my^2)*(1-((B*(C^2)+A*(D^2)-2*C*D*E)/(A*B-(E^2))))
85
86  #Condition to compare Est8(P) with Est8(DE)
87  cond <- ((B*(C^2)+A*(D^2)-2*C*D*E)/(A*B-(E^2)))
88  -(((1-(1/8)*1*c2_x)^2)/(1+1*c2_z*(1-(Ro_ZX^2))))-0.25*1*c2_x
89
90  #Empirical results
91  #Simulation of 5000 replicas of estimates
92  ...
93
94  #Results
95  res <- rbind(res,c(N,n,Ro_YX,Ro_ZX,
96  c_x,c_y,c_z,k1,k2,w1,w2,
97  z1,z2,mx,my,mz,
```



```

98         med_est7,med_est8,
99         bias7i,bias8i,
100        emp_mse7,mse7i,
101        emp_mse8,mse8i,
102        cond))
103    }
104    colnames(res) <- c("N","n","RhoXY","RhoZX",
105                    "Cx","Cy","Cz","k1","k2","w1","w2",
106                    "z1","z2","mX","mY","mZ",
107                    "Est7","Est8",
108                    "BIAS7I","BIAS8I",
109                    "EMP_MSE7","MSE7I",
110                    "EMP_MSE8","MSE8I",
111                    "COND")
112    return(res)
113 }
114
115 #Package for generation
116 require(MASS)
117
118 #Import data
119 data_yx <- read.table("IVNEI2010.txt",sep="\t",dec="," ,header = T)
120 #Study variable (purchase, millions of euros)
121 data_yx <- data_yx[data_yx$MES>=10,]
122 Y <- data_yx[,5]/1000
123 #Auxiliary variable, correlated with Y (turnover, millions of euros)
124 X <- data_yx[:,4]
125
126 #Data application
127 N <- dim(data_yx)[1]
128 res <- proj_improved_exp_real(Y,X,N)
129
130 #Export data
131 write.table(res,"chapter5_ne_results.txt",sep="\t",dec="," ,row.names=FALSE)

```




Improved Mean Estimation of a Sensitive Variable Using Auxiliary Information in Stratified Sampling

Abstract

Sousa et al. (2010) and Gupta et al. (2012) have recently introduced ratio estimator and regression estimators for the mean of a sensitive variable which perform better than the ordinary mean estimator based on a Randomized Response Technique (RRT). In the present study we extend these estimators to the stratified sampling setting.

The performance of the proposed estimators is compared to the exiting estimators both theoretically and through a simulation study. We also apply the proposed estimators to some real data.

Submitted as: SOUSA, R., GUPTA, S., SHABBIR, J. & REAL, P. C. 2013. Improved Mean Estimation of a Sensitive Variable Using Auxiliary Information in Stratified Sampling. *Journal of Statistics and Management Systems*.

6.1 Introduction

The main goal of this paper is to extend the results of Sousa et al. (2010) and Gupta et al. (2012) to the case of stratified sampling. It is assumed that the study variable is sensitive and the auxiliary variable is non-sensitive.

Many authors have presented ratio and regression estimators when both the study variable Y and the auxiliary variable X are directly observable. These include Kadilar and Cingi (2005), Kadilar et al. (2007), Shabbir and Gupta (2007, 2010) and Nangsue (2009). Gupta and Shabbir (2008) have suggested a general class of ratio estimators when the population parameters of the auxiliary variable are known. These estimators have also been extended by Kadilar and Cingi (2003) to stratified random sampling. In an attempt to improve the estimators, Kadilar and Cingi (2005), Shabbir and Gupta (2005, 2006) and Singh and Vishwakarma (2008) have suggested new ratio estimators in stratified random sampling. Koyuncu and Kadilar (2008, 2009) have proposed a family of combined-type estimators in stratified random sampling based on the family of estimators proposed by Khoshnevisan et al. (2007). Recently Koyuncu and Kadilar (2010) have suggested a family of estimators in stratified random sampling following Kadilar and Cingi (2003).

Some studies on estimation of the mean have been submitted with different sampling schemes, such as Sahoo et al. (2009) and Singh and Kumar (2011) in a two-stage sampling scheme and recently by Singh and Solanki (2012) in a systematic sampling design.

This paper suggests a combined ratio estimator and a combined regression estimator of population mean of a sensitive variable using non-sensitive auxiliary information, using Randomized Response Technique (RRT) methodology (Gupta et al., 2002 and 2010; Warner, 1965) in stratified sampling. The Bias and the Mean Square Error (MSE) of the suggested estimators are derived. Both theoretical and empirical findings support the reliability of the present study.

6.2 Terminology

We denote the finite population by $U = \{U_1, U_2, \dots, U_N\}$. Consider a stratified random sample s (Cochran, 1977), selected from U with sampling rate $f = \frac{n}{N}$. The study population is divided into L strata with strata sizes N_h , such that $\sum_{h=1}^L N_h = N$ ($h = 1, \dots, L$).

Let Y be the sensitive study variable which cannot be observed directly. Let X be a non-sensitive auxiliary variable which is strongly correlated with Y . Let S be a scrambling random variable independent of Y and X . The particular values of S are unknown to the interviewer but its distribution is known. The respondent is asked to report an additively scrambled response for Y given by $Z = Y + S$ and also asked to provide a true response for X .

Consider a stratified random sample of size n be drawn from U such that the sample size in the h^{th} stratum is n_h and $\sum_{h=1}^L n_h = n$. Let y_{hi} and x_{hi} respectively be the values of the study variable Y and the auxiliary variable X in the h^{th} stratum, with $i = 1, \dots, n_h$.

Let $\bar{y}_{st} = \sum_{h=1}^L W_h \bar{y}_h$, $\bar{x}_{st} = \sum_{h=1}^L W_h \bar{x}_h$ and $\bar{z}_{st} = \sum_{h=1}^L W_h \bar{z}_h$ be the stratified sample means, where $\bar{y}_h = \frac{1}{n_h} \sum_{i=1}^{n_h} y_{hi}$, $\bar{x}_h = \frac{1}{n_h} \sum_{i=1}^{n_h} x_{hi}$ and $\bar{z}_h = \frac{1}{n_h} \sum_{i=1}^{n_h} z_{hi}$ are the stratum sample means corresponding to population stratum means $\bar{Y}_h = E(Y_h)$, $\bar{X}_h = E(X_h)$ and $\bar{Z}_h = E(Z_h)$ and $W_h = \frac{N_h}{N}$ are the known stratum weights.

To estimate $\bar{Y} = \sum_{h=1}^L W_h \bar{Y}_h$ we assume that $\bar{X} = \sum_{h=1}^L W_h \bar{X}_h$ is known. Let $\bar{Z} = \sum_{h=1}^L W_h \bar{Z}_h$ be the population mean for the scrambled variable Z .

To discuss the properties of different estimators, we define the following error terms. Let $e_{0st} = \frac{\bar{z}_{st} - \bar{Z}}{\bar{Z}}$ and $e_{1st} = \frac{\bar{x}_{st} - \bar{X}}{\bar{X}}$, $e_{2st} = \frac{s_{xst}^2 - S_{xst}^2}{S_{xst}^2}$ and $e_{3st} = \frac{s_{zxt}^2 - S_{zxt}^2}{S_{zxt}^2}$ such that $E(e_{ist}) = 0$, $i = 0, 1, 2, 3$.

6.3 Estimators Review

Below we list some existing mean estimators for simple random sampling.

(i) Ordinary sample mean:

$$\hat{\mu}_y = \bar{z}. \tag{6.1}$$

$$MSE(\hat{\mu}_y) = \frac{1-f}{n} (S_y^2 + S_s^2), \tag{6.2}$$

where $S_y^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{Y})^2$ and $S_s^2 = \frac{1}{N-1} \sum_{i=1}^N (s_i - \bar{S})^2$.

(ii) Sousa et al. (2010) ratio estimator:

$$\hat{\mu}_R = \bar{z} \frac{\bar{X}}{\bar{x}}. \tag{6.3}$$

The *Bias* and *MSE* of $\hat{\mu}_R$ to first degree of approximation are given by

$$Bias(\hat{\mu}_R) \cong \frac{1-f}{n} \bar{Y} (C_x^2 - \rho_{zx} C_z C_x) \tag{6.4}$$

and

$$MSE(\hat{\mu}_R) \cong \frac{1-f}{n} \bar{Y}^2 (C_z^2 + C_x^2 - 2\rho_{zx} C_z C_x), \tag{6.5}$$

where $C_z^2 = C_y^2 + \frac{S_s^2}{\bar{Y}^2}$, $\rho_{zx} = \frac{\rho_{yx}}{\sqrt{1 + \frac{S_s^2}{S_y^2}}}$ and C_z , C_y and C_x are the coefficients of variation of Z , Y and X , respectively.

(iii) Gupta et al. (2012) regression estimator:

$$\hat{\mu}_{Reg} = \bar{z} + \hat{\beta}_{zx} (\bar{X} - \bar{x}), \quad (6.6)$$

where $\hat{\beta}_{zx} = \frac{S_{zx}}{S_x^2} = \frac{S_{yx}}{S_x^2}$ is the sample regression coefficient between Z and X , $S_{yx} = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})$ and $S_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$.

The *Bias* and *MSE* of $\hat{\mu}_{Reg}$ to first degree of approximation, are given by

$$Bias(\hat{\mu}_{Reg}) \cong -\beta_{zx} \left(\frac{1-f}{n} \right) \left\{ \frac{\mu_{12}}{\mu_{11}} - \frac{\mu_{03}}{\mu_{02}} \right\} \quad (6.7)$$

and

$$MSE(\hat{\mu}_{Reg}) \cong \left(\frac{1-f}{n} \right) \bar{Y}^2 C_z^2 (1 - \rho_{zx}^2), \quad (6.8)$$

where $\mu_{rs} = \frac{1}{N-1} \sum_{i=1}^N (z_i - \bar{Z})^r (x_i - \bar{X})^s$.

For a stratified random sample the usual combined sample mean, ignoring the auxiliary information, is given by

$$\hat{\mu}_{Yst} = \bar{z}_{st}, \quad (6.9)$$

which is an unbiased estimator of population mean \bar{Y} .

The *MSE* of $\hat{\mu}_{Yst}$ is given by

$$MSE(\hat{\mu}_{Yst}) = \sum_{h=1}^L W_h^2 \gamma_h \{ S_{yh}^2 + S_{sh}^2 \}, \quad (6.10)$$

where $\gamma_h = \left(\frac{1}{n_h} - \frac{1}{N_h} \right)$, $S_{yh}^2 = \frac{1}{N_h-1} \sum_{i=1}^{N_h} (y_{hi} - \bar{Y}_h)^2$ and $S_{sh}^2 = \frac{1}{N_h-1} \sum_{i=1}^{N_h} (s_{hi} - \bar{S}_h)^2$.

The remainder of the paper is as follows. In Section 6.4, we introduce a combined ratio estimator and compare it to the ordinary mean estimator and to the ratio estimator (Sousa et al., 2010), considering the *MSE* as an accuracy indicator. In Section 6.5, we propose a combined regression estimator and compare it with other estimators as well as with the regression estimator proposed by Gupta et al. (2012). We present an empirical study in Section 6.6 and a numerical example in Section 6.7 to support the proposed methodology. Section 6.8 provides some concluding remarks.

6.4 Proposed combined ratio estimator

We propose the following combined ratio estimator

$$\hat{\mu}_{Rst} = \bar{z}_{st} \left(\frac{\bar{X}}{\bar{x}_{st}} \right). \quad (6.11)$$

Using Taylor's approximation and retaining terms of order up to 2, (6.11) can be rewritten as

$$\hat{\mu}_{Rst} - \bar{Z} \cong \bar{Z} \{e_{0st} - e_{1st} + e_{1st}^2 - e_{0st}e_{1st}\}. \quad (6.12)$$

Under the assumption of bivariate normality (see Sukhatme and Sukhatme, 1984), we have:

$$E(e_{0st}^2) = \sum_{h=1}^L W_h^2 \gamma_h C_{zh}^2, \quad E(e_{1st}^2) = \sum_{h=1}^L W_h^2 \gamma_h C_{xh}^2, \quad E(e_{0st}e_{1st}) = \sum_{h=1}^L W_h^2 \gamma_h C_{zxh},$$

where $C_{zxh} = \rho_{zxh} C_{zh} C_{xh}$, $C_{zh}^2 = C_{yh}^2 + \left(\frac{S_{sh}^2}{\bar{Y}^2}\right)$ and $\rho_{zxh} = \frac{\rho_{yxh}}{\sqrt{1 + \left(\frac{S_{sh}^2}{S_{yh}^2}\right)}}$.

Using $\bar{Z} = \bar{Y}$ in (6.12), the Bias of $\hat{\mu}_{Rst}$ to first degree of approximation is given by

$$Bias(\hat{\mu}_{Rst}) \cong \bar{Y} \sum_{h=1}^L W_h^2 \gamma_h (C_{xh}^2 - C_{zxh}). \quad (6.13)$$

Using (6.12), the MSE of $\hat{\mu}_{Rst}$, correct up to first order of approximation, is given by

$$MSE(\hat{\mu}_{Rst}) = E \{ \hat{\mu}_{Rst} - \bar{Y} \}^2 \cong \bar{Y}^2 E \{ e_{0st} - e_{1st} + e_{1st}^2 - e_{0st}e_{1st} \}^2.$$

So, if an independent simple random sample is drawn in each stratum, we have

$$MSE(\hat{\mu}_{Rst}) \cong \bar{Y}^2 \sum_{h=1}^L W_h^2 \gamma_h \{ C_{zh}^2 + C_{xh}^2 - 2C_{zxh} \}. \quad (6.14)$$

It can be observed that $MSE(\hat{\mu}_{Rst}) < MSE(\hat{\mu}_{Yst})$ if

$$\sum_{h=1}^L W_h^2 \gamma_h C_{zxh} - \frac{1}{2} \sum_{h=1}^L \gamma_h C_{xh}^2 > 0. \quad (6.15)$$

On the other hand, comparing this estimator with different sampling methods $MSE(\hat{\mu}_{Rst}) < MSE(\hat{\mu}_R)$ if

$$\sum_{h=1}^L W_h^2 \gamma_h \{ C_{zh}^2 + C_{xh}^2 - 2\rho_{zxh} C_{zh} C_{xh} \} < \frac{1-f}{n} \{ C_z^2 + C_x^2 - 2\rho_{zx} C_z C_x \}, \quad (6.16)$$

a condition that can be ensured by a suitable stratification.

6.5 Proposed combined regression estimator

Assuming linear relationship between Y and X , we propose the following combined regression estimator for the population mean of Y

$$\hat{\mu}_{Regst} = \bar{z}_{st} + \hat{\beta}_c (\bar{X} - \bar{x}_{st}), \quad (6.17)$$

where $\hat{\beta}_c = \frac{\sum_{h=1}^L W_h^2 \gamma_h S_{zxh}}{\sum_{h=1}^L W_h^2 \gamma_h S_{xh}^2}$ is the sample regression coefficient between Z and X and $Z = Y + S$ is the scrambled response on Y .

Using Taylor's approximation and retaining terms of order up to 2, (6.17) can be rewritten as

$$\hat{\mu}_{Regst} - \bar{Z} \cong \bar{Z} e_{ost} - \beta_c \bar{X} [e_{1st} + e_{1st} e_{3st} - e_{1st} e_{2st}], \quad (6.18)$$

where $\beta_c = \frac{\sum_{h=1}^L W_h^2 \gamma_h S_{zxh}}{\sum_{h=1}^L W_h^2 \gamma_h S_{xh}^2}$ is the population regression coefficient between Z on X .

From Mukhopadhyay (1998, p. 123) and considering a random sample selected from each population stratum we can deduce:

$$E(e_{1st} e_{2st}) = \frac{1}{\bar{X}} \sum_{h=1}^L W_h^2 \gamma_h \frac{\mu_{03h}}{\mu_{02h}} \quad \text{and} \quad E(e_{1st} e_{3st}) = \frac{1}{\bar{X}} \sum_{h=1}^L W_h^2 \gamma_h \frac{\mu_{12h}}{\mu_{11h}},$$

$$\text{where } \mu_{rsh} = \frac{1}{N_h - 1} \sum_{i=1}^{N_h} (z_{hi} - \bar{Z}_h)^r (x_{hi} - \bar{X}_h)^s.$$

Recognizing that $\bar{Z} = \bar{Y}$ in Equation (6.18), the Bias and MSE of $\hat{\mu}_{Regst}$, are given by

$$\text{Bias}(\hat{\mu}_{Regst}) \cong - \sum_{h=1}^L W_h^2 \gamma_h \beta_c \left\{ \frac{\mu_{12h}}{\mu_{11h}} - \frac{\mu_{03h}}{\mu_{02h}} \right\} \quad (6.19)$$

and

$$\text{MSE}(\hat{\mu}_{Regst}) \cong \bar{Y}^2 \sum_{h=1}^L W_h^2 \gamma_h C_{zh}^2 (1 - \rho_c^2), \quad (6.20)$$

$$\text{where } \rho_c = \frac{\sum_{h=1}^L W_h^2 \gamma_h S_{zxh}}{\sqrt{\sum_{h=1}^L W_h^2 \gamma_h C_{zh}^2} \sqrt{\sum_{h=1}^L W_h^2 \gamma_h C_{xh}^2}}.$$

It can be verified easily that

(i) $\text{MSE}(\hat{\mu}_{Regst}) < \text{MSE}(\hat{\mu}_{Yst})$ if

$$\sum_{h=1}^L W_h^2 \gamma_h C_{zh}^2 \rho_c^2 > 0. \quad (6.21)$$

(ii) $MSE(\hat{\mu}_{Regst}) < MSE(\hat{\mu}_{Rst})$ if

$$\left(\sqrt{\sum_{h=1}^L W_h^2 \gamma_h C_{xh}^2} - \frac{\sum_{h=1}^L W_h^2 \gamma_h C_{zxh}}{\sqrt{\sum_{h=1}^L W_h^2 \gamma_h C_{xh}^2}} \right)^2 > 0. \quad (6.22)$$

These two conditions will always hold true indicating that, up to first order of approximation, the regression estimator performs better than ordinary mean and ratio estimators in stratified random sampling also, as they did in the case of simple random sampling.

On the other hand, we can say that

(iii) $MSE(\hat{\mu}_{Regst}) < MSE(\hat{\mu}_{Reg})$ if

$$\sum_{h=1}^L W_h^2 \gamma_h C_{zh}^2 (1 - \rho_c^2) < \frac{1-f}{n} C_z^2 (1 - \rho_{zx}^2), \quad (6.23)$$

a condition that can be ensured by a suitable stratification.

6.6 A Simulation Study

In this section, we present a simulation study with particular focus on comparing the performance of the proposed combined estimators $\hat{\mu}_{Rst}$ and $\hat{\mu}_{Regst}$ to the RRT mean estimator $\hat{\mu}_{Yst}$ and to the corresponding estimators in simple random sampling (Sousa et al., 2010; Gupta et al., 2012). For this purpose we rely on *Bias* and *MSE*, correct up to first order of approximation.

We considered three bivariate normal populations with different covariance matrices to represent the distribution of (Y, X) . The scrambling variable S is taken to be a normal distribution with mean equal to zero and standard deviation equal to 10% of the standard deviation of X . The reported scrambled response on Y is given by $Z = Y + S$.

All of the simulated populations have theoretical mean of $[Y, X]$ as $\mu = [5, 5]$ and covariance matrices as given below.

Population 1

$$N = 1000$$

$$\Sigma = \begin{bmatrix} 9 & 3.2 \\ 3.2 & 4 \end{bmatrix}, \rho_{XY} = 0.5333.$$

Population 2

$$N = 1000$$

$$\Sigma = \begin{bmatrix} 6 & 3.3 \\ 3.3 & 3 \end{bmatrix}, \rho_{XY} = 0.7778.$$

Population 3

$$N = 1000$$

$$\Sigma = \begin{bmatrix} 5 & 3 \\ 3 & 2 \end{bmatrix}, \rho_{XY} = 0.9487.$$

For each population we considered five sample sizes: $n = 30, 60, 150$ and 300 .

The population is divided in two strata according to a certain criteria set for the auxiliary variable. The sample size from each stratum is based on the *Neyman* allocation. We compare the results of stratified random sampling with the corresponding results of simple random sampling.

Table 6.1 below gives the empirical and theoretical *MSE*'s for the proposed combined estimators based on first order approximation. We use the following expression to find the Percent Relative Efficiency (*PRE*) of study estimators as compared to the ordinary sample mean:

$$PRE = \frac{MSE(\hat{\mu}_{Yst})}{MSE(\hat{\mu}_{\alpha})} \times 100,$$

where $\alpha = R_{st}, Reg_{st}$. This measure is calculated using first degree of approximation for *MSE* per unit estimator. We estimate the empirical *MSE* using 5000 samples of size n and considering the average of all the observed values.

Table 6.1: Empirical and Theoretical *MSE*, for the RRT mean estimator, ratio estimator (underlined) and regression estimator (bold); and corresponding *PRE* relative to the RRT mean estimator.

Population				MSE Estimation			
<i>N</i>	<i>N_h</i>	ρ_{XY}	ρ_{XYh}	<i>n</i>	Empirical <i>MSE</i>	Theoretical <i>MSE</i>	<i>PRE</i>
1000	$N_1 = 550$ $N_2 = 450$	0.5395	$\rho_{XY1} = 0.5397$ $\rho_{XY2} = 0.5410$	30	0.3039	0.2897	100.00
					<u>0.2227</u>	<u>0.2104</u>	<u>137.70</u>
					0.2395	0.2077	139.51
				60	0.1403	0.1404	100.00
					<u>0.1024</u>	<u>0.1019</u>	<u>137.77</u>
					0.1339	0.0993	141.34
				150	0.0533	0.0508	100.00
					<u>0.0380</u>	<u>0.0369</u>	<u>137.72</u>
					0.0488	0.0359	141.31
				300	0.0208	0.0209	100.00
					<u>0.0156</u>	<u>0.0152</u>	<u>137.73</u>
					0.0187	0.0148	141.00
	$N_1 = 550$ $N_2 = 450$	0.7827	$\rho_{XY1} = 0.7888$ $\rho_{XY2} = 0.7868$	30	0.2028	0.1932	100.00
					<u>0.0803</u>	<u>0.0763</u>	<u>253.17</u>
					0.0791	0.0792	244.04
				60	0.0937	0.0936	100.00
					<u>0.0371</u>	<u>0.0370</u>	<u>253.29</u>
					0.0430	0.0367	255.33
				150	0.0355	0.0339	100.00
					<u>0.0137</u>	<u>0.0134</u>	<u>253.21</u>
					0.0150	0.0131	258.38
				300	0.0139	0.0139	100.00
					<u>0.0057</u>	<u>0.0055</u>	<u>253.21</u>
					0.0061	0.0054	256.64
$N_1 = 550$ $N_2 = 450$	0.9501	$\rho_{XY1} = 0.9522$ $\rho_{XY2} = 0.9478$	30	0.1688	0.1608	100.00	
				<u>0.0346</u>	<u>0.0323</u>	<u>497.12</u>	
				0.0168	0.0179	899.95	
			60	0.0783	0.0779	100.00	
				<u>0.0157</u>	<u>0.0157</u>	<u>496.80</u>	
				0.0084	0.0079	979.71	
			150	0.0296	0.0282	100.00	
				<u>0.0059</u>	<u>0.0057</u>	<u>496.88</u>	
				0.0030	0.0028	1002.03	
			300	0.0115	0.0116	100.00	
				<u>0.0024</u>	<u>0.0023</u>	<u>496.98</u>	
				0.0012	0.0012	993.07	

According to the results in Table 6.1, all the percent relative efficiencies are greater than 100, indicating that the proposed combined estimators perform better than the ordinary mean estimator. The use of auxiliary information provides a gain for a stratified random sample. Therefore, the proposed estimators increases the accuracy since there is a significant correlation between *X* and *Y*.

These results for the stratified estimators agree with the Sousa et al. (2010) and Gupta et al. (2012) findings for a simple random sampling. Clearly the gain with regression estimator is substantial when correlation between the primary and auxiliary variables is high.

6.7 Numerical Example

We now compare the performances of the proposed combined estimators using a real data set. The data come from a sample from the survey on Information and Communication Technologies (ICT) usage in enterprises in 2010 with seat in Portugal (Smilhily and Storm, 2010). This survey intends to promote the development of the national statistical system in the information society and to contribute to a deeper knowledge about the usage of ICT by enterprises. The target population covers all industries with one and more persons employed in the sections of economic activity C (Manufacturing) to N (Administrative and support service activities) and S (Other service activities), from NACE¹ Rev. 2 (Eurostat, 2008).

The ICT survey has an extensive plan of indicators, so the use of auxiliary information on the sampling stage is essential to get a stratified random sample as a proper representation for the population. The enterprises commercialization is directly related to their turnover, so this auxiliary variable is usually used for stratification. In our example we consider three strata: the first one is enterprises with less than 10 million (in euros) of turnover, the second between 10 and less than 30 million of turnover, and the third with 30 million or more of turnover.

In our application the variable of interest Y is the purchase orders in 2010, collected by the ICT survey in that year. This is typically a confidential variable for enterprises, only known from business surveys. On the other hand, the auxiliary variable X is the turnover of each enterprise which is known for all the population and annually available with the statistical institutes as administrative information.

The purchase orders information was collected in the ICT survey and is known for a sample of 1698 small and medium enterprises (at least 10 and no more than 100 employees) in 2010. For this study, these 1698 enterprises are considered as our population. The scrambling variable S is taken to be a normal random variable with mean equal to zero. Given the high magnitude of the auxiliary variable X , we consider a standard deviation of S equal to 1% of the standard deviation of X , that is $\sigma_S = 0.01\sigma_X$. The reported scrambled response is given by $Z = Y + S$ (the purchase order value plus a random quantity). The variables Y and X are strongly correlated so we can take advantage of this correlation by using the combined ratio and regression estimators.

The variables X and Y are expressed in millions of Euros. We test our stratified sample estimators with random sample of sizes $n = 100, 250$ and 500 . The sample size of each stratum is allocated proportionally to the dimension of strata population.

¹NACE is derived from the French title "Nomenclature générale des Activités économiques dans les Communautés Européennes" (Statistical classification of economic activities in the European Communities).

Population Characteristics:

Stratum	N_h	ρ_{XYh}	μ_{Yh}	σ_{Yh}	μ_{Xh}	σ_{Xh}	Population
1	979	0.7802	2.15	2.46	3.12	2.68	$N = 1698, \rho_{XY} = 0.9368,$ $\beta_{YX} = 0.8284, \mu_Y = 14.44,$ $\sigma_Y = 22.39, \mu_X = 17.97, \sigma_X = 25.31$
2	362	0.7952	16.67	6.86	20.31	6.02	
3	357	0.8408	45.88	30.21	56.33	30.18	

Table 6.2 below presents the results for the empirical MSE estimates, the theoretical estimates, correct up to first degree of approximation, and the PRE of combined ratio and regression estimators relative to the ordinary sample mean in the stratified sample. For both sampling designs, we estimate the empirical MSE using 5000 samples of size n selected from the population.

We also show the Design Effect ($Deff$) comparing the efficiency of study estimators in stratified sample (Str) relative to the ordinary sample mean in simple random sample (SRS):

$$Deff = \frac{MSE(\hat{\mu}_Y)}{MSE(\hat{\mu}_\alpha)} \times 100,$$

where $\alpha = Y_{st}, R_{st}, Reg_{st}$.

Table 6.2: Empirical, theoretical MSE , PRE for the ratio estimator (underlined) and for the regression estimator (bold) relative to the RRT mean estimator and PRE for the simple random sample (SRS) relative to the stratified sample (Str).

Population			SRS		Str			$Deff^2$	
N	N_h	ρ_{XY}	n	Empirical MSE	Theoretical MSE	Empirical MSE	Theoretical MSE		PRE^1
1698	$N_1 = 979$ $N_2 = 362$ $N_3 = 357$	0.9368	100	5.6291	5.6291	1.9699	1.9373	100.00	290.56
				<u>0.8027</u>	<u>0.8027</u>	<u>0.5535</u>	<u>0.6577</u>	294.56	855.88
				0.6894	0.6894	0.6290	0.5303	365.35	1061.57
			250	2.0110	2.0403	0.6994	0.6948	100.00	293.66
				<u>0.2882</u>	<u>0.2909</u>	<u>0.2110</u>	<u>0.2397</u>	<u>289.91</u>	<u>851.35</u>
				0.2493	0.2499	0.2265	0.1863	373.00	1095.35
500	0.8502	0.8440	0.2917	0.2903	100.00	290.71			
	<u>0.1186</u>	<u>0.1204</u>	<u>0.0855</u>	<u>0.0992</u>	<u>292.66</u>	<u>850.81</u>			
	0.1022	0.1034	0.0882	0.0785	369.91	1075.37			

¹ MSE comparison of study estimators relative to the ordinary sample mean in Str .

² MSE comparison of study estimators in Str relative to the ordinary sample mean in SRS.

According to the results in Table 6.2, all of the percent relative efficiencies are greater than 100, so the proposed estimators perform better than the ordinary RRT mean estimator which does not use auxiliary information. Moreover, there is clear reduction of MSE if we compare the results based on stratification to those based on simple random sampling. The $Deff$ shows an increase in efficiency by using the stratified sample.

Taking into account the large correlation between the variable of interest Y and the

auxiliary variable X , the proposed estimators have similar gain regardless of the sample size. However, the gain is more evident in a simple random sample because the stratification already significantly reduces the MSE values.

6.8 Conclusions

In the survey research context, using auxiliary information can be essential to improve the accuracy of estimates, mainly when we have to deal with sensitive variables. We can observe from this study that the estimation of the mean of a sensitive variable can be improved by using a non-sensitive auxiliary variable.

The ratio and the regression estimators perform better than the RRT mean estimator in both simple random sampling and stratified sampling also. Although both the ratio and regression estimators perform better than the ordinary RRT mean estimator, the improvement is much larger with the regression estimator.

Regarding the efficiency, the results indicate that the proposed estimators become more and more efficient as the coefficient of correlation increases. When the study and the auxiliary variables are strongly correlated the proposed estimators, particularly the combined regression estimator performs much better, regardless the sample size. These results agree with the findings of Sousa et al. (2010) and Gupta et al. (2012) in simple random sampling.

All of the study estimators show better performance than the ordinary RRT sample mean. Nevertheless, the gain in accuracy is stronger in the simple random sampling because the stratification already reduces the MSE value for the RRT mean estimator.

The main conclusion of this study is that the advantage of using the RRT in the presence of auxiliary information still holds in the context of stratified sampling.

References

- COCHRAN, W.G. 1997. *Sampling Techniques*, 3rd Ed., New York, Wiley Eastern Ltd.
- EUROSTAT. 2008. NACE Rev. 2 - Statistical classification of economic activities in the European Community. *Official Publications of the European Communities*, 112-285 and 306-311.
- GUPTA, S. N., GUPTA, B. C. & SINGH, S. 2002. Estimation of sensitivity level of personal interview survey questions. *Journal of Statistical Planning and Inference*, 100, 239-247.
- GUPTA, S. & SHABBIR, J. 2008. On improvement in estimating the population mean in simple random sampling. *Journal of Applied Statistics*, 35(5), 559-566.

- GUPTA, S., SHABBIR, J. & SEHRA, S. 2010. Mean and sensitivity estimation in optional randomized response models. *Journal of Statistical Planning and Inference*, 140(10), 2870-2874.
- GUPTA, S., SHABBIR, J., SOUSA, R. & REAL, P. C. 2012. Estimation of the Mean of a Sensitive Variable in the Presence of Auxiliary Information. *Communications in Statistics - Theory and Methods*, 41(13-14), 2394-2404.
- KADILAR, C., CANDAN, M. & CINGI, H. 2007. Ratio estimators using robust regression. *Haceteppe Journal of Mathematics and Statistics*, 36(2), 81-188.
- KADILAR, C. & CINGI, H. 2003. Ratio estimator in stratified sampling. *Biometrical Journal*, 45, 218-225.
- KADILAR, C. & CINGI, H. 2005. A new estimator in stratified random sampling. *Communication in Statistics-Theory and Methods*, 34, 597-602.
- KHOSHNEVISAN, M., SINGH, R., CHAUHAN, P., SAWAN, N. & SMARANDACHE, F. 2007. A general family of estimators for estimating population mean using known value of some population parameter(s). *Far East Journal of Theoretical Statistics*, 22, 181-191.
- KOYUNCU, N. & KADILAR, C. 2008. Ratio and product estimators in stratified random sampling. *Journal of Statistical Planning and Inference*, 139(8), 2552-2558.
- KOYUNCU, N. & KADILAR, C. 2009. Family of estimators of population mean using two auxiliary variables in stratified random sampling. *Communications in Statistics: Theory and Methods*, 38(14), 2398-2417.
- KOYUNCU, N. & KADILAR, C. 2010. On the family of estimators of population mean in stratified random sampling. *Pakistan Journal of Statistics*, 26(2), 427-443.
- MUKHOPADHYAY, P. 1998. *Theory and Methods of Survey Sampling*, New Delhi, Prentice-Hall of India.
- NANGSUE, N. 2009. Adjusted Ratio and Regression Type Estimators for Estimation of Population Mean when some Observations are missing. *World Academy of Science, Engineering and Technology*, 53, 781-784.
- SAHOO, L. N., MAHAPATRA, N. & SENAPATI, S. C. 2009. A Generalized Method of Estimation for Two-stage Sampling Using Two Auxiliary Variables. *Journal of Statistical Theory and Practice*, 3(4), 831-839.
- SHABBIR, J. & GUPTA, S. 2005. Improved ratio estimators in stratified sampling. *American Journal of Mathematical and Management Sciences*, 25, 293-311.
- SHABBIR, J. & GUPTA, S. 2006. A new estimator of population mean in stratified sampling. *Communication in Statistics Theory and Methods*, 35, 1201-1209.

SHABBIR, J. & GUPTA, S. 2007. On improvement in variance estimation using auxiliary information. *Communication in Statistics-Theory and Methods*, 36(12), 2177-2185.

SHABBIR, J. & GUPTA, S. 2010. Estimation of the finite population mean in two-phase sampling when auxiliary variables are attribute. *Hacettepe Journal of Mathematics and Statistics*, 39(1), 121-129.

SINGH, H. P. & KUMAR, S. 2011. A General Family of Estimators of Finite Population Ratio, Product and Mean Using Two Phase Sampling Scheme in the Presence of Non-Response. *Journal of Statistical Theory and Practice*, 2(4), 677-692.

SINGH, H. P. & SOLANKI, R. S. 2012. An Efficient Class of Estimators for the Population Mean Using Auxiliary Information in Systematic Sampling. *Journal of Statistical Theory and Practice*, 6(2), 274-285.

SINGH, H.P. & VISHWAKARMA, G. K. 2008. A family of estimators of population mean using auxiliary information in stratified sampling. *Communication in Statistics-Theory and Methods*, 37(7), 1038-1050.

SMILHILY, M. & STORM, H. 2010. ICT usage in enterprises - 2009. *Eurostat Publications*, Issue 1.

SOUSA, R., SHABBIR, J., REAL, P. C. & GUPTA, S. 2010. Ratio estimation of the mean of a sensitive variable in the presence of auxiliary information. *Journal of Statistical Theory and Practice*, 4(3), 495-507.

SUKHATME, P.V. & SUKHATME, B.V. 1984. *Sampling theory of surveys with applications*, 3rd Ed., Ames, Iowa, Iowa State University Press.

WARNER, S. L. 1965. Randomized response: a survey technique for elimination evasive answer bias. *Journal of American Statistical Association*, 60, 63-69.

Appendix E - R Routines

Listing 6.1: R Code for Simulation Study of Proposed Estimator in Chapter 6

```

1
2 proj3_ratio_st_NeymanAlloc <- function(N, sigma, mu, L)
3 {
4
5   #Generation of a bivariate normal population
6   data_yx <- mvrnorm(10000, mu, sigma)
7   data_yx <- data.frame(data_yx)
8   colnames(data_yx) <- c("Y", "X")
9
10  indices1 <- round(runif(550, 0, 10000))
11  data_yx1 <- data_yx[indices1,]
12  data_yx1$ST <- 1
13  indices2 <- round(runif(450, 0, 10000))
14  data_yx2 <- data_yx[indices2,]
15  data_yx2$ST <- 2
16
17  data_yx <- rbind(data_yx1, data_yx2)
18
19  #Study variable
20  Y <- data_yx[,1]
21  #Auxiliary variable, correlated with Y
22  X <- data_yx[,2]
23  #Stratum
24  ST <- data_yx[,3]
25
26  #Scrambling variable independent of Y and X, with mean=0
27  S <- rnorm(N, mean=0, sd=0.1*sd(X))
28  #Scrambled response
29  Z <- Y+S
30
31  #Population
32  univ <- data.frame(cbind(Y=Y, S=S, Z=Z, X=X,
33    ST=ST, N RAND=runif(N)))
34  univ <- univ[order(univ$ST, univ$N RAND),]
35
36  #Coefficients of correlation
37  Ro_YX <- cor(Y, X)
38  Ro_ZX <- Ro_YX/sqrt(1+(var(S)/var(Y)))
39  Ro_YXh <- by(cbind(univ$Y, univ$X),
40    univ$ST, function(x) {cor(x[,1], x[,2])})
41  Ro_ZXh <- Ro_YXh/sqrt(1+(by(univ$S,
42    univ$ST, var)/by(univ$Y, univ$ST, var)))
43
44  #Population means
45  Mz <- mean(univ$Z)
46  Mx <- mean(univ$X)

```

```

47 My <- mean(univ$Y)
48
49 #Information
50 SY <- sd(univ$Y)
51 SYh <- by(univ$Y, univ$ST, sd)
52 SZ <- sd(univ$Z)
53 SZh <- by(univ$Z, univ$ST, sd)
54 SX <- sd(univ$X)
55 SXh <- by(univ$X, univ$ST, sd)
56
57 #Samples dimension
58 dim_samp <- c(30, 60, 150, 300)
59
60 #Information for the population
61 Nh <- by(univ$Z, univ$ST, length)
62 wh <- Nh/N
63
64 res <- NULL
65 for (i in 1:length(dim_samp))
66 {
67
68   #sample dimension
69   n <- dim_samp[i]
70
71   #sample with Neyman Allocation
72   n_total <- 0
73   samp <- NULL
74   for (l in 1:L)
75   {
76     n_aux <- round(n * ((length(univ$Z[univ$ST==l])
77       *sd(univ$Y[univ$ST==l]))/sum(Nh*by(univ$Y, univ$ST, sd))))
78     if (l==L) {n_aux<-n-n_total}
79     n_total <- n_total+n_aux
80     samp <- rbind(samp, univ[univ$ST==l, ][1:n_aux, ])
81   }
82
83   #Sampling rate for each stratum
84   nh <- by(samp$Z, samp$ST, length)
85   fh <- nh/Nh
86   gh <- (1-fh)/nh
87   f <- n/N
88
89   #Sampling mean for each stratum
90   mzh <- by(samp$Z, samp$ST, mean)
91   myh <- by(samp$Y, samp$ST, mean)
92   mxh <- by(samp$X, samp$ST, mean)
93   msh <- by(samp$S, samp$ST, mean)
94
95   #Population mean for each stratum
96   Mzh <- by(univ$Z, univ$ST, mean)

```

```

97   Myh <- by(univ$Y, univ$ST, mean)
98   Mxh <- by(univ$X, univ$ST, mean)
99   Msh <- by(univ$S, univ$ST, mean)
100
101   #Sampling mean for each stratum
102   szh <- by(samp$Z, samp$ST, sd)
103   sxh <- by(samp$X, samp$ST, sd)
104
105   mu11 <- cbind(sum((univ$Z[univ$ST==1]-Mzh[1])
106                 *(univ$X[univ$ST==1]-Mxh[1]))/(Nh[1]-1),
107                sum((univ$Z[univ$ST==2]-Mzh[2])
108                 *(univ$X[univ$ST==2]-Mxh[2]))/(Nh[2]-1))
109   mu12 <- cbind(sum((univ$Z[univ$ST==1]-Mzh[1])
110                 *((univ$X[univ$ST==1]-Mxh[1])^2))/(Nh[1]-1),
111                sum((univ$Z[univ$ST==2]-Mzh[2])
112                 *((univ$X[univ$ST==2]-Mxh[2])^2))/(Nh[2]-1))
113   mu02 <- cbind(sum((univ$X[univ$ST==1]-Mxh[1])^2)/(Nh[1]-1),
114                sum((univ$X[univ$ST==2]-Mxh[2])^2)/(Nh[2]-1))
115   mu03 <- cbind(sum((univ$X[univ$ST==1]-Mxh[1])^3)/(Nh[1]-1),
116                sum((univ$X[univ$ST==2]-Mxh[2])^3)/(Nh[2]-1))
117
118   #Ratio
119   R <- mean(univ$X)/mean(samp$X)
120
121   #Ordinary meam
122   est1 <- sum(wh*mzh)
123   #Ratio estimator
124   est2 <- sum(wh*mzh)*(Mx/sum(wh*mxh))
125   #Regression estimator
126   betac <- sum((wh^2)*gh*Ro_ZXh*szh*sxh)/sum((wh^2)*gh*Ro_ZXh*(sxh^2))
127   est3 <- sum(wh*mzh)+betac*(Mx-sum(wh*mxh))
128
129   #Coefficient of variation
130   c_xh <- by(univ$X, univ$ST, sd)/Mxh
131   c_yh <- by(univ$Y, univ$ST, sd)/Myh
132   c2_xh <- c_xh^2
133   c2_yh <- c_yh^2
134   c_zh <- by(univ$Z, univ$ST, sd)/Mzh
135   c2_zh <- c_zh^2
136
137   #Bias of ratio estimator - 1st degree approximation
138   bias2i <- My*sum((wh^2)*gh*(c2_xh-Ro_ZXh*c_zh*c_xh))
139   #Bias of regression estimator - 1st degree approximation
140   bias3i <- -sum(c((wh^2)*gh*betac)*((mu12/mu11)-(mu03/mu02)))
141
142   #Mean Square Error of 1st estimator (ordinal mean)
143   mse1 <- sum((wh^2)*(gh*(by(univ$Y, univ$ST, var)
144                             +by(univ$S, univ$ST, var))))
145   #Mean Square Error of ratio estimator
146   #1st degree approximation

```

```

147 mse2i <- (My^2)*sum((wh^2)*gh*(c2_zh+c2_xh-2*Ro_ZXh*c_zh*c_xh))
148 #Mean Square Error of regression estimator
149 #1st degree approximation
150 rhoc <- sum((wh^2)*gh*Ro_ZXh*szh*sxh)
151 / (sqrt(sum((wh^2)*gh*(szh^2)))*sqrt(sum((wh^2)*gh*(sxh^2))))
152 mse3i <- (My^2)*sum((wh^2)*gh*c2_zh*(1-(rhoc^2)))
153
154 #Empirical results
155 #Simulation of 5000 replicas of estimates
156 ...
157
158 #Results
159 res <- rbind(res, cbind(Nh, N, nh, n,
160                         Ro_YXh, Ro_ZXh,
161                         SY, SYh, SX, SXh,
162                         Mxh, Mx, Myh, My, Mzh, Mz,
163                         med_est1, med_est2, med_est3,
164                         bias2i, bias3i,
165                         emp_mse1, mse1,
166                         emp_mse2, mse2i,
167                         emp_mse3, mse3i))
168 }
169 colnames(res) <- c("Nh", "N", "nh", "n",
170                  "RhoXYh", "RhoZXh",
171                  "SY", "SYh", "SX", "SXh",
172                  "MXh", "MX", "MYh", "MY", "MZh", "MZ",
173                  "Est1", "Est2", "Est3",
174                  "BIAS2I", "BIAS3I",
175                  "EMP_MSE1", "MSE1",
176                  "EMP_MSE2", "MSE2I",
177                  "EMP_MSE3", "MSE3I")
178 return(res)
179 }
180
181 #Package for generation
182 require(MASS)
183
184 #Parameters
185 #Population dimension
186 N <- 1000
187 #Variance-Covariance matrix
188 sigma1 <- matrix(c(9, 3.2, 3.2, 4), 2, 2)
189 sigma2 <- matrix(c(6, 3.3, 3.3, 3), 2, 2)
190 sigma3 <- matrix(c(5, 3, 3, 2), 2, 2)
191 #Mean vector
192 mu <- c(5, 5)
193 #Number of strata
194 L <- 2
195
196 res <- NULL

```

```
197 for (i in 1:length(N))
198 {
199   res <- rbind(res,proj3_ratio_st_NeymanAlloc(N[i],sigma1,mu,L))
200   res <- rbind(res,proj3_ratio_st_NeymanAlloc(N[i],sigma2,mu,L))
201   res <- rbind(res,proj3_ratio_st_NeymanAlloc(N[i],sigma3,mu,L))
202 }
203
204 write.table(res,"chapter6_ss_results.txt",sep="\t",dec="," ,row.names=FALSE)
```

Listing 6.2: R Code for Numerical Example of Proposed Estimator in Chapter 6

```

1 proj3_ratio_st_real <- function(Y,X,N)
2 {
3
4
5   L <- 3
6   data_yx <- data.frame(cbind(Y,X))
7   colnames(data_yx) <- c("Y","X")
8
9   #Strata
10  data_yx$ST <- 0
11  data_yx$ST <- ifelse(data_yx$X<10,1,
12                      ifelse(data_yx$X>=10 & data_yx$X<30,2,
13                              ifelse(data_yx$X>=30,3,0)))
14
15  data_yx <- data_yx[order(data_yx$ST),]
16  Y <- data_yx$Y
17  X <- data_yx$X
18  ST <- data_yx$ST
19
20  S<-NULL
21  #Scrambling variable independent of Y and X, with mean=0
22  for (s in 1:L)
23  {
24    S <- c(S,rnorm(sum(ST==s),mean=0,sd=0.01*sd(X[ST==s])))
25  }
26  #Scrambled response
27  Z <- Y+S
28
29  #Population
30  univ <- data.frame(cbind(Y=Y,S=S,Z=Z,X=X,ST=ST,NRAND=runif(N)))
31  univ <- univ[order(univ$ST,univ$NRAND),]
32
33  #Coefficients of correlation
34  Ro_YX <- cor(Y,X)
35  Ro_ZX <- Ro_YX/sqrt(1+(var(S)/var(Y)))
36  Ro_YXh <- by(cbind(univ$Y,univ$X),univ$ST,
37              function(x) {cor(x[,1],x[,2])})
38  Ro_ZXh <- Ro_YXh/sqrt(1+(by(univ$S,univ$ST,var)
39                             /by(univ$Y,univ$ST,var)))
40
41  #Population means
42  Mz <- mean(univ$Z)
43  Mx <- mean(univ$X)
44  My <- mean(univ$Y)
45
46  #Information
47  SY <- sd(univ$Y)

```

```

48 SYh <- by(univ$Y, univ$ST, sd)
49 SZ <- sd(univ$Z)
50 SZh <- by(univ$Z, univ$ST, sd)
51 SX <- sd(univ$X)
52 SXh <- by(univ$X, univ$ST, sd)
53
54 #Samples dimension
55 dim_samp <- c(100,250,500)
56
57 #Information for the population
58 Nh <- by(univ$Z, univ$ST, length)
59 wh <- Nh/N
60 nL <- Nh[L]
61
62 res <- NULL
63 for (i in 1:length(dim_samp))
64 {
65   #sample dimension
66   n <- dim_samp[i]
67
68   #sample with Neyman Allocation
69   n_total <- 0
70   samp <- NULL
71   for (l in 1:L)
72   {
73     n_aux <- round(n*(Nh[l]/N))
74     n_total <- n_total+n_aux
75     samp <- rbind(samp, univ[univ$ST==l, ][1:n_aux, ])
76   }
77
78   #Sampling rate for each stratum
79   nh <- by(samp$Z, samp$ST, length)
80   fh <- nh/Nh
81   gh <- (1-fh)/nh
82   f <- n/N
83
84   #Sampling mean for each stratum
85   mzh <- by(samp$Z, samp$ST, mean)
86   myh <- by(samp$Y, samp$ST, mean)
87   mxh <- by(samp$X, samp$ST, mean)
88   msh <- by(samp$S, samp$ST, mean)
89
90   #Population mean for each stratum
91   Mzh <- by(univ$Z, univ$ST, mean)
92   Myh <- by(univ$Y, univ$ST, mean)
93   Mxh <- by(univ$X, univ$ST, mean)
94   Msh <- by(univ$S, univ$ST, mean)
95
96   #Sampling mean for each stratum
97   szh <- by(samp$Z, samp$ST, sd)

```

```

98   sxh <- by(samp$X, samp$ST, sd)
99
100  mu11 <- cbind(sum((univ$Z[univ$ST==1]-Mzh[1])*(univ$X[univ$ST==1]
101    -Mxh[1]))/(Nh[1]-1), sum((univ$Z[univ$ST==2]-Mzh[2])
102    *(univ$X[univ$ST==2]-Mxh[2]))/(Nh[2]-1),
103    sum((univ$Z[univ$ST==3]-Mzh[3])*(univ$X[univ$ST==3]
104    -Mxh[3]))/(Nh[3]-1))
105  mu12 <- cbind(sum((univ$Z[univ$ST==1]-Mzh[1])*((univ$X[univ$ST==1]
106    -Mxh[1])^2))/(Nh[1]-1), sum((univ$Z[univ$ST==2]-Mzh[2])
107    *((univ$X[univ$ST==2]-Mxh[2])^2))/(Nh[2]-1),
108    sum((univ$Z[univ$ST==3]-Mzh[3])*((univ$X[univ$ST==3]
109    -Mxh[3])^2))/(Nh[3]-1))
110  mu02 <- cbind(sum((univ$X[univ$ST==1]-Mxh[1])^2)/(Nh[1]-1),
111    sum((univ$X[univ$ST==2]-Mxh[2])^2)/(Nh[2]-1),
112    sum((univ$X[univ$ST==3]-Mxh[3])^2)/(Nh[3]-1))
113  mu03 <- cbind(sum((univ$X[univ$ST==1]-Mxh[1])^3)/(Nh[1]-1),
114    sum((univ$X[univ$ST==2]-Mxh[2])^3)/(Nh[2]-1),
115    sum((univ$X[univ$ST==3]-Mxh[3])^3)/(Nh[3]-1))
116
117  #Ratio
118  R <- mean(univ$X)/mean(samp$X)
119
120  #Ordinary meam
121  est1 <- sum(wh*mzh)
122  #Ratio estimator
123  est2 <- sum(wh*mzh)*(Mx/sum(wh*mhx))
124  #Regression estimator
125  betac <- sum((wh^2)*gh*Ro_ZXh*szh*sxh)/sum((wh^2)*gh*Ro_ZXh*(sxh^2))
126  est3 <- sum(wh*mzh)+betac*(Mx-sum(wh*mhx))
127
128  #Coefficient of variation
129  c_xh <- by(univ$X, univ$ST, sd)/Mxh
130  c_yh <- by(univ$Y, univ$ST, sd)/Myh
131  c2_xh <- c_xh^2
132  c2_yh <- c_yh^2
133  c_zh <- by(univ$Z, univ$ST, sd)/Mzh
134  c2_zh <- c_zh^2
135
136  #Bias of ratio estimator - 1st degree approximation
137  bias2i <- My*sum((wh^2)*gh*(c2_xh-Ro_ZXh*c_zh*c_xh))
138  #Bias of regression estimator - 1st degree approximation
139  bias3i <- -sum(c((wh^2)*gh*betac)*((mu12/mu11)-(mu03/mu02)))
140
141  #Mean Square Error of 1st estimator (ordinal mean)
142  mse1 <- sum((wh^2)*(gh*(by(univ$Y, univ$ST, var)+by(univ$S, univ$ST, var))))
143  #Mean Square Error of ratio estimator - 1st degree approximation
144  mse2i <- (My^2)*sum((wh^2)*gh*(c2_zh+c2_xh-2*Ro_ZXh*c_zh*c_xh))
145  #Mean Square Error of regression estimator - 1st degree approximation
146  rhoc <- sum((wh^2)*gh*Ro_ZXh*szh*sxh)
147  / (sqrt(sum((wh^2)*gh*(szh^2)))*sqrt(sum((wh^2)*gh*(sxh^2))))

```



```

148     mse3i <- (My^2) *sum((wh^2) *gh*c2_zh*(1-(rhoc^2)))
149
150     #Empirical results
151     #Simulation of 5000 replicas of estimates
152     ...
153
154     #Results
155     res <- rbind(res,cbind(Nh,N,nh,n,
156                           Ro_YXh,Ro_ZXh,
157                           SY,SYh,SX,SXh,
158                           Mxh,Mx,Myh,My,Mzh,Mz,
159                           med_est1,med_est2,med_est3,
160                           bias2i,bias3i,
161                           emp_mse1,mse1,
162                           emp_mse2,mse2i,
163                           emp_mse3,mse3i))
164   }
165   colnames(res) <- c("Nh","N","nh","n",
166                     "RhoXYh","RhoZXh",
167                     "SY","SYh","SX","SXh",
168                     "MXh","MX","MYh","MY","MZh","MZ",
169                     "Est1","Est2","Est3",
170                     "BIAS2I","BIAS3I",
171                     "EMP_MSE1","MSE1",
172                     "EMP_MSE2","MSE2I",
173                     "EMP_MSE3","MSE3I")
174   return(res)
175 }
176
177 #Package for generation
178 require(MASS)
179
180 #Import data
181 data_yx <- read.table("IUTICE10_BA.txt",sep="\t",dec=".",header = T)
182 data_yx <- data_yx[data_yx$NPS>=10 & data_yx$NPS<150,]
183 data_yx <- data_yx[data_yx$sturn<=200,]
184 #Study variable (purchase, millions of euros)
185 Y <- data_yx$purch
186 #Auxiliary variable, correlated with Y (turnover, millions of euros)
187 X <- data_yx$sturn
188 #Data application
189 N <- dim(data_yx)[1]
190
191 res <- proj3_ratio_st_real(Y,X,N)
192 #Export data
193 write.table(res,"chapter6_ne_results.txt",sep="\t",dec=".",row.names=FALSE)

```




General Discussion

7.1 Summary

Our thesis work is based on the improvement of the mean estimation of sensitive variables (Edwards, 1957; Groves et al., 2004). In the sampling literature (Cochran, 1997; Mukhopadhyay, 1998; Särndal et al., 1997; Sukhatme and Sukhatme, 1984), researchers have proposed several estimators which use auxiliary information in order to improve their performance. Over the chapters of this thesis we have proposed different estimators which combine the Randomized Response Technique (RRT) method (Eichhorn, 1983; Warner, 1965) with the use of auxiliary information.

In section 7.2 we present a numerical example that aims to make a comparison of the performance of the main proposed estimators. For that purpose we conduct a study with a real dataset and we show the numerical results for the *Bias* and Mean Square Error (*MSE*), as well as graphic evidence which illustrates the performance of each estimator in terms of estimation precision.

7.2 Comparison of the main study estimators

In this section we conduct a study with a real dataset with particular focus on comparing the performance of the main estimators proposed in this thesis, using the *Bias* and the *MSE* results as the criteria.

Consider a real dataset concerning enterprises for the Monthly Economic Survey (MES) in Portugal. The survey is conducted to provide an accurate picture of business

trends of enterprises. It provides short-term indicators on a monthly basis compiled for four sectors: industry, retail trade, construction and service sector.

Generally, the enterprises do not want to report the value of their orders. This is typically a confidential variable for enterprises, only known from business surveys. Nevertheless, every year the entity responsible for MES, the Statistics Portugal [1], provides administrative information with the value of the orders for the previous year. Thus, in our numerical example, we consider the value of orders in 2009 as a sensitive variable and the value of orders in 2008 as an auxiliary variable.

Let Y be the annual orders amount in 2009 collected by the MES in that year. The auxiliary variable X is the annual orders amount in 2008, available from business data registers. The variables Y and X are strongly correlated so we can take advantage of this correlation by using an auxiliary variable. We take the 608 business survey respondents common between 2008 and 2009 as our population. For the RRT part, let S be a normal random variable with mean equal to zero and standard deviation equal to 10% of the standard deviation of X . The reported response is given by $Z = Y + S$ (the orders amount in 2008 plus a random quantity). The summary statistics about the populations are given below.

Population Characteristics:

$N = 608, \rho_{XY} = 0.9447$
$\mu_X = 21357.69, \mu_Y = 17828.2$ (in thousands of Euros)
$\sigma_X = 65874.83, \sigma_Y = 57489.53$ and $\sigma_{XY} = 3577597688$

We use the following samples sizes in our simulation study: $n = 50, 100, 200$ and 300 .

In this study we compare the results for the ordinary RRT sample mean ($\hat{\mu}_Y$) to the main estimators proposed in this study: the ratio estimator ($\hat{\mu}_R$) (Sousa et al., 2010), the regression estimator ($\hat{\mu}_{Reg}$) (Gupta et al., 2012), the generalized regression-cum-ratio estimator ($\hat{\mu}_{GRR}$) (Gupta et al., 2012), the generalized regression-cum-exponential estimator ($\hat{\mu}_{exp1}$) (Koyuncu et al., 2013) and the improved exponential estimator ($\hat{\mu}_{IE}$) (Gupta et al., 2013).

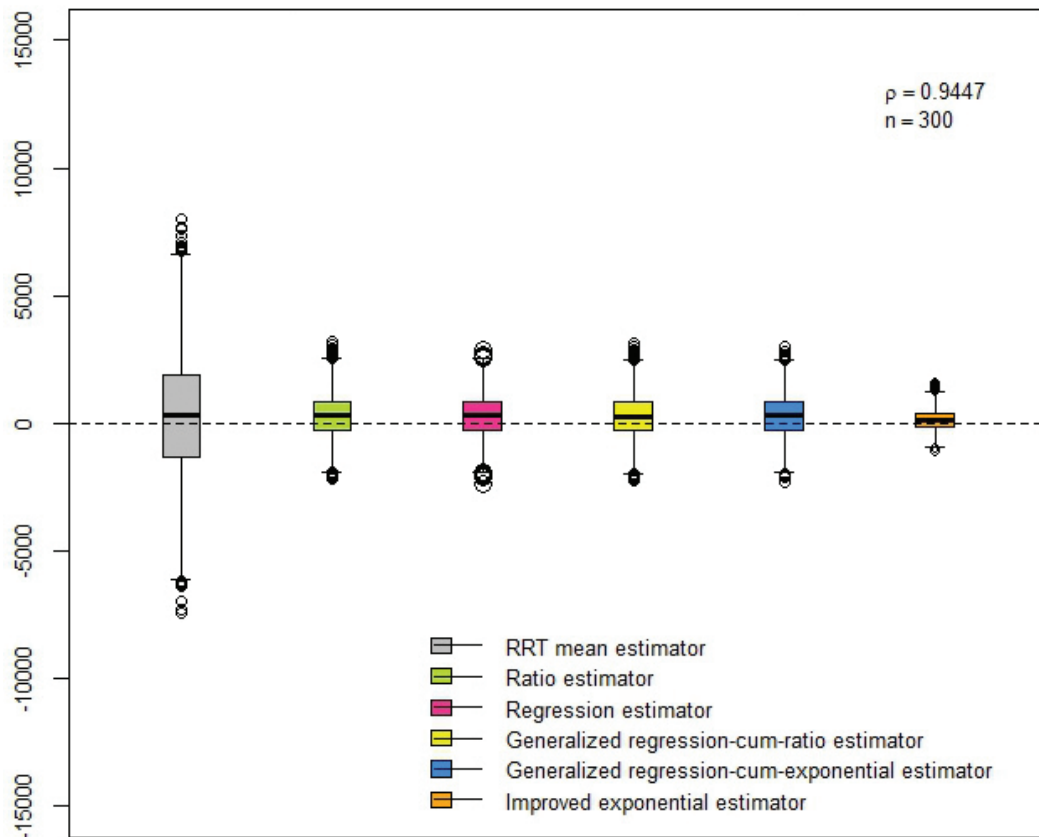
In Table 7.1 below we present the the theoretical *ARB* results for all the estimators in comparison.

Table 7.1: Theoretical *ARB* for the estimators in comparison.

Population		Estimator	Theoretical <i>ARB</i>			
N	ρ_{XY}		$n = 50$	$n = 100$	$n = 200$	$n = 300$
608	0.9447	$(\hat{\mu}_R)$	0.0022	0.0010	0.0004	0.0002
		$(\hat{\mu}_{Reg})$	0.0080	0.0036	0.0015	0.0007
		$(\hat{\mu}_{GRR})$	0.0225	0.0104	0.0042	0.0021
		$(\hat{\mu}_{exp1})$	0.0215	0.0093	0.0036	0.0018
		$(\hat{\mu}_{IE})$	0.0042	0.0022	0.0009	0.0005

The *ARB* results show that it is not always the case that estimators with better performance in terms of accuracy are the best performing in terms of *Bias* as well. However, the improved exponential estimator ($\hat{\mu}_{IE}$) manages to combine great precision results with reduction in *Bias* compared to other estimators which use auxiliary information, such as the ratio estimator ($\hat{\mu}_R$), the regression estimator ($\hat{\mu}_{Reg}$) and the study generalized exponential estimators ($\hat{\mu}_{exp1}$ and $\hat{\mu}_{IE}$).

From our 5000 samples, selected for each sample size and for each estimator, we take the empirical *Bias* and we draw a graph, presented in Figure 7.1, which shows the *Bias* distribution for all the generated estimates.

Figure 7.1: Distribution of empirical *Bias*

According to the graph in Figure 7.1 all the estimators have *Bias* about zero but it is the improved exponential estimator which shows less dispersion, with results closer to zero. Despite being unbiased, the RRT mean estimator's greater empirical *Bias* relative to the estimators which use auxiliary information, also based on RRT, is obvious in this graph indicating that just a RRT version of the mean estimator might not be enough.

Table 7.2 below gives empirical and theoretical *MSE*'s based on the first order of approximation for all the estimators considered here. We estimate the empirical *MSE* using 5000 samples of various sizes selected from the study population. We use the following expression to find the Percent Relative Efficiency (*PRE*) of ratio, regression, generalized regression-cum-ratio, generalized regression-cum-exponential and improved exponential estimators as compared to the RRT mean estimator:

$$PRE = \frac{MSE(\hat{\mu}_Y)}{MSE(\hat{\mu}_\alpha)} \times 100,$$

where $\alpha = R, Reg, GRR, exp1, IE$.

Table 7.2: Empirical *MSE*, theoretical *MSE* correct up to 1st order approximation and *PRE* for all the estimators in comparison relative to the RRT mean estimator.

Population			MSE Estimation			
<i>N</i>	ρ_{XY}	<i>n</i>	Estimator	Empirical	Theoretical	<i>PRE</i>
608	0.9447	50	$\hat{\mu}_Y$	63985642.27	61487318.16	100.00
			$\hat{\mu}_R$	7098375.83	7357491.69	835.71
			$\hat{\mu}_{Reg}$	7520669.52	7349015.26	836.67
			$\hat{\mu}_{GRR}$	6682678.20	7148757.06	860.11
			$\hat{\mu}_{exp1}$	6308486.16	6721357.39	914.81
			$\hat{\mu}_{IE}$	1633458.04	1334326.31	4608.12
		100	$\hat{\mu}_Y$	27999936.04	27988850.92	100.00
			$\hat{\mu}_R$	3593960.07	3349109.12	835.71
			$\hat{\mu}_{Reg}$	3531015.22	3345250.67	836.67
			$\hat{\mu}_{GRR}$	3431317.25	3307434.75	846.24
			$\hat{\mu}_{exp1}$	3252540.69	3213576.01	870.96
			$\hat{\mu}_{IE}$	794022.08	686993.08	4074.11
		200	$\hat{\mu}_Y$	11444042.86	11239617.30	100.00
			$\hat{\mu}_R$	1475109.02	1344917.84	835.71
			$\hat{\mu}_{Reg}$	1412070.68	1343368.38	836.67
			$\hat{\mu}_{GRR}$	1416483.71	1337528.93	840.33
			$\hat{\mu}_{exp1}$	1370565.25	1322001.04	850.20
			$\hat{\mu}_{IE}$	0333343.39	292135.71	3847.40
		300	$\hat{\mu}_Y$	5728956.33	5656539.42	100.00
			$\hat{\mu}_R$	791784.85	676854.07	835.71
			$\hat{\mu}_{Reg}$	752752.90	676074.28	836.67
			$\hat{\mu}_{GRR}$	765577.73	674615.91	838.48
			$\hat{\mu}_{exp1}$	740274.33	670651.05	843.44
			$\hat{\mu}_{IE}$	175396.47	149770.32	3776.81

According to the *MSE* results in Table 7.2 the component regression shows performance gains. As expected and shown in Chapter 5 the best performance comes from

improved exponential estimator because of its large reduction in MSE .

From our 5000 samples, selected for each sample size and for each estimator, we take the empirical MSE and we draw a graph, presented in Figure 7.2, which shows the precision distribution for all the generated estimates.

According to the graph in Figure 7.2 the use of auxiliary information significantly reduces the magnitude of MSE , particularly in the improved exponential estimator which presents results closer to zero.

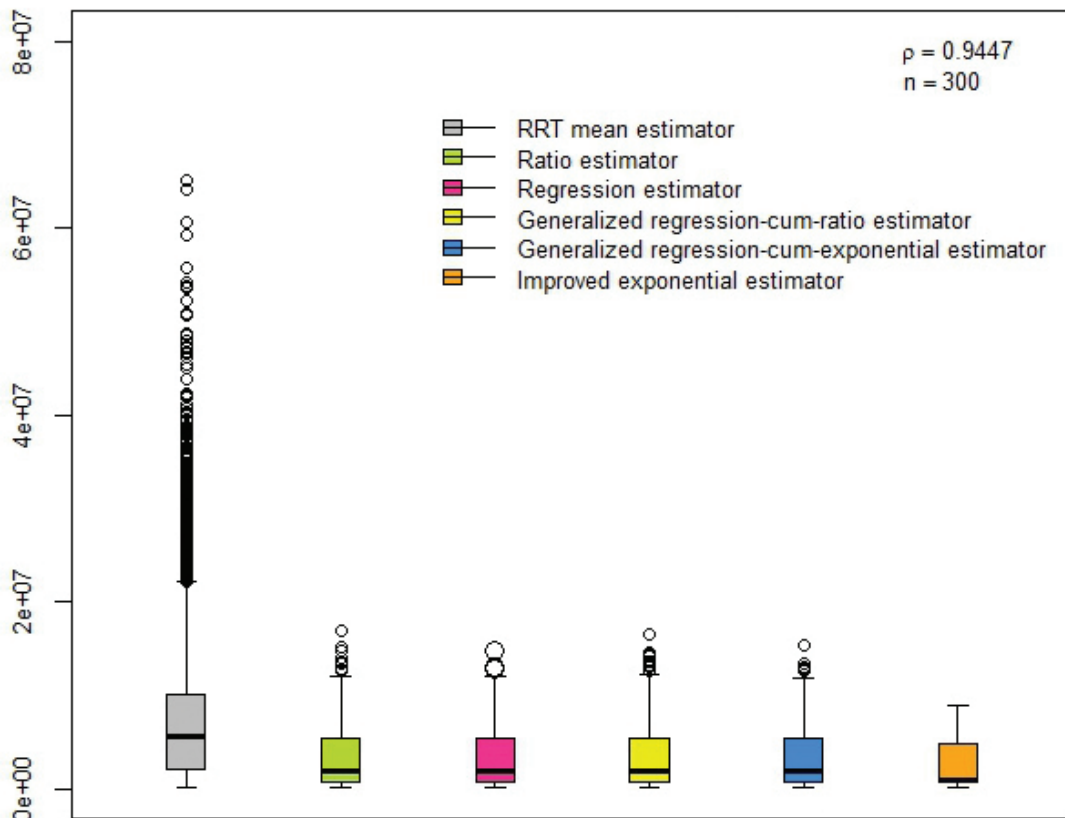


Figure 7.2: Distribution of empirical MSE

7.3 Final Remarks

The aim of this project was to develop new methodologies that can potentially improve the mean estimation in the presence of auxiliary information. These new methodologies were proposed in Chapters 2 to 6 and were compared with each other and with the ordinary RRT mean estimator which does not use the auxiliary information. For that purpose we studied theoretically the proposed estimators, deriving the expressions for

the *Bias* and Mean Square Error (*MSE*) correct up to first or second order approximations. Also, R routines [2] were developed for an extensive study by using real data and simulated data for all the estimators under study.

One of the main conclusions was the use of auxiliary information significantly reduces the magnitude of *MSE*, providing a gain for the parameter estimation based on RRT, just as in the context of direct estimation of non-sensitive parameters.

We concluded from this study that the estimation of the mean of a sensitive variable can be improved further by using a correlated non-sensitive auxiliary variable.

When there is a high correlation between the study variable and the auxiliary variable the regression estimator performs better than ratio estimator.

We also found some exponential type estimators more efficient than the ratio and regression type estimators.

We showed that the advantage of using the RRT in the presence of auxiliary information still holds with other sampling designs, such as the stratified sampling (Sousa et al., 2013).

Even though during the thesis project we have tested many new estimators and compared them to existing estimators in literature, we only present those who showed an effective improvement relative to the existing estimators or to the estimators previously proposed for us.

We are aware that in this area there is still much more to explore and part of our future work plans consists of studying other combinations of estimators, as well as the application to different sampling designs and with different techniques which provides confidence to the respondents when they have to answer to sensitive questions. Also, we would like to plan and implement a survey to a group of respondents who struggle with sensitive questions in order to evaluate the performance of proposed estimators with a real application of a RRT.

Even with the natural constraints we face in this kind of research, this study was a great challenge, both in the theoretical context as well as in the practical context of parameter estimation in the presence of auxiliary information.

References

- COCHRAN, W.G. 1997. *Sampling Techniques*, 3rd Ed., New York, Wiley Eastern Ltd.
- EDWARDS, A. L. 1957. *The social desirability variable in personality assessment and research*, New York: Dryden, Praeger.
- EICHHORN, B. H. & HAYRE, L. S. 1983. Scrambled randomized response methods for

obtaining sensitive quantitative data. *Journal of Statistical Planning and Inference*, 7, 307-316.

GROVES, R. M., FOWLER, F. J., COUPER, M. P., LEPKOWSKI, J. M., SINGER, E. & TOURANGEAU, R. 2004. *Survey Methodology*, New York, Wiley.

GUPTA, S., SHABBIR, J., SOUSA, R. & REAL, P. C. 2012. Estimation of the Mean of a Sensitive Variable in the Presence of Auxiliary Information. *Communications in Statistics - Theory and Methods*, 41(13-14), 2394-2404.

GUPTA, S., SHABBIR, J., SOUSA, R. & REAL, P. C. 2013. Improved exponential type estimators of the mean of a sensitive variable in the presence of non-sensitive auxiliary information. (*submitted*)

KOYUNCU, N., GUPTA, S. & SOUSA, R. 2013. Exponential type estimators of the mean of a sensitive variable in the presence of non-sensitive auxiliary information. *Communications in Statistics - Simulation and Computation*. (*accepted*).

MUKHOPADHYAY, P. 1998. *Theory and Methods of Survey Sampling*, New Delhi, Prentice-Hall of India.

SÄRDNAL, C.-E., SWENSSON, B. & WRETMAN, J. 1997. *Model assisted survey sampling*, Springer series in statistics.

SOUSA, R., GUPTA, S., SHABBIR, J. & REAL, P. C. 2013. Improved Mean Estimation of a Sensitive Variable Using Auxiliary Information in Stratified Sampling. *Journal of Statistics and Management Systems*. (*submitted*).

SOUSA, R., SHABBIR, J., REAL, P. C. & GUPTA, S. 2010. Ratio estimation of the mean of a sensitive variable in the presence of auxiliary information. *Journal of Statistical Theory and Practice*, 4(3), 495-507.

SUKHATME, P.V. & SUKHATME, B.V. 1984. *Sampling theory of surveys with applications*, 3rd Ed., Ames, Iowa, Iowa State University Press.

WARNER, S. L. 1965. Randomized response: a survey technique for elimination evasive answer bias. *Journal of American Statistical Association*, 60, 63-69.

[1] Statistics Portugal: http://www.ine.pt/xportal/xmain?xpid=INE&xpgid=ine_main&xlang=en.

[2] The R Project for Statistical Computing: www.r-project.org.

Appendix F - R Routines

Listing 7.1: R Code for Numerical Example of Proposed Estimator in Chapter 7

```

1
2 comparison_chap7 <- function(Y,X,N)
3 {
4   #Coefficient of correlation between Y and X
5   Ro_YX <- cor(Y,X)
6
7   #Scrambling variable independent of Y and X, with mean=0
8   S <- rnorm(N,mean=0,sd=0.1*sd(X))
9   #Scrambled response
10  Z <- Y+S
11
12  #Coefficient of correlation between Z and X
13  Ro_ZX <- Ro_YX/sqrt(1+(var(S)/var(Y)))
14
15  #population
16  univ <- data.frame(cbind(Y=Y,S=S,Z=Z,X=X,NRAND=runif(N)))
17  univ <- univ[order(univ$NRAND),]
18
19  #Mean of Y
20  mz <- mean(univ$Z)
21  mx <- mean(univ$X)
22  my <- mean(univ$Y)
23
24  mu11 <- sum((univ$Z-mz)*(univ$X-mx))/(N-1)
25  mu12 <- sum((univ$Z-mz)*((univ$X-mx)^2))/(N-1)
26  mu02 <- sum((univ$X-mx)^2)/(N-1)
27  mu03 <- sum((univ$X-mx)^3)/(N-1)
28
29  beta_zx <- Ro_YX*(sd(univ$Y)/sd(univ$X))
30
31  #Samples dimension
32  dim_samp <- c(50,100,200,300)
33
34  #Initialize the variables...
35
36  for (i in 1:length(dim_samp))
37  {
38    #sample dimension
39    n <- dim_samp[i]
40    #sample
41    samp <- univ[1:n,]
42    #Sampling rate
43    f <- n/N
44
45    #Ordinary meam
46    est1 <- mean(samp$Z)

```

```

47   #Ratio estimator
48   est2 <- mean(samp$Z) * (mx/mean(samp$X))
49   #Regression estimator
50   est3 <- mean(samp$Z) + beta_zx * (mx - mean(samp$X))
51
52   #Coefficient of variation
53   c_x <- sd(univ$X)/mx
54   c_y <- sd(univ$Y)/my
55   c2_x <- c_x^2
56   c2_y <- c_y^2
57   c2_z <- c2_y + (var(univ$S)/(my^2))
58   c_z <- sqrt(c2_z)
59
60   l <- (1-f)/n
61
62   #Generalized Regression-cum-ratio Estimator
63   k1 <- (1 - ((1-f) * c2_x/n)) / (1 - ((1-f)/n) * (c2_x - c2_z * (1 - (Ro_ZX^2))))
64   k2 <- (my/mx) * (1 + k1 * ((Ro_ZX * c_z/c_x) - 2))
65   est5 <- (k1 * mean(samp$Z) + k2 * (mx - mean(samp$X)))
66           * (mx/mean(samp$X))
67
68   #Generalized regression-cum-exponential type Estimator
69   w1 <- (1 - (1 * c2_x/8)) / (1 + 1 * c2_z * (1 - (Ro_ZX^2)))
70   w2 <- (my/mx) * (0.5 - w1 * (1 - (Ro_ZX * c_z/c_x)))
71   est7 <- (w1 * mean(samp$Z) + w2 * (mx - mean(samp$X)))
72           * exp((mx - mean(samp$X)) / (mx + mean(samp$X)))
73
74   #2nd Improved Exponential Estimator
75   A <- 1 + 1 * c2_z + 1 * c2_x - 2 * 1 * Ro_ZX * c_z * c_x
76   B <- 1 + 1 * c2_x
77   C <- 1 + (3/8) * 1 * c2_x - 0.5 * 1 * Ro_ZX * c_z * c_x
78   D <- 1 + (3/8) * 1 * c2_x
79   E <- 1 + 1 * c2_x - 1 * Ro_ZX * c_z * c_x
80   z1 <- (B * C - D * E) / (A * B - (E^2))
81   z2 <- my * (A * D - C * E) / (A * B - (E^2))
82   est8 <- (z1 * mean(samp$Z) + z2)
83           * exp((mx - mean(samp$X)) / (mx + mean(samp$X)))
84
85   #Mean Square Error of 1st estimator (ordinal mean)
86   mse1 <- ((1-f)/n) * (var(univ$Y) + var(univ$S))
87
88   #Bias of ratio estimator - 1st degree approximation
89   bias2i <- ((1-f)/n) * my * (c2_x - Ro_ZX * c_z * c_x)
90   #Mean Square Error of ratio estimator - 1st degree approximation
91   mse2i <- ((1-f)/n) * (my^2) * (c2_z + c2_x - 2 * Ro_ZX * c_z * c_x)
92
93   #Bias of regression estimator - 1st degree approximation
94   bias3i <- -beta_zx * ((1-f)/n) * ((mu12/mu11) - (mu03/mu02))
95   #Mean Square Error of regression estimator
96   #1st degree approximation

```

```

97 mse3i <- ((1-f)/n)*(my^2)*c2_z*(1-(Ro_ZX^2))
98
99 #Bias of generalized regression-cum-ratio estimator
100 #1st degree approximation
101 bias5i <- (k1-1)*my+k1*my*((1-f)/n)*(c2_x-Ro_ZX*c_z*c_x)
102         +k2*mx*((1-f)/n)*c2_x
103 #Mean Square Error of generalized regression-cum-ratio estimator
104 #1st degree approximation
105 mse5i <- ((k1-1)^2)*(my^2)+(k1^2)*(my^2)
106         *((1-f)/n)*(c2_z+3*c2_x-4*Ro_ZX*c_z*c_x)
107         +(k2^2)*(mx^2)*((1-f)/n)*c2_x-2*k1*(my^2)
108         *((1-f)/n)*(c2_x-Ro_ZX*c_z*c_x)
109         -2*k2*my*mx*((1-f)/n)*c2_x-2*k1*k2*my*mx
110         *((1-f)/n)*(Ro_ZX*c_z*c_x-2*c2_x)
111
112 #Bias of generalized exponential type estimator
113 #1st degree approximation
114 bias7i <- (w1-1)*my+w1*my*1*((3/8)*c2_x-0.5*Ro_ZX*c_z*c_x)
115         +w2*mx*1*c2_x
116 mse7i <- (my^2)*((1-0.25*1*c2_x)-(((1-(1/8)*1*c2_x)^2)
117         /(1+1*c2_z*(1-(Ro_ZX^2))))))
118
119 #Bias of improved exponential estimator 2
120 #1st degree approximation
121 bias8i <- (z1-1)*my+z1*my*((3/8)*1*c2_x-0.5*1*Ro_ZX*c_z*c_x)
122         +z2*(1+(3/8)*1*c2_x)
123 #Mean Square Error of improved exponential estimator 2
124 #1st degree approximation
125 mse8i <- (my^2)*(1-((B*(C^2)+A*(D^2)-2*C*D*E)/(A*B-(E^2))))
126
127 #Empirical results
128 #Simulation of 5000 replicas of estimates
129 ...
130
131 #Graphics
132 #*****
133 emp_bias <- emp-my
134 emp_arb <- abs(emp_bias)/my
135 emp_mse <- apply(emp, 2, var)+(emp_bias^2)
136 emp_res <- rbind(emp_res, rbind(cbind(N, n, Ro_YX, 1,
137     emp_bias[,1], emp_arb[,1], emp_mse[,1]),
138     cbind(N, n, Ro_YX, 2, emp_bias[,2], emp_arb[,2], emp_mse[,2]),
139     cbind(N, n, Ro_YX, 3, emp_bias[,3], emp_arb[,3], emp_mse[,3]),
140     cbind(N, n, Ro_YX, 5, emp_bias[,5], emp_arb[,5], emp_mse[,5]),
141     cbind(N, n, Ro_YX, 7, emp_bias[,7], emp_arb[,7], emp_mse[,7]),
142     cbind(N, n, Ro_YX, 8, emp_bias[,8], emp_arb[,8], emp_mse[,8])))
143 colnames(emp_res) <- c("N", "n", "RhoXY", "EST",
144     "EMP_BIAS", "EMP_ARB", "EMP_MSE")
145 #*****
146

```

```

147   #Results
148   res <- rbind(res, c(N, n, Ro_YX, Ro_ZX,
149                     c_x, c_y, c_z, k1, k2, w1, w2,
150                     z1, z2, mx, my, mz,
151                     med_est1, med_est2, med_est3,
152                     med_est5, med_est7, med_est8,
153                     bias2i, bias3i, bias5i,
154                     bias7i, bias8i,
155                     emp_mse1, mse1, emp_mse2, mse2i,
156                     emp_mse3, mse3i, emp_mse5, mse5i,
157                     emp_mse7, mse7i, emp_mse8, mse8i))
158   }
159   colnames(res) <- c("N", "n", "RhoXY", "RhoZX",
160                    "Cx", "Cy", "Cz", "k1", "k2", "w1", "w2",
161                    "z1", "z2", "mX", "mY", "mZ",
162                    "Est1", "Est2", "Est3",
163                    "Est5", "Est7", "Est8",
164                    "BIAS2I", "BIAS3I", "BIAS5I",
165                    "BIAS7I", "BIAS8I",
166                    "EMP_MSE1", "MSE1", "EMP_MSE2", "MSE2I",
167                    "EMP_MSE3", "MSE3I", "EMP_MSE5", "MSE5I",
168                    "EMP_MSE7", "MSE7I", "EMP_MSE8", "MSE8I")
169   return(list(res=res, emp_res=emp_res))
170 }
171
172 #Package for generation
173 require(MASS)
174 #Import data
175 data_yx <- read.table("ENC0809.txt", sep="\t", dec=",", header = T)
176 #Study variable (orders in 2009, thousands of euros)
177 Y <- data_yx[,3]
178 #Auxiliary variable (orders in 2009, thousands of euros)
179 X <- data_yx[,2]
180
181 #Data application
182 N <- dim(data_yx)[1]
183 res <- comparison_chap7(Y, X, N)
184
185 res_exp <- res[[1]]
186 #Export data
187 write.table(res_exp, "chapter7_ne_results.txt", sep="\t", dec=",", row.names=FALSE)

```