



Universidade Federal de Pernambuco
Centro de Ciências Exatas e da Natureza
Departamento de Estatística

Pós-graduação em Estatística

Inferência sob planos amostrais de cadastro duplo

Hemílio Fernandes Campos Coêlho

Tese de Doutorado

Recife

23 de fevereiro de 2011

Universidade Federal de Pernambuco
Centro de Ciências Exatas e da Natureza
Departamento de Estatística

Hemílio Fernandes Campos Coêlho

Inferência sob planos amostrais de cadastro duplo

Trabalho apresentado ao Programa de Pós-graduação em Estatística do Departamento de Estatística da Universidade Federal de Pernambuco como requisito parcial para obtenção do grau de Doutor em Estatística.

Orientador: Cristiano Ferraz

Recife

23 de fevereiro de 2011

Catálogo na fonte
Bibliotecária Jane Souto Maior, CRB4-571

Coelho, Hemílio Fernandes Campos
Inferência sob planos amostrais de cadastro duplo /
Hemílio Fernandes Campos Coelho - Recife: O Autor, 2011.
xv, 144 p. : il., fig., tab.

Orientador: Cristiano Ferraz.
Tese (doutorado) Universidade Federal de Pernambuco.
CCEN. Estatística, 2011.

Inclui bibliografia e apêndices.

1. Estatística aplicada. 2. Amostragem. 3. Inferência. I.
Ferraz, Cristiano (orientador). II. Título.

310

CDD (22. ed.)

MEI2011 – 024

Universidade Federal de Pernambuco
Pós-Graduação em Estatística

23 de fevereiro de 2011

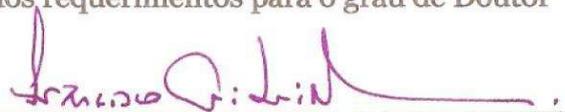
Nós recomendamos que a tese de doutorado de autoria de

Hemílio Fernandes Campos Coêlho

intitulada

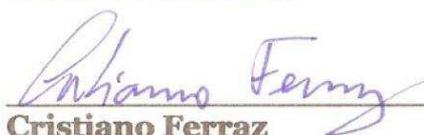
“Inferência sob planos amostrais de cadastro duplo”

seja aceita como cumprimento parcial dos requerimentos para o grau de Doutor em Estatística.



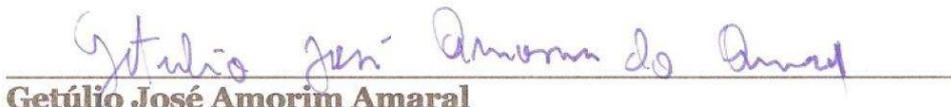
Coordenador da Pós-Graduação em Estatística

Banca Examinadora:



Cristiano Ferraz

orientador



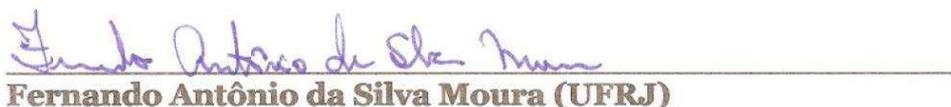
Getúlio José Amorim Amaral



Júlio da Motta Singer (USP)



Marcel de Toledo Vieira (UFJF)



Fernando Antônio da Silva Moura (UFRJ)

Este documento será anexado à versão final da dissertação.

*Dedico este trabalho à minha esposa,
Merciany Rodrigues Ferreira.*

Agradecimentos

A realização deste trabalho é uma prova da Fé em Deus. Uma Fé que nos momentos mais difíceis surge como uma força maior que nos move pra frente sem deixar aumentar a dor nem sangrar as feridas. Fé no Deus que nos mostra que para Ele nada é impossível, e que nos dá condições na vida para enfrentar todas as adversidades, sejam elas quais forem. Portanto, meu primeiro agradecimento é para Deus, por ter me dado uma missão importante neste mundo, que sempre estará condicionada pelas minhas ações.

Aos meus pais, por toda a força e incentivo dada desde a infância, em especial à minha mãe, Wilma, uma das pernambucanas mais fantásticas que existe. Uma guerreira que incondicionalmente luta no dia-a-dia pelo melhor para seus filhos, e que sempre esteve presente nos momentos em que mais precisei de seu apoio para seguir firme nos estudos.

À minha esposa Merciany Rodrigues, por ser meu ágape. Ágape em grego significa o **amor incondicional**. Nos últimos dez anos, desde a graduação até o fim do Doutorado, ela esteve ao meu lado como amiga, namorada, noiva e esposa, sendo a minha fortaleza e ágape. Nos momentos em que apenas a emoção transbordava em mim, ela vinha e se tornava a razão para o meu equilíbrio e perseverança.

Aos meus irmãos, pelos constantes momentos engraçados na vida de irmãos e também por toda a força dada nessa jornada.

Aos meus sogros, Izaías e Mercedes, pelo exemplo de Fé e pelos valorosos conselhos.

À minha avó, Maria José, que carinhosamente chamo de “Mãe Bezinha”, pelas valorosas lições de vida.

Ao meu tio Jadiel Sobreira (in memorian), que mesmo não estando neste mundo há quase seis anos, ainda transmite importantes lições para toda a família pelo exemplo que sempre será em todos os aspectos. Agradeço também a todos os primos que acompanharam minha jornada.

Aos meus tios maternos, em especial Nilma e Lindoval, pela constante alegria e apoio.

Ao professor Cristiano Ferraz, que foi mais que um orientador. O professor Cristiano é um verdadeiro amigo e um exemplo de profissional, pois confiou plenamente na minha capacidade de desenvolver este trabalho e sempre me incentivou a ter perseverança e jamais desistir diante dos difíceis desafios que surgem na vida.

Aos professores Marcel Vieira e Damião Nóbrega, pelas importantes contribuições dadas na qualificação do Doutorado. Ao professor Damião, um grande agradecimento por ter me recebido na UFRN durante dois dias e ter contribuído de forma extraordinária para o desenvolvimento desta tese.

Aos estatísticos Jeffrey Bailey e Chadd Crouse, do National Agricultural Statistics Service of the United States Department of Agriculture (NASS/USDA), pelo material fornecido sobre

o trabalho desenvolvido com cadastro duplo nos Estados Unidos.

Aos grandes amigos da Estatística - UFPE: Abraão David, Carlos Renato, Claudyvan Paiva, Emmanuelle Araújo, Leila Rameh, Luiz Medeiros, Marcelo Rodrigo, Oscar Raposo, Raul Siqueira, Romero Filho, Robson Florêncio, Tadeu Rodrigues, Vanessa Santos, Lucemberg Pedrosa, Camilla Rocha, Hélio, Lucianna, Mariana Batista e Natália Pires, por todos os bons momentos vivenciados.

Aos grandes amigos da pós-graduação em Estatística, em especial Artur Lemonte, Basíli-des, Carlos Raphael, Daniel Neyra, Fábio Bayer, Fábio Fajardo, Nataly Monroy, Jane, Jeremias Leão, Juliana Pires, Luz Marina, Maria Lídia, Rejane Brito, Rita Lima, Sílvio Fernando, Tarciana Liberal, Tatiene Souza, Themis Abensur, e Víncius Quintas. Cada um fez a diferença nos momentos de estudo, alegria, preocupação, perseverança.

Ao mestrando Oscar Lopez da Universidade Nacional da Colômbia, pela disponibilidade e atenção nas discussões sobre o tema desta tese.

Aos amigos professores que tive a honra de conhecer durante o prazeroso período em que trabalhei na UFPI, no saudoso Campus Ministro Reis Velloso, cidade de Parnaíba, Piauí: Fuad Hazime, Patrícia Hazime, André Perinotto, Cíntia Martins, Marcelo Filgueiras, Erik Rodarte, Roberto Ramos, Alexandro Marinho, Marcelo Rêgo, Gustavo Portela, Anna Carolina, Fábio Motta, Marcelo Coerltjens, Patrícia Coerltjens, Carla Eiras, France Keiko, André Luiz, Jand Venes, Valmária, Cleide e Fernando.

Aos professores do Departamento de Estatística da UFPE, em especial Maria Cristina Falcão Raposo, Sylvio José Pereira dos Santos, Klaus Leite Pinto Vasconcellos, Francisco Cribari Neto e Cláudia Regina Oliveira de Paiva Lima, pelo apoio e importantes orientações.

Aos amigos professores do Departamento de Estatística da UFPB, em especial Ulisses Umbelino, Eufrásio Neto e José Carlos, pelos valorosos conselhos e apoio para a finalização do Doutorado.

À secretária da pós, Valéria Bittencourt, uma profissional da mais alta competência e uma grande amiga que tive a honra de conhecer na pós.

Ao secretário da graduação, Lódino Neto, pela amizade e pelas constantes conversas e apostas sobre o futebol pernambucano.

Aos funcionários do Departamento em Estatística, em especial Jimmy, Maurício e Cícero.

Aos meus alunos.

À CAPES, pelo apoio financeiro nos dois anos iniciais do doutorado.

Resumo

A abordagem de cadastro duplo envolve um levantamento amostral onde dois cadastros são utilizados com o propósito de fornecer maior cobertura e identificar elementos de uma única população-alvo. Tal abordagem tem sido vastamente utilizada na literatura em situações onde um único cadastro não consegue fornecer cobertura completa da população-alvo, ou ainda, quando há diferenças de custo de amostragem realizada em cada um dos cadastros disponíveis. Esta tese estuda estratégias de inferência assistida por modelos aplicadas à abordagem de cadastro duplo, propondo estimadores do tipo regressão generalizado. Variáveis auxiliares são consideradas disponíveis em ambos os cadastros e utilizadas no processo de estimação para aumentar a precisão de estimativas. As propriedades de consistência, centralidade assintótica e erro quadrático médio assintótico dos novos estimadores propostos são apresentadas e uma simulação é realizada para analisar sua eficiência relativa a estimadores já propostos na literatura, sob a situação em que um plano de amostragem aleatória simples é aplicado em cada um dos cadastros. Os estimadores propostos têm potencialidade de aplicação direta em diversas áreas, como a de pesquisa agropecuária.

Palavras-chave: Cadastro duplo, estimador regressão, consistência, centralidade assintótica.

Abstract

A dual frame sampling approach is defined as a survey sample design with two sample frames providing coverage for the same target population. Under this context, we propose several generalized regression estimators are proposed and we proved that they are asymptotically design-unbiased, design-consistent and have an asymptotic mean square error under some conditions. Their performance is compared with estimators found in literature and investigated for several dual frame scenarios in a Monte Carlo study for simple random sampling design. The proposed estimators are motivated by their potential application to many areas, such as agricultural research, health and ecology.

key words: Dual frame, generalized regression estimator, design consistency, agricultural research, asymptotic design unbiased, asymptotic mean square error.

“Enquanto eu puder trabalhar eu trabalho, porque o trabalho me distrai, não é sacrifício. Cada problema que aparece é um esforço para resolver que me agrada muito e que ainda consigo fazer.”

–OSCAR NIEMEYER, ao completar 103 anos em 2010.

“Provai e vede como o Senhor é bom. Feliz quem encontra Nele o seu refúgio.”

–SALMO 34.

LISTA DE TABELAS

2.1	Situações possíveis para uso da abordagem de Cadastro Duplo	p. 13
2.2	Notação utilizada na estratégia de Hartley	p. 14
2.3	Valores da razão $\text{Var}(\bar{y}_H) / \text{Var}(\bar{y}')$ quando $\sigma_B^2 / \sigma_a^2 = 16$	p. 18
2.4	Valores da razão $\text{Var}(\bar{y}_H) / \text{Var}(\bar{y}')$ quando $\sigma_B^2 / \sigma_a^2 = 4$	p. 18
6.1	Tamanhos populacionais utilizados para a , b e ab	p. 69
6.2	Avaliação das estratégias de Hartley, Fuller, Bankier e PML para $N = 4500$, $N_{ab} = 500$ e desvios unitários.	p. 74
6.3	Avaliação das estratégias de estimador regressão para $N = 4500$, $N_{ab} = 500$ e desvios unitários.	p. 75
6.4	Avaliação das estratégias de Hartley, Fuller, Bankier e PML para $N = 4500$, $N_{ab} = 1000$ e desvios unitários.	p. 77
6.5	Avaliação das estratégias de estimador regressão para $N = 4500$, $N_{ab} = 1000$ e desvios unitários.	p. 78
6.6	Avaliação das estratégias de Hartley, Fuller, Bankier e PML para $N = 4500$, $N_{ab} = 1500$ e desvios unitários.	p. 80
6.7	Avaliação das estratégias de estimador regressão para $N = 4500$, $N_{ab} = 1500$ e desvios unitários.	p. 81
6.8	Avaliação das estratégias de Hartley, Fuller, Bankier e PML para $N = 4500$, $N_{ab} = 500$ e desvios diferentes.	p. 83
6.9	Avaliação das estratégias de estimador regressão propostos para $N = 4500$, $N_{ab} = 500$ e desvios diferentes.	p. 84
6.10	Avaliação das estratégias de Hartley, Fuller, Bankier e PML para $N = 4500$, $N_{ab} = 1000$ e desvios diferentes.	p. 86
6.11	Avaliação das estratégias de estimador regressão propostos para $N = 4500$, $N_{ab} = 1000$ e desvios diferentes.	p. 87

- 6.12 Avaliação das estratégias de Hartley, Fuller, Bankier e PML para $N = 4500$,
 $N_{ab} = 1500$ e desvios diferentes. p. 89
- 6.13 Avaliação das estratégias de estimador regressão propostos para $N = 4500$,
 $N_{ab} = 1500$ e desvios diferentes. p. 90

LISTA DE FIGURAS

1.1 Cadastro duplo fornecendo cobertura completa para a população-alvo	p. 2
1.2 Cadastro duplo para melhorar a relação custo-benefício	p. 2
1.3 Domínios a , b e ab induzidos pela abordagem de cadastro duplo	p. 3
1.4 Exemplo de seleção de segmento de área	p. 6
1.5 Segmento escolhido e subdivisão do segmento em estratos	p. 6
2.1 Quantidades amostrais geradas pela estratégia de Hartley	p. 14
2.2 Quantidades amostrais geradas pela estratégia BLG	p. 22
4.1 Cenário (a) $\mathcal{A}_1 \subset \mathcal{A}_2$, $\mathcal{B}_1 \subset \mathcal{B}_2$ e $\mathcal{A}_1 \cap \mathcal{B}_1 \subset \mathcal{A}_2 \cap \mathcal{B}_2$	p. 40
6.1 Desvio Padrão dos estimadores propostos para $N = 4500$, $N_{ab} = 500$ e desvios unitários.	p. 76
6.2 Desvio Padrão das estratégias de estimador regressão para $N = 4500$, $N_{ab} = 1000$ e desvios unitários.	p. 79
6.3 Desvio Padrão das estratégias de estimador regressão para $N = 4500$, $N_{ab} = 1500$ e desvios unitários.	p. 82
6.4 Desvio Padrão das estratégias de estimador regressão para $N = 4500$, $N_{ab} = 500$ e desvios diferentes.	p. 85
6.5 Desvio Padrão das estratégias de estimador regressão para $N = 4500$, $N_{ab} = 1000$ e desvios diferentes.	p. 88
6.6 Desvio Padrão das estratégias de estimador regressão para $N = 4500$, $N_{ab} = 1500$ e desvios diferentes.	p. 91
6.7 Densidade estimada dos estimadores propostos para $N = 4500$ e $N_{ab} = 500$	p. 92
6.8 Densidade estimada dos estimadores propostos para $N = 4500$ e $N_{ab} = 1000$	p. 92
6.9 Densidade estimada dos estimadores propostos para $N = 4500$ e $N_{ab} = 1500$	p. 93
6.10 Densidade Estimada de $\hat{N}_{ab,s}$, com valor de referência $N_{ab} = 500$	p. 94
6.11 Densidade Estimada de $\hat{N}_{ab,s}$, com valor de referência $N_{ab} = 1000$	p. 95
6.12 Densidade Estimada de $\hat{N}_{ab,s}$, com valor de referência $N_{ab} = 1500$	p. 96

6.13	Densidade Estimada de $\hat{N}_{ab,PML}$, com valor de referência $N_{ab} = 500$	p. 97
6.14	Densidade Estimada de $\hat{N}_{ab,PML}$, com valor de referência $N_{ab} = 1000$	p. 98
6.15	Densidade Estimada de $\hat{N}_{ab,PML}$, com valor de referência $N_{ab} = 1500$	p. 99
1	Cenário (b) $\mathcal{A}_1 \subset \mathcal{A}_2$ e $\mathcal{B}_1 \subset \mathcal{A}_2$, e $\mathcal{A}_1 \cup \mathcal{B}_1 \subset \mathcal{A}_2$	p. 138
2	Cenário (c) $\mathcal{A}_1 \subset \mathcal{A}_2$ e $\mathcal{A}_1 \cup \mathcal{B}_1 \subset \mathcal{A}_2$, com $(\mathcal{A}_1 \cup \mathcal{B}_1) \cap \mathcal{B}_2 = \emptyset$	p. 138
3	Cenário (d) $\mathcal{A}_1 \subset \mathcal{A}_2 \subset \mathcal{B}_1 \subset \mathcal{B}_2$	p. 139
4	Cenário (e) $(\mathcal{A}_1 \cup \mathcal{A}_2) \subset \mathcal{A}_2$ e $\mathcal{A}_2 \subset \mathcal{B}_2$	p. 139
5	Cenário (f) $(\mathcal{A}_1 \cup \mathcal{B}_1) \subset \mathcal{A}_2 \cap \mathcal{B}_2$	p. 140
6	Cenário (g) $\mathcal{A}_1 \subset \mathcal{B}_1$ e $(\mathcal{A}_1 \cup \mathcal{B}_1) \subset \mathcal{A}_2 \cap \mathcal{B}_2$	p. 140

SUMÁRIO

1	Introdução	p. 1
1	A abordagem de cadastro duplo	p. 1
2	Aplicações	p. 3
2.1	Agricultura	p. 3
2.2	Ecologia	p. 8
2.3	Saúde	p. 9
3	Contribuições desta tese	p. 10
4	Organização da tese	p. 11
2	Inferência em uma abordagem de cadastro duplo	p. 12
1	Estratégia de Hartley(1962)	p. 12
2	Estratégia de Fuller & Burmeister(1972)	p. 19
3	Estratégia de Bankier, Lepkowski & Groves(1986)	p. 22
4	Estratégia de Máxima Pseudo-Verossimilhança	p. 25
5	Estratégia de Singh & Wu (2003)	p. 28
3	Estimadores do tipo regressão em uma abordagem de cadastro duplo	p. 29
1	Introdução	p. 29
2	Estimador regressão sob a estratégia de Hartley	p. 30
3	Estimador regressão sob a estratégia de Fuller & Burmeister	p. 32
4	Estimador regressão sob a estratégia BLG	p. 33
5	Estimador regressão sob a estratégia de Máxima Pseudo-Verossimilhança	p. 36
4	Propriedades Estatísticas dos Estimadores	p. 39
1	Introdução	p. 39

2	Centralidade assintótica e consistência do estimador sob a abordagem de cadastro duplo	p. 41
3	Centralidade assintótica e consistência do estimador regressão generalizado	p. 44
4	Erro quadrático médio assintótico dos estimadores regressão sob a abordagem de cadastro duplo	p. 46
4.1	Estratégia de Hartley	p. 49
4.2	Estratégia de Fuller & Burmeister	p. 51
4.3	Estratégia BLG	p. 53
4.4	Estratégia de Máxima Pseudo-Verossimilhança	p. 54
5	Distribuições assintóticas dos estimadores propostos	p. 55
5	Inferência para casos particulares do estimador regressão generalizado	p. 57
1	Estimadores do tipo razão sob a estratégia de Hartley	p. 57
1.1	Estimador razão RH1	p. 57
1.2	Estimador razão RH2	p. 60
2	Estimadores do tipo razão sob a estratégia de Fuller & Burmeister	p. 61
2.1	Estimador razão FB1	p. 61
2.2	Estimador razão RFB2	p. 63
3	Estimadores do tipo razão sob a estratégia BLG	p. 64
4	Estimadores do tipo razão sob a estratégia de Máxima-Pseudo Verossimilhança	p. 65
6	Avaliação Numérica	p. 69
1	Discussão metodológica	p. 69
2	Resultados da Simulação	p. 71
2.1	Desempenho dos estimadores propostos sob o plano de amostragem aleatória simples	p. 71
2.2	Estudo sobre a distribuição de $\hat{N}_{ab,s}$ e $\hat{N}_{ab,PML}$	p. 73
7	Considerações Finais	p. 100
	Apêndice A - Resumo de conceitos importantes	p. 103
1	Ordens de Magnitude de sequências de números reais (O e o)	p. 103
1.1	Teorema 1	p. 103
2	Desigualdade de Markov	p. 104
3	Convergência Estocástica	p. 104
3.1	Ordens de Magnitude de sequências de números reais (O_p(.) e o_p(.))	p. 104
3.2	Teorema 2	p. 104
4	Tipos de convergência estocástica	p. 105

Apêndice B - Demonstrações dos resultados propostos	p. 106
Apêndice C - Programas utilizados	p. 118
Apêndice D - Cenários de sequências de populações finitas sob a abordagem de cadastro duplo	p. 138

CAPÍTULO 1

Introdução

1 A abordagem de cadastro duplo

A abordagem de cadastro duplo é definida como um levantamento onde dois cadastros, denotados por A e B , são utilizados simultaneamente para identificar elementos de uma população-alvo. Nesta abordagem, amostras aleatórias independentes, denotadas por S_A e S_B , são obtidas de cada cadastro, respectivamente, através de planos amostrais com probabilidades $p(S_A)$ e $p(S_B)$ possivelmente diferentes. Hartley (1962) provavelmente foi o primeiro a apresentar uma estratégia de estimação que trouxe vantagens na relação custo-benefício em favor da abordagem de cadastro duplo. Motivados por esta contribuição inicial, vários autores contribuíram para o desenvolvimento e para a melhoria de estratégias de estimação sob essa abordagem, como Bankier, Lepkowski & Groves (1986), Lund (1968), Fuller & Burmeister (1972), Skinner (1991), Skinner & Rao (1996), Bosecker & Ford (1976), Haines & Pollock (1998), Carfagna (2004), Alpizar-Jara, Pollock & Haines (2005), Lohr (2007) e Lu (2007). No contexto de estratégias de estimação assistida por modelos, a lista de autores parece ser mais reduzida. Singh & Wu (2003) propuseram estimadores do tipo regressão modificados, baseados em métodos de calibração, e Coelho (2007) apresentou estratégias de estimação assistida por modelos baseadas em estimadores do tipo razão.

Existem inúmeras situações em que é possível utilizar a abordagem de cadastro duplo. Por exemplo, quando um dos cadastros não tem um grau de cobertura desejável da população-alvo, porém existe um segundo cadastro, que em conjunto com o primeiro, fornece cobertura completa. Esta situação é representada pela figura 1.1.

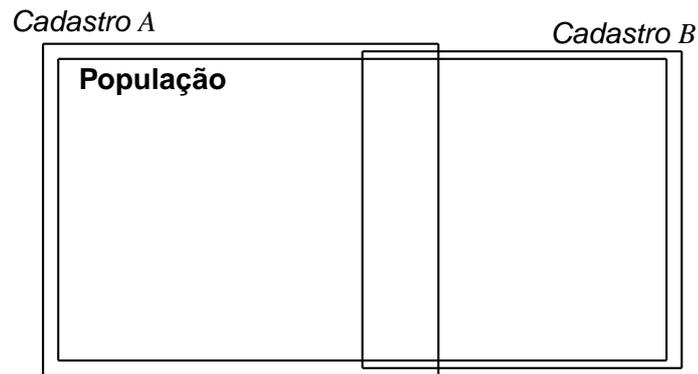


Figura 1.1: Cadastro duplo fornecendo cobertura completa para a população-alvo

Outro exemplo está na situação em que um cadastro cobre por completo a população-alvo, mas tem alto custo de seleção da amostra. Nesse caso, se houver um outro cadastro disponível e de abrangência menor que o primeiro, é possível utilizar ambos os cadastros de modo a tornar vantajosa a relação custo-benefício de seleção de elementos. Esta situação é representada pela figura 1.2.

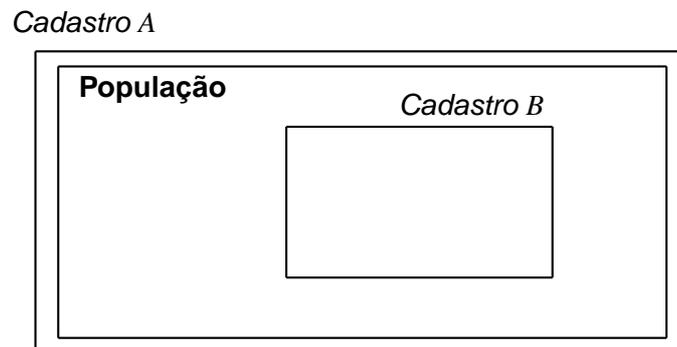


Figura 1.2: Cadastro duplo para melhorar a relação custo-benefício na seleção de elementos da população-alvo

Denote por \mathcal{A} o conjunto de elementos pertencentes ao cadastro A e \mathcal{B} o conjunto de elementos pertencentes ao cadastro B . O uso de dois cadastros induz a existência de no máximo três domínios disjuntos, a , b e ab , tais que $a = \mathcal{A} \cap \mathcal{B}^c$, $b = \mathcal{A}^c \cap \mathcal{B}$ e $ab = \mathcal{A} \cap \mathcal{B}$. Quando $\mathcal{B} \subset \mathcal{A}$ apenas dois domínios são gerados, como ilustra a figura 1.3.(b).

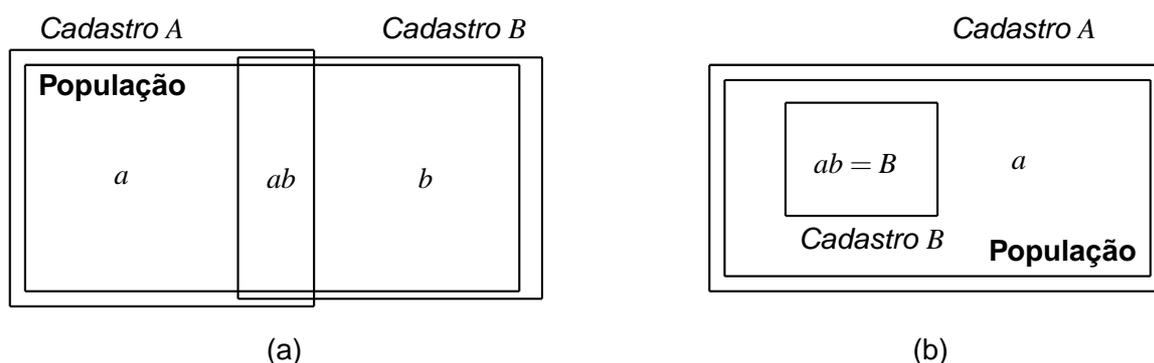


Figura 1.3: Domínios a , b e ab induzidos pela abordagem de cadastro duplo

A abordagem de cadastro duplo é um caso particular da abordagem de cadastros múltiplos em que F cadastros são utilizados ($F \geq 2$). Em geral, tal abordagem implica a existência de $2^F - 1$ domínios possíveis de serem identificados.

2 Aplicações

Situações reais em que a abordagem de cadastro duplo pode ser aplicada são encontradas na literatura, abrangendo diversas áreas de aplicação. A seguir, serão apresentados estudos que utilizaram a abordagem com sucesso ou que ainda estudam as suas potencialidades.

2.1 Agricultura

1. Pesquisa Agropecuária Nacional - Canadá

Armstrong (1979) estudou o uso da abordagem de cadastro duplo na Divisão de Agropecuária e Estatística do Canadá na condução da Pesquisa Agropecuária Nacional, baseada em um cadastro do tipo área. Devido a problemas neste tipo de cadastro, como a vulnerabilidade a altas taxas de não-resposta e a não-atualização ao longo dos anos, era necessário verificar se outras técnicas poderiam ser mais eficientes na produção de estimativas agropecuárias. Para isto, uma das províncias ligadas à produção agrícola, a de New Brunswick no ano de 1978 foi escolhida para a realização de um estudo piloto através da abordagem de cadastro duplo. Foram considerados dois cadastros:

- Um cadastro do tipo *área* (listagem de segmentos de área cobrindo um determinado território), que fornece total cobertura da população-alvo;
- e um cadastro do tipo *específico* (listagem de elementos pertencentes a população da província).

A complexidade encontrada no uso da abordagem neste cenário estava na forma de tratamento do cadastro do tipo área. Por ser um cadastro que contemplava grandes áreas, garantia a total cobertura da população-alvo. Neste cadastro, as grandes fazendas (que possuíam a maioria das informações de interesse) identificadas em censos anteriores sempre eram incluídas na pesquisa com probabilidade 1, para que e na amostra obtida não se corresse o risco das informações destas fazendas gerarem superestimativas dos parâmetros de interesse. Além disso, uma eventual exclusão de uma destas grandes fazendas poderia comprometer seriamente a qualidade das estimativas de interesse. Os demais elementos do cadastro de área eram selecionados com base em planos amostrais complexos, como o de conglomerado em dois estágios.

No estudo, o chamado estimador *screening* foi empregado, e tem este nome porque para ser utilizado, é necessário identificar quais elementos da amostra obtida do cadastro de área estão no cadastro de lista. Após a identificação, os elementos em duplicata são desprezados. Considerando A como sendo o cadastro do tipo área e B o cadastro do tipo específico, tem-se que $B \subset A$ (ver figura 1.3.(b)). Considerando uma amostragem aleatória simples em ambos os cadastros, o estimador *screening* para o total populacional $t_y = t_{ya} + t_{yB}$ é dado por:

$$\hat{t}_y = N_a \bar{y}_a + N_B \bar{y}_B,$$

onde t_{ya} é o total populacional do domínio a , t_{yB} é o total populacional do cadastro B , N_a é o tamanho populacional do domínio a , N_B é o tamanho populacional do cadastro B , \bar{y}_a é a média amostral do domínio a e \bar{y}_B é a média amostral do cadastro B . A variância do estimador *screening* é dada por

$$\text{Var}(\hat{t}_y) = (N_a)^2 \text{Var}(\bar{y}_a) + (N_B)^2 \text{Var}(\bar{y}_B)$$

onde $\text{Var}(\bar{y}_a)$ e $\text{Var}(\bar{y}_B)$ são as variâncias dos estimadores para a média no domínio a e cadastro B , respectivamente. Foram obtidos resultados satisfatórios na produção de estimativas agropecuárias em New Brunswick quando estes foram comparados aos da pesquisa realizada apenas com o cadastro de área. Esta experiência foi uma das primeiras a apresentar potencial de aplicação da abordagem de cadastro duplo em pesquisas agropecuárias.

2. National Agricultural Statistics Service (NASS-USA)

As pesquisas do setor agropecuário nos Estados Unidos são desenvolvidas pelo *United States Department of Agriculture* (USDA) através do *National Agricultural Statistics Service* (NASS), que fornece oficialmente toda a informação estatística referente à agricultura. O NASS utiliza métodos de pesquisa baseados em múltiplos cadastros. Desde 1950, abordagens como a de cadastro duplo fornecem grande contribuição para as pesquisas agropecuárias do país.

O programa desenvolvido pelo NASS, intitulado *Multiple Frame Agricultural Survey*, foi iniciado em 1986 em 27 estados daquele país e conduzido nos meses de junho, julho, agosto, setembro, janeiro e março, com o intuito de obter estimativas referentes a áreas produtivas, número de unidades agrícolas em atividade, área total ocupada por unidades agrícolas, produção total das unidades agrícolas, dentre outros. São considerados os seguintes tipos de cadastros:

- Cadastro de área;
- Cadastros específicos (listagens de elementos da população-alvo).

O cadastro do tipo área é resultante da combinação dos cadastros de área construídos para cada estado, independentemente. A cada 10 ou 20 anos o cadastro é atualizado. Já os cadastros específicos são regularmente atualizados, e utilizados em combinação com o cadastro de área. Existem aproximadamente 2 milhões de unidades agrícolas nos Estados Unidos, e os cadastros específicos, quando construídos em grande escala, conseguem uma cobertura de aproximadamente 1,1 milhões de unidades agrícolas.

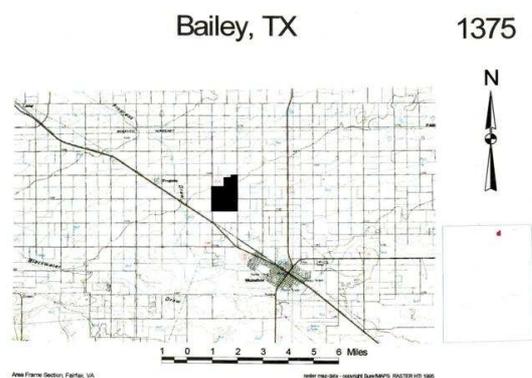
As informações coletadas para composição do cadastro específico são divididas em duas partes. A primeira é referente aos dados para identificar, localizar e contactar fazendeiros e empresários do setor: nome, endereço, telefone, estado, distrito, município, razão social, número de identificação de empregador, etc. A segunda parte contém dados sobre a fazenda e/ou negócio. Uma vez construído, o cadastro é constantemente atualizado através de pesquisas feitas ao longo do ano, e através do censo agropecuário. O cadastro específico é construído com múltiplos propósitos, dentre os quais é possível citar os seguintes:

- i.* identificação e classificação de fazendeiros e/ou empresários do setor que estão diretamente ligados a determinados produtos de interesse;
- ii.* estratificação para identificar grupos de fazendeiros e produtores que possuem determinado produto em comum, para que se possa coletar amostras em cada grupo;
- iii.* pesquisas específicas, ou seja a coleta de dados é feita de forma a atender uma população-alvo em particular.
- iv.* possibilidade de se atribuir probabilidades de seleção para cada elemento do cadastro

O cadastro do tipo área do NASS é construído a partir de toda a área dos Estados Unidos, que é dividida em segmentos de terra através de delimitações, associando fazendas, produção, animais, etc. com a área pertencente a cada segmento. Imagens de satélite, mapas topográficos, sistemas de informação geográfico e fotos aéreas também são utilizados para construção dos cadastros de área. A vantagem desse cadastro está na cobertura completa de

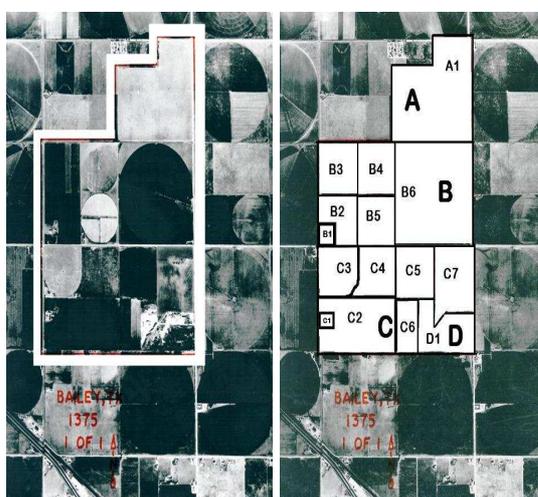
todas as unidades agrícolas em atividade. As figuras (1.4) e (1.5) ilustram um exemplo de segmento de área selecionado pelo NASS. Com o passar dos anos, novas tecnologias ajudaram a construção dos cadastros do tipo área.

Figura 1.4: Exemplo de seleção de segmento de área



Fonte: National Agricultural Statistics Service (NASS)

Figura 1.5: Segmento escolhido e subdivisão do segmento em estratos



Fonte: National Agricultural Statistics Service (NASS)

3. Pesquisas Agropecuárias do Departamento Administrativo Nacional de Estatística (Colômbia)

As pesquisas do setor agropecuário na Colômbia são desenvolvidas pelo *Departamento Administrativo Nacional de Estadística* (DANE) em função de seu papel como coordenador do *Sistema Estadístico Nacional* (SEN) daquele país. O DANE trabalha atualmente com o objetivo de fortalecer e consolidar o SEN através das chamadas “estatísticas estratégicas”, na consolidação das informações obtidas através das pesquisas agropecuárias e adoção de instrumentos de coleta e iniciativas mais eficientes no sentido de melhorar a qualidade da informação, sua disponibilidade e acessibilidade. Tudo isto é focado na necessidade de obtenção de informações estatísticas que atendam à demanda e aos múltiplos propósitos das pesquisas implementadas na região.

A abordagem de cadastro duplo não está sendo utilizada no momento nas pesquisas colombianas. Porém está sendo considerada como possível alternativa para melhorar a relação custo-benefício. López (2010) talvez seja o primeiro a apresentar um estudo sobre a abordagem naquele país, pois utilizou os dados provenientes da chamada *Encuesta Nacional Agropecuaria* (ENA), pesquisa realizada anualmente para analisar a eficiência de estimadores propostos na literatura diante dos múltiplos propósitos das pesquisas agropecuárias realizadas na Colômbia, como:

- Estimar a área cultivada para os principais produtos de caráter permanente ou temporário;
- Determinar a implantação de políticas voltadas para o meio-ambiente;
- Estabelecer as situações em que será necessária uma intervenção técnica;
- Determinar o uso e aproveitamento do solo;
- Conhecer a tecnologia utilizada nos sistemas produtivos dos principais produtos agrícolas (provisório ou permanentes)
- Manejo tecnológico do gado e controle da produção de leite.

2.2 Ecologia

1. Estimação de ninhos de águias norte-americanas

Haines & Pollock (1998), interessados no monitoramento de animais silvestres, consideraram a abordagem de cadastro duplo para estimar o total de ninhos bem-sucedidos (ou seja, que contém pelo menos um casal com um filhote) da ave símbolo e uma das mais conhecidas dos Estados Unidos, a águia norte-americana, que se encontra ameaçada de extinção. O monitoramento da espécie é importante pois a quantidade de ninhos serve como indicador para a qualidade do ecossistema, e até mesmo para a qualidade da água da região de estudo, uma vez que a espécie constrói ninhos em regiões com grande quantidade de água de boa qualidade. A águia norte-americana é muito sensível a toxinas e a presença destas no habitat da espécie é prejudicial dependendo do tipo de presa que capturam. Por outro lado, os filhotes que ainda não estão prontos para deixar os ninhos são bastante prejudicados. Em regiões onde são encontrados ninhos vazios, análises foram feitas e infelizmente nestes ninhos foram encontrados resíduos destas toxinas. Ovos mais frágeis e baixa produção de ovos são exemplos de informações obtidas a partir da análise dos ninhos desta espécie de ave.

O conhecimento de informações sobre a qualidade do ecossistema utilizando essa espécie de águia é baseado necessariamente em várias variáveis, como mortalidade, quantidade de determinado tipo de toxina e número de ovos por ninho por exemplo. A necessidade de uso da abordagem de cadastro duplo neste caso foi importante porque dois tipos de cadastros estavam disponíveis:

- Um cadastro *específico*, que continha informações como a localização dos ninhos, por exemplo. Apesar de bastante eficiente, este tipo de cadastro pode ser desatualizado com facilidade, principalmente se forem levadas em consideração variáveis como tempo, uma vez que novos ninhos podem vir a ser construídos em épocas distintas.
- Um cadastro de área, que continha informações sobre os limites geográficos das regiões da espécie, possuindo a vantagem de fornecer completa cobertura da população de interesse e a desvantagem de proporcionar alto custo de seleção de elementos para amostras a serem coletadas.

Com o intuito de minimizar as deficiências individuais destes cadastros, a abordagem de cadastro duplo foi utilizada com o objetivo de estimar o total de ninhos de águia, e apresentou resultados satisfatórios. Neste estudo, sempre foi assumido que o cadastro de área estava completo, e o que era sempre atualizado era o cadastro tipo específico. Além disso, o uso desta abordagem revelou outra regra importante no comportamento da espécie: a desconsideração da teoria que estabelecia a permanência de um casal de aves em um ninho construído. Ficou comprovado que um casal de aves mantém vários ninhos em seu território, revelando que a

variável *número de territórios ocupados* se tornou mais efetiva do que o número de ninhos bem-sucedidos.

2.3 Saúde

1. Pesquisa Nacional de Idosos (EUA)

Choudhry, Park & Li (2002) consideraram a abordagem de cadastro duplo a partir de informações da Pesquisa Nacional de Idosos realizada em 2000, conduzida para obter estimativas da população de idosos. A pesquisa foca, em especial, grupos populacionais de interesse do Departamento de Assistência ao Idoso (*Department of Veteran Affairs - DVA*) dos Estados Unidos, de modo a estruturar o planejamento de programas destinados à serviços de assistência geral para idosos.

O DVA tinha interesse em obter estimativas dos grupos populacionais em questão (mulheres, afro-americanos, e latino-americanos). A população-alvo considerada na pesquisa foi todo o conjunto de idosos residentes nos Estados Unidos e em Porto Rico, e com base em critérios determinados pelo DVA, a população é então estratificada com base nos grupos populacionais de interesse. Existiam dois cadastros disponíveis:

- Um cadastro específico, com números de telefones de idosos identificados de acordo com Casady & Lepkowski (1991); Casady & Lepkowski (1993) e Potter *et al.* (1991).
- Um cadastro específico constituído pela combinação de duas fontes administrativas, a primeira com informações de idosos cadastrados nos serviços de assistência à saúde e a segunda com informações do registro nacional de seguridade social e pensões.

O primeiro cadastro, apesar de conter a quase totalidade das informações, em geral apresentava grande taxa de não-resposta ou de não-conclusão de uma entrevista, e alto custo para obtenção das informações, uma vez que se faz necessário um tempo elevado ao telefone para realização das entrevistas. Já o segundo cadastro, apesar de incompleto, não apresentava as deficiências do primeiro. Com isso, a ideia dos autores foi combinar ambos os cadastros de modo a alcançar uma maior precisão nas estimativas de interesse, com custo o mais reduzido possível. Ao utilizar a abordagem de cadastro duplo, foi verificada uma redução significativa do custo quando se comparou a abordagem com um plano amostral aplicado apenas no primeiro cadastro.

2. Pesquisa Nacional de Imunização (EUA)

Shin, Molinari & Wolter (2008) consideraram a abordagem de cadastro duplo na Pesquisa Nacional de Imunização (National Immunization Survey - NIS) realizada pelo Centro de Pesquisas de Opiniões Nacionais (National Opinion Research Center - NORC), de modo a fornecer informações para os centros de Controle e Prevenção de Doenças (Center for Disease Control and Prevention - CDC). Este tipo de pesquisa monitora a campanha de vacinação em crianças de 19 a 35 meses de vida. A cada ano, esta pesquisa realiza uma série de aproximadamente 24000 entrevistas em domicílios selecionados. Nesta pesquisa dois cadastros foram utilizados:

- Um cadastro de telefones registrados em um ou mais bancos de telefones do país;
- Um cadastro específico de domicílios que continham pelo menos uma criança.

De todos os telefones cadastrados, menos de 25% destes eram associados a domicílios e aproximadamente 68% dos telefones em funcionamento estavam presentes neste cadastro. Os autores encontraram evidência de que a aplicação da abordagem de cadastro duplo na pesquisa NIS representa uma alternativa de baixo custo quando comparada à pesquisa feita apenas no cadastro A, principalmente em situações em que é necessário utilizar um plano de amostragem estratificada com base no custo de seleção de unidades amostrais em cada estrato.

3 Contribuições desta tese

Quando uma ou mais variáveis auxiliares estão disponíveis em uma abordagem com um único cadastro, seu uso pode melhorar consideravelmente as estimativas dos parâmetros, através de estimadores do tipo regressão. Coelho (2007) mostrou que, como era de se esperar, o uso de estimadores do tipo regressão também pode melhorar consideravelmente estimativas em uma abordagem de cadastro duplo.

O objetivo desta tese de doutorado é apresentar estratégias de inferência utilizando formas gerais de estimadores do tipo regressão, adaptados à abordagem de cadastro duplo e compará-las de modo a informar qual destas estratégias tem melhor desempenho em situações onde a informação de uma ou mais variáveis auxiliares está disponível. Estas estratégias permitem propor diferentes tipos de estimadores de regressão, dentre os quais os propostos por Coelho (2007) são casos particulares. A investigação conduzida contempla a proposição de novos estimadores, bem como o estudo de suas propriedades estatísticas, com respeito a variância, consistência com respeito ao plano amostral, centralidade assintótica e erro quadrático médio assintótico, e distribuição assintótica. Um estudo de simulação foi realizado para avaliar os estimadores sob o caso específico de amostragem aleatória simples empregada nos dois cadastros, através do método de simulação de Monte Carlo.

4 Organização da tese

A presente tese está dividida em sete capítulos. Neste primeiro capítulo foi apresentada uma introdução sobre a abordagem de cadastro duplo. O capítulo 2 apresenta uma revisão das principais estratégias de estimação sob a abordagem de cadastro duplo já propostas na literatura, as quais não utilizam a informação de variáveis auxiliares no processo de estimação. O capítulo 3 apresenta, como contribuição original, estimadores do tipo regressão para a abordagem de cadastro duplo e suas propriedades de variância. O capítulo 4 apresenta, ainda como contribuição original desta tese, o estudo da consistência, centralidade assintótica, erro quadrático médio assintótico e distribuição assintótica dos estimadores propostos. O capítulo 5 apresenta casos particulares dos estimadores apresentados no capítulo 3 e suas propriedades de variância sob os planos de amostragem aleatória simples e estratificada. O capítulo 6 apresenta a avaliação computacional dos estimadores propostos, através do método de simulação de Monte Carlo. O ambiente de programação, análise de dados e gráficos R em sua versão 2.12.0 foi a plataforma computacional escolhida para a avaliação computacional. Por fim, no capítulo 7 são apresentadas as considerações finais deste trabalho.

CAPÍTULO 2

Inferência em uma abordagem de cadastro duplo

1 Estratégia de Hartley(1962)

Hartley (1962) avaliou a eficiência de uso da abordagem de cadastro duplo em um estudo realizado no *Statistical Laboratory* da *Iowa State University* sobre efeitos da industrialização na agricultura, sob a perspectiva de um plano de amostragem aleatória simples aplicado em cada cadastro. Os responsáveis pela pesquisa eram do Departamento de Economia e Sociologia da mesma universidade. Os cadastros utilizados nessa pesquisa foram os seguintes:

- Um cadastro do tipo área rural, utilizado para seleção de produtores rurais;
- Um cadastro do tipo específico, constituído por empregados de um empresa automotiva e que também eram produtores rurais.

O uso simultâneo destes cadastros para a cobertura da população-alvo forneceu uma experiência bem sucedida no que diz respeito ao custo-benefício na seleção de elementos. De modo a facilitar a compreensão do leitor a respeito do processo de inferência sob a abordagem de cadastro duplo, a descrição do estimador de Hartley (1962) será realizada de forma detalhada. Considere que dois cadastros A e B estejam disponíveis e que juntos forneçam completa cobertura para a população-alvo, ilustrado pela figura 1.3(a). Para a implementação da abordagem, duas condições são necessárias:

1. Todos os elementos da população-alvo devem ser identificados pela união dos cadastros;
2. Pode-se verificar se qualquer elemento de um cadastro pertence ou não ao outro.

Denotando por $U = \mathcal{A} \cup \mathcal{B}$ o conjunto de elementos da população-alvo e considerando a notação de domínios apresentada na figura 1.3(a), Hartley (1962) faz uma descrição dos possíveis cenários gerados por um levantamento amostral realizado sob a abordagem de cadastro duplo. Estes cenários dependem basicamente da disponibilidade de informações populacionais, como mostra a tabela (2.1) a seguir.

Tabela 2.1: Situações possíveis para uso da abordagem de Cadastro Duplo

Tipo de Informação Disponível	CENÁRIOS			
	1	2	3	4
Tamanho dos domínios e dos cadastros	Tamanhos dos domínios e cadastros conhecidos	Tamanhos dos domínios e cadastros conhecidos	Apenas os tamanhos dos cadastros conhecidos	Apenas as magnitudes relativas dos cadastros conhecidas
Possibilidade de alocação da amostra	Alocação de amostra aos domínios	Alocação de amostra aos cadastros	Alocação de amostra aos cadastros	Alocação de amostra aos cadastros

É possível notar que $U = a \cup b \cup ab$, como mostra a figura 1.3(a). Os cenários 1 e 2 ilustram a situação em que os tamanhos dos domínios são conhecidos, o que implica que os tamanhos populacionais nos cadastros também o são. O cenário 1 permite que a amostra seja alocada a cada domínio, o que ilustra o caso de um plano sob estratificação, enquanto que o cenário 2 admite a alocação da amostra apenas aos cadastros. Os cenários 3 e 4 só permitem que a alocação da amostra seja feita nos cadastros. No cenário 3, apenas os tamanhos populacionais nos cadastros são conhecidos, enquanto o cenário 4 apresenta como única informação disponível o tamanho relativo dos cadastros, fornecendo informações bastante restritivas, tornando possível apenas a estimação de médias populacionais. Assim, todo o desenvolvimento nesta tese será referente aos cenários 2 e 3. A tabela (2.2) fornece a notação utilizada para cálculos feitos a partir das amostras obtidas dos cadastros A e B . Nela, as quantidades $n_{ab(A)}$, $\hat{t}_{yab(A)}$ e $\tilde{y}_{ab(A)}$ são funções de elementos na amostra que pertencem ao domínio ab e que foram selecionados do cadastro A . As quantidades $n_{ab(B)}$, \hat{t}_{yab}^B e \tilde{y}_{ab}^B são definidas analogamente e estão relacionadas ao cadastro B .

Tabela 2.2: Notação utilizada na estratégia de Hartley

Quantidades de Interesse	CADASTRO		Domínio		
	<i>A</i>	<i>B</i>	<i>a</i>	<i>b</i>	<i>ab</i>
População	<i>A</i>	<i>B</i>	<i>A</i>	<i>B</i>	$A \cap B$
Tamanho da População	N_A	N_B	N_a	N_b	N_{ab}
Total populacional	t_{yA}	t_{yB}	t_{ya}	t_{yb}	t_{yab}
Média populacional	μ_{yA}	μ_{yB}	μ_{ya}	μ_{yb}	μ_{yab}
Amostra	S_A	S_B	S_A	S_B	S_{ab}
Tamanho da amostra	n_A	n_B	n_a	n_b	$n_{ab(A)}$ $n_{ab(B)}$
Estimador do Total	\hat{t}_{yA}	\hat{t}_{yB}	\hat{t}_{ya}	\hat{t}_{yb}	\hat{t}_{yab}^A \hat{t}_{yab}^B
Média amostral	\bar{y}_A	\bar{y}_B	\bar{y}_a	\bar{y}_b	\bar{y}_{ab}^A \bar{y}_{ab}^B

Considerando a implementação de um plano de amostragem aleatória simples em cada cadastro, por exemplo, a ideia de Hartley pode ser representada pela figura a seguir.

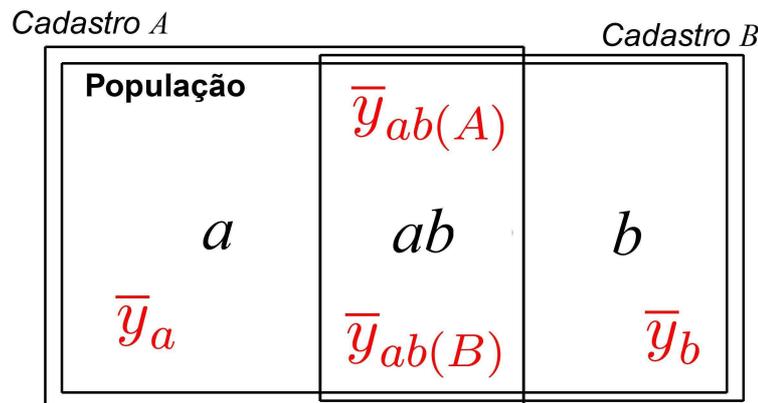


Figura 2.1: Quantidades amostrais geradas pela estratégia de Hartley

Nesta tese, a menos que seja dito o contrário, y_k será sempre considerado como o valor da variável de interesse y associada ao elemento k na população-alvo. O método de estimação desenvolvido por Hartley (1962), que considera o cenário 2, é proposto com base em uma variável y_k^* definida da seguinte forma:

$$y_k^* = \begin{cases} y_k, & \text{se } k \in a \text{ ou } k \in b \\ py_k, & \text{se } k \in ab \text{ e } k \in \mathcal{A} \\ (1-p)y_k, & \text{se } k \in ab \text{ e } k \in \mathcal{B}, \end{cases}$$

onde p é uma constante de ponderação para identificação dos elementos de cada cadastro na população, tal que $0 \leq p \leq 1$. Dessa forma, é possível reescrever o total populacional, t_y , da

seguinte maneira:

$$\begin{aligned}
 t_y &= t_{ya} + t_{yb} + t_{yab} \\
 &= t_{ya} + t_{yb} + (p + 1 - p)t_{yab} \\
 &= t_{ya} + pt_{yab} + t_{yb} + (1 - p)t_{yab} \\
 &= \sum_{k \in \mathcal{A}} y_{kA}^* + \sum_{k \in \mathcal{B}} y_{kB}^* \\
 &= t_{yA}^* + t_{yB}^*. \tag{2.1}
 \end{aligned}$$

Observando a estrutura da expressão (2.1), Hartley enfatiza o fato de que t_y pode ser estimado utilizando estimadores para cada um dos domínios a , b e ab . É possível também representar t_y ou a média populacional, μ , da seguinte forma:

$$\begin{aligned}
 t_y &= t_{yA} + t_{yB} - t_{yab} \\
 &= t_{yA} + t_{yB} - (pt_{yab} + (1 - p)t_{yab}) \tag{2.2}
 \end{aligned}$$

$$\mu = N^{-1}t_y = N^{-1}(t_{yA} + t_{yB} - t_{yab}), \tag{2.3}$$

onde $N = N_A + N_B - N_{ab}$. Nesta tese, o estimador de Hartley será motivado pelas expressões (2.2) e (2.3). Dessa forma, considerando \hat{t}_{yA} , \hat{t}_{yB} e \hat{t}_{yab} como estimadores do total para os cadastros A , B e para o domínio ab , respectivamente, e denotando por \bar{y}_H o estimador de Hartley para μ , tem-se que

$$\bar{y}_H = N^{-1}(\hat{t}_{yA} + \hat{t}_{yB} - \hat{t}_{yab}). \tag{2.4}$$

Quando um plano de amostragem aleatória simples é empregado em cada cadastro (caso originalmente considerado em Hartley (1962)),

$$\hat{t}_{yA} = N_A \bar{y}_A, \quad \hat{t}_{yB} = N_B \bar{y}_B \quad \text{e} \quad \hat{t}_{yab} = pN_{ab} \bar{y}_{ab(A)} + (1 - p)N_{ab} \bar{y}_{ab(B)}.$$

Ao substituir estas quantidades em (2.4), a expressão do estimador de Hartley é

$$\bar{y}_H = N^{-1} \left[N_A \bar{y}_A + N_B \bar{y}_B - \left(pN_{ab} \bar{y}_{ab(A)} + (1 - p)N_{ab} \bar{y}_{ab(B)} \right) \right]$$

$$= w_A \bar{y}_A + w_B \bar{y}_B - \left(p w_{ab} \bar{y}_{ab(A)} + (1-p) w_{ab} \bar{y}_{ab(B)} \right), \quad (2.5)$$

onde $w_A = N_A/N$, $w_B = N_B/N$ e $w_{ab} = N_{ab}/N$. Em uma situação onde estimadores de Horvitz-Thompson são utilizados para estimar as quantidades populacionais nos cadastros e domínios, tem-se que

$$\hat{t}_{yA} = \sum_{k \in S_A} \frac{y_k}{\pi_k}, \quad \hat{t}_{yB} = \sum_{k \in S_B} \frac{y_k}{\pi_k} \quad \text{e} \quad \hat{t}_{yab} = p \sum_{k \in S_A \cap ab} \frac{y_k}{\pi_k} + (1-p) \sum_{k \in S_B \cap ab} \frac{y_k}{\pi_k}.$$

Substituindo estes valores em (2.4), obtém-se

$$\begin{aligned} \bar{y}_H &= N^{-1} \left[\sum_{k \in S_A} \frac{y_k}{\pi_k^A} + \sum_{k \in S_B} \frac{y_k}{\pi_k^B} - \left(p \sum_{\substack{k \in ab \\ k \in S_A}} \frac{y_k}{\pi_k^A} + (1-p) \sum_{\substack{k \in ab \\ k \in S_B}} \frac{y_k}{\pi_k^B} \right) \right] \\ &= N^{-1} \left[N_A \sum_{k \in S_A} \frac{y_k}{N_A \pi_k^A} + N_B \sum_{k \in S_B} \frac{y_k}{N_B \pi_k^B} - \left(p N_{ab} \sum_{\substack{k \in ab \\ k \in S_A}} \frac{y_k}{N_{ab} \pi_k^A} + (1-p) \sum_{\substack{k \in ab \\ k \in S_B}} \frac{y_k}{N_{ab} \pi_k^B} \right) \right] \\ &= w_A \bar{y}_{HT(A)} + w_B \bar{y}_{HT(B)} - \left(p w_{ab} \bar{y}_{HT(ab)}^A + (1-p) w_{ab} \bar{y}_{HT(ab)}^B \right), \end{aligned} \quad (2.6)$$

onde

- π_k^A é a probabilidade de inclusão do elemento k do cadastro A na amostra;
- π_k^B é a probabilidade de inclusão do elemento k do cadastro B na amostra;
- $\bar{y}_{HT(A)}$ é o estimador de Horvitz-Thompson para a média populacional do cadastro A ;
- $\bar{y}_{HT(B)}$ é o estimador de Horvitz-Thompson para a média populacional do cadastro B ;
- $\bar{y}_{HT(ab)}^A$ é o estimador de Horvitz-Thompson para a média populacional do domínio ab , no cadastro A e
- $\bar{y}_{HT(ab)}^B$ é o estimador de Horvitz-Thompson para a média populacional do domínio ab , no cadastro B .

Nota-se que o estimador de Hartley para a média populacional μ é uma média ponderada entre os estimadores de Horvitz-Thompson para a média em cada cadastro e também para a interseção entre eles.

As variâncias aproximadas de \bar{y}_H e dos outros estimadores que serão apresentados nesta tese podem ser expressas em termos matriciais na forma

$$A\text{Var}(\hat{\theta}) = \mathbf{d}^T \Sigma \mathbf{d}, \quad (2.7)$$

onde \mathbf{d}^T é vetor coluna e Σ é a matriz de variâncias e covariâncias do estimador $\hat{\theta}$, diagonal em blocos, pois planos amostrais são aplicados independentemente em cada cadastro. Os elementos de Σ são representados pelas variâncias dos estimadores referentes a cadastros e domínios para cada estimador. Para o estimador de Hartley, tem-se que

$$\mathbf{d} = \begin{bmatrix} w_A \\ pw_{ab} \\ w_B \\ (1-p)w_{ab} \end{bmatrix} \quad \text{e} \quad \Sigma_H = \begin{pmatrix} \Sigma_A & \mathbf{0} \\ \mathbf{0} & \Sigma_B \end{pmatrix} = \begin{pmatrix} \sigma_A^2 & \sigma_{A;ab(A)} & 0 & 0 \\ \sigma_{A;ab(A)} & \sigma_{ab(A)}^2 & 0 & 0 \\ 0 & 0 & \sigma_B^2 & \sigma_{B;ab(B)} \\ 0 & 0 & \sigma_{B;ab(B)} & \sigma_{ab(B)}^2 \end{pmatrix}.$$

Os termos σ_A^2 e σ_B^2 representam as variâncias dos estimadores de Horvitz-Thompson para a média populacional nos cadastros A e B . Por exemplo,

$$\sigma_A^2 = \frac{1}{N_A^2} \sum_{k \in \mathcal{A}} \sum_{l \in \mathcal{A}} \Delta_{kl}^A \frac{y_k}{\pi_k^A} \frac{y_l}{\pi_l^A},$$

com $\Delta_{kl}^A = \pi_{kl}^A - \pi_k^A \pi_l^A$, sendo π_{kl}^A a probabilidade de inclusão de segunda ordem dos elementos k e l do cadastro A . $\sigma_{ab(A)}^2$ e $\sigma_{ab(B)}^2$ representam as variâncias dos estimadores para a média populacional do domínio ab nos cadastros A e B respectivamente, e $\sigma_{A;ab(A)}$ $\sigma_{B;ab(B)}$ representam as covariâncias entre os estimadores calculados dentro dos cadastros A e B respectivamente. Por exemplo,

$$\sigma_{A;ab(A)} = \text{Cov}(\bar{y}_{HT(A)}; \bar{y}_{HT(ab)}^A) = \frac{1}{N_A} \frac{1}{N_{ab}} \sum_{k \in \mathcal{A}} \sum_{l \in ab} \Delta_{kl}^A \frac{y_k}{\pi_k^A} \frac{y_l}{\pi_l^A}.$$

Σ_H é uma matriz diagonal em blocos, pois planos amostrais são aplicados independentemente em cada cadastro. Hartley mostrou ainda que o valor de p que minimiza a expressão (2.7) é dado por

$$p = \frac{N_b (f^B)^{-1} w_B^{ab} \sigma_{ab(B)}^2}{N_a (f^A)^{-1} w_A^{ab} \sigma_{ab(A)}^2 + N_b (f^B)^{-1} w_B^{ab} \sigma_{ab(B)}^2}, \quad (2.8)$$

onde $w_A^{ab} = N_{ab}/N_A$, $w_B^{ab} = N_{ab}/N_B$, $f_A = n_A/N_A$ e $f_B = n_B/N_B$. Quando o valor de ao menos uma das covariâncias apresentadas for muito alto, o valor ótimo de p (2.8) pode ficar fora do intervalo $[0, 1]$. Um estimador não-viesado para $\text{Var}(\bar{y}_H)$ pode ser obtido ao substituir as

quantidades populacionais pelos respectivos estimadores de Horvitz-Thompson.

Hartley (1962) mostrou que a abordagem de cadastro duplo é mais eficiente do que a abordagem de selecionar uma amostra de um único cadastro, para o caso em que $B \subset A$. O critério escolhido para a comparação foi a razão

$$\text{Var}(\bar{y}_H) / \text{Var}(\bar{y}')$$

onde \bar{y}' é o estimador para a média populacional do cadastro A e \bar{y}_H é o estimador de Hartley para μ sob a abordagem de cadastro duplo. Considerando c_A e c_B como os custos de seleção de elementos dos cadastros A e B respectivamente, e σ_B^2 e σ_a^2 como as variâncias dos estimadores obtidos no cadastro B e domínio a respectivamente, as tabelas a seguir apresentam a redução da variância do estimador de Hartley em função de c_B/c_A e N_B/N_A , para valores fixos de σ_B^2/σ_a^2 .

Tabela 2.3: Valores da razão $\text{Var}(\bar{y}_H) / \text{Var}(\bar{y}')$ quando $\sigma_B^2/\sigma_a^2 = 16$

c_B/c_A	N_B/N						
	0.5	0.6	0.7	0.8	0.9	0.95	1
0.01	0.096	0.076	0.059	0.045	0.031	0.024	0.010
0.05	0.154	0.134	0.118	0.102	0.086	0.075	0.050
0.10	0.206	0.188	0.174	0.160	0.143	0.131	0.100
0.20	0.288	0.278	0.269	0.261	0.248	0.237	0.200
0.30	0.359	0.356	0.355	0.353	0.347	0.338	0.300
0.40	0.423	0.428	0.435	0.440	0.441	0.436	0.400
0.50	0.483	0.496	0.510	0.524	0.533	0.532	0.500
1.00	0.096	0.784	0.836	0.889	0.944	0.972	1.000

Tabela 2.4: Valores da razão $\text{Var}(\bar{y}_H) / \text{Var}(\bar{y}')$ quando $\sigma_B^2/\sigma_a^2 = 4$

c_B/c_A	N_B/N						
	0.5	0.6	0.7	0.8	0.9	0.95	1
0.01	0.259	0.201	0.152	0.108	0.066	0.044	0.010
0.05	0.340	0.284	0.234	0.186	0.137	0.107	0.050
0.10	0.404	0.352	0.304	0.257	0.205	0.172	0.100
0.20	0.500	0.456	0.415	0.372	0.322	0.287	0.200
0.30	0.576	0.540	0.507	0.472	0.426	0.393	0.300
0.40	0.640	0.613	0.588	0.561	0.523	0.493	0.400
0.50	0.696	0.678	0.661	0.642	0.614	0.589	0.500
1.00	0.900	0.914	0.932	0.953	0.976	0.988	1.000

Por exemplo, para $\sigma_B^2/\sigma_a^2 = 4$, na situação em que $c_B/c_A = 0.30$ e $N_B/N = 0.7$, observa-se que o uso da abordagem de cadastro duplo comparada à abordagem de uso de apenas um cadastro apresenta um valor da razão igual a 0.507.

Exemplo de aplicação da abordagem de Hartley

Nesta seção será apresentado um exemplo de estimação através da estratégia de Hartley. Considere uma população com $N = 169587$ elementos gerados aleatoriamente. De modo a exemplificar o uso da abordagem, foram considerados cadastros A e B , de acordo com a situação descrita pela figura 1.3(a), da seguinte forma:

1. Cadastro A : constituído do elemento 1 ao elemento 84793, com $N_A = 84793$;
2. Cadastro B : constituído do elemento 74793 ao elemento 169587, com $N_B = 84794$.
3. Domínio ab : constituído por $N_{ab} = 10000$ elementos, comuns aos cadastros gerados.

O parâmetro de interesse é a média populacional $\mu = 21.1$. Através do plano de amostragem aleatória simples, duas amostras independentes, de tamanho 1000 (cadastro A) e 800 (cadastro B) foram selecionadas, obtendo-se

$$\bar{y}_A = 22.56, \quad \bar{y}_a = 21.35, \quad \bar{y}_{ab(A)} = 31.91$$

$$\bar{y}_B = 22.57, \quad \bar{y}_b = 24.04, \quad \bar{y}_{ab(B)} = 12.67$$

$$w_{ab} = 10000/169587 = 0.0589$$

Logo, basta substituir estas quantidades na forma do estimador de Hartley, bem como na forma do valor ótimo de p , para se obter um estimador para a média populacional.

2 Estratégia de Fuller & Burmeister(1972)

A estratégia de Hartley é utilizada nas situações em que N_{ab} é conhecido (cenário 2, tabela 2.1.). Quando isto não acontece, o seguinte estimador proposto por Fuller & Burmeister (1972) pode ser utilizado. De uma forma geral, utilizando estimadores do tipo Horvitz-Thompson, tem-se

$$\begin{aligned} \bar{y}_{FB} &= N^{-1} \{ \hat{t}_{ya} + \hat{t}_{yb} + \beta_1 \hat{t}_{yab}^A + (1 - \beta_1) \hat{t}_{yab}^B + \beta_2 (\hat{N}_{ab}^A - \hat{N}_{ab}^B) \} \\ &= \hat{w}_a \bar{y}_{HT(a)} + \hat{w}_b \bar{y}_{HT(b)} + \beta_1 \hat{w}_{ab(A)} \bar{y}_{HTab(A)} + \hat{w}_{ab(B)} (1 - \beta_1) \bar{y}_{HTab(B)} + \\ &+ N^{-1} \beta_2 (\hat{N}_{ab}^A - \hat{N}_{ab}^B), \end{aligned} \tag{2.9}$$

em que

- $\hat{w}_a = \frac{\hat{N}_a}{N}$, $\hat{w}_b = \frac{\hat{N}_b}{N}$, $\hat{w}_{ab(A)} = \frac{\hat{N}_{ab}^A}{N}$ e $\hat{w}_{ab(B)} = \frac{\hat{N}_{ab}^B}{N}$, em que, por exemplo, $\hat{N}_a = \sum_{k \in S_a} (1/\pi_k^A)$;
- \hat{t}_{ya} é o estimador de Horvitz-Thompson para o total populacional do domínio a ;
- \hat{t}_{yb} é o estimador de Horvitz-Thompson para o total populacional do domínio b ;
- \hat{t}_{yab}^A é o estimador de Horvitz-Thompson para o total populacional do domínio ab no cadastro A ;
- \hat{t}_{yab}^B é o estimador de Horvitz-Thompson para o total populacional do domínio ab no cadastro B ;
- \hat{N}_{ab}^A e \hat{N}_{ab}^B são estimadores para o tamanho populacional do domínio ab em cada cadastro, respectivamente.
- β_1 e β_2 são constantes de ponderação, escolhidas como os valores que minimizam a variância do estimador de Fuller & Burmeister, que será apresentada mais adiante.

O estimador de Fuller & Burmeister (1972) é função de estimadores da forma $\hat{Q} = \hat{w}_d \bar{y}_d$. Através do método de linearização de Taylor é possível obter um pseudo-estimador para que seja possível o cálculo da variância aproximada de \hat{Q} . Tem-se que

$$Q_L = w\mu + \sum_{i=1}^2 b_i(\hat{\theta}_i - \theta_i), \text{ onde:}$$

$$Q = w_d \mu_d, \quad \hat{\theta}_1 = \hat{w}_d \quad \hat{\theta}_2 = \hat{\mu}_d$$

$$b_1 = \left. \frac{\partial \hat{Q}}{\partial \hat{w}} \right|_{(\hat{w}, \hat{\mu})=(w, \mu)} = \mu \quad \text{e}$$

$$b_2 = \left. \frac{\partial \hat{Q}}{\partial \hat{\mu}} \right|_{(\hat{w}, \hat{\mu})=(w, \mu)} = w.$$

Assim,

$$\begin{aligned} \hat{Q}_{(L)} &= Q + \mu_d(\hat{w}_d - w_d) + w_d(\hat{\mu}_d - \mu_d) \\ &= Q + \mu_d \hat{w}_d - \mu w_d + w_d \hat{\mu}_d - w_d \mu_d \\ &= \mu_d \hat{w}_d + w_d \hat{\mu}_d - w_d \mu_d, \end{aligned}$$

onde $w_d = N_d/N$ e μ_d é a média populacional do domínio d . O resultado obtido implica um pseudo-estimador para \bar{y}_{FB} , dado por

$$\begin{aligned} \bar{y}_{FB(L)} = & (\mu_a \hat{w}_a + w_a \hat{\mu}_a - w_a \mu_a) + (\mu_b \hat{w}_b + w_b \hat{\mu}_b - w_b \mu_b) \\ & + \beta_1 (\mu_{ab(A)} \hat{w}_{ab(A)} + w_{ab(A)} \hat{\mu}_{ab(A)} - w_{ab(A)} \mu_{ab(A)}) \\ & + (1 - \beta_1) (\mu_{ab(B)} \hat{w}_{ab(B)} + w_{ab(B)} \hat{\mu}_{ab(B)} - w_{ab(B)} \mu_{ab(B)}) + N^{-1} \beta_2 (\hat{N}_{ab}^A - \hat{N}_{ab}^B), \end{aligned}$$

Considerando a expressão (2.7), a variância do estimador de Fuller & Burmeister é dada em função de

$$\mathbf{d} = \begin{bmatrix} w_a \\ w_{ab(A)} \beta_1 \\ w_b \\ w_{ab(B)} (1 - \beta_1) \\ N^{-1} \beta_2 \\ N^{-1} \beta_2 \end{bmatrix} \quad (2.10)$$

e

$$\Sigma_{FB} = \begin{pmatrix} \sigma_a^2 & \sigma_{a;ab(A)} & 0 & 0 & 0 & 0 \\ \sigma_{a;ab(A)} & \sigma_{ab(A)}^2 & 0 & 0 & 0 & 0 \\ 0 & 0 & \sigma_b^2 & \sigma_{b;ab(B)} & 0 & 0 \\ 0 & 0 & \sigma_{b;ab(B)} & \sigma_{ab(B)}^2 & 0 & 0 \\ 0 & 0 & 0 & 0 & \sigma_{N_{A(ab)}}^2 & 0 \\ 0 & 0 & 0 & 0 & 0 & \sigma_{N_{B(ab)}}^2 \end{pmatrix} \quad (2.11)$$

Na matriz Σ_{FB} , por exemplo, $\sigma_{A(a)}^2 = \text{Var}(\bar{y}_{HT(a)})$, e os demais termos são definidos como anteriormente. Sob um plano de amostragem aleatória simples, por exemplo, o estimador de Fuller & Burmeister (1972) é dado por

$$\bar{y}_{FB} = \hat{w}_a \bar{y}_a + \hat{w}_b \bar{y}_b + \hat{w}_{ab,s} \bar{y}_{ab}, \quad (2.12)$$

onde

$$\bar{y}_{ab} = [n_{ab(A)} (1 - f_B) + n_{ab(B)} (1 - f_A)]^{-1} [(1 - f_B) n_{ab(B)} \bar{y}_{ab(A)} + (1 - f_A) \bar{y}_{ab(B)}],$$

$$w_{ab,s} = \hat{N}_{ab,s}/N, f_A = n_A/N_A \text{ e } f_B = n_B/N_B$$

Na expressão, $\hat{N}_{ab,s}$ é definida como a menor raiz da seguinte equação quadrática em x

$$\begin{aligned} & (n_A g_B + n_B g_A) x^2 - (n_A N_B g_B + n_B N_A g_A + n_{ab(A)} N_A g_B + n_{ab(B)} N_B g_A) x \\ & + (n_{ab(A)} g_B + n_{ab(B)} g_A) N_A N_B = 0. \end{aligned} \tag{2.13}$$

Fuller & Burmeister provaram que a variância de $\hat{N}_{ab,s}$ é dada por

$$\text{Var}(\hat{N}_{ab,s}) = \frac{N_a N_b N_{ab} g_A g_B}{n_A n_b g_B + n_B n_a g_A} + \mathbf{O}(1) \text{ e } E(\hat{N}_{ab,s}) = N_{ab} + \mathbf{O}\left(\frac{1}{n}\right),$$

onde $\mathbf{O}(\cdot)$ indica ordem de magnitude, como usado em análise real¹.

3 Estratégia de Bankier, Lepkowski & Groves(1986)

A estratégia de estimação sob a abordagem de cadastro duplo adotada por Bankier, Lepkowski & Groves (1986), descrita aqui simplesmente por estratégia BLG, atribui pesos, às unidades amostrais de A e B . A ideia dos autores foi adotar uma abordagem que envolvesse apenas quantidades amostrais referentes aos cadastros A e B , como mostra a figura (2.2) a seguir.

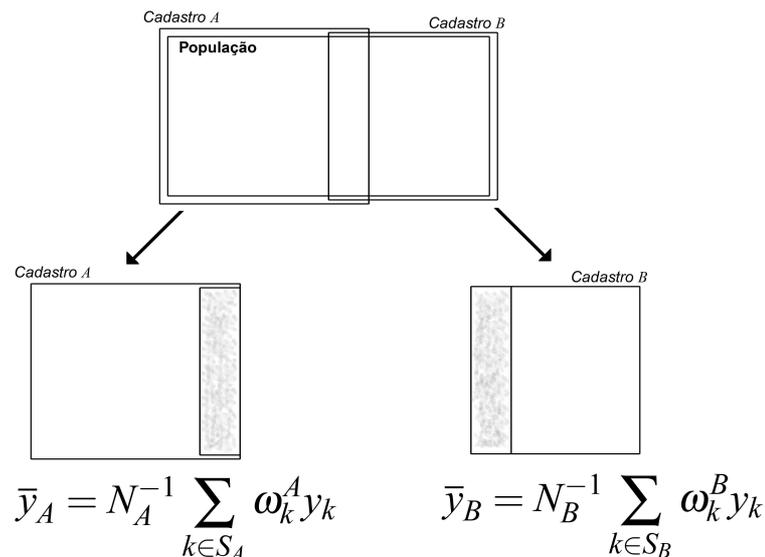


Figura 2.2: Quantidades amostrais geradas pela estratégia BLG

¹Para um maior aprofundamento, o leitor pode consultar Fuller (1996). Uma breve revisão sobre o tema é apresentada no apêndice A.

Denotando por ω_k^C o peso associado a um cadastro C , tem-se que

$$\omega_k^C = \begin{cases} \frac{1}{\pi_k^A}, & k \in a \\ \frac{1}{\pi_k^B}, & k \in b \\ \frac{1}{\pi_k^A + \pi_k^B}, & k \in ab \end{cases}$$

Quando há o interesse em estimar por exemplo uma média populacional denotada por μ , tem-se que o estimador BLG é dado por

$$\begin{aligned} \bar{y}_{BLG} &= N^{-1} \left(\sum_{k \in S_A} \omega_k^A y_k + \sum_{k \in S_B} \omega_k^B y_k \right) = w_A \frac{\sum_{k \in S_A} \omega_k^A y_k}{N_A} + w_B \frac{\sum_{k \in S_B} \omega_k^B y_k}{N_B} \\ &= N^{-1} \left(\sum_{k \in a} \omega_k y_k + \sum_{k \in b} \omega_k y_k + \sum_{k \in ab} \omega_k y_k \right) = w_a \bar{y}_a + w_b \bar{y}_b + w_{ab} \bar{y}_{ab} \end{aligned} \quad (2.14)$$

com $w_A = N_A/N$, $w_B = N_B/N$, $w_a = N_a/N$, $w_b = N_b/N$ e $w_{ab} = N_{ab}/N$. A variância de \bar{y}_{BLG} é obtida a partir da expressão (2.7), onde

$$\mathbf{d} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \text{e} \quad \Sigma_{BLG} = \begin{pmatrix} \sigma_{yA}^2 & 0 \\ 0 & \sigma_{yB}^2 \end{pmatrix}$$

Na matriz Σ_{BLG} , $\sigma_{yA}^2 = \text{Var} \left(\sum_{k \in \mathcal{A}} \omega_k^A y_k \right)$ e $\sigma_{yB}^2 = \text{Var} \left(\sum_{k \in \mathcal{B}} \omega_k^B y_k \right)$. É possível obter um estimador não-viesado para a variância do estimador de Bankier, substituindo em Σ_{BLG} as quantidades populacionais pelos seus respectivos estimadores de Horvitz-Thompson.

A partir de (2.14) Bankier, Lepkowski & Groves (1986) propuseram ainda um estimador iterativo chamado *raking ratio estimator*, dado por

$$\bar{y}_{BLG(r)} = N^{-1} \left(\sum_{k \in S_A} \omega_{ka}^{(r)} y_k + \sum_{k \in S_B} \omega_{kb}^{(r)} y_k + \sum_{k \in ab} \omega_{kab}^{(r)} y_k \right). \quad (2.15)$$

A ideia é que se façam r ajustes à constante ω_k , onde

$$\begin{aligned} \omega_{ka}^{(r)} &= \frac{N_a}{\hat{N}^{A(r-1)}} \omega_{ka}^{(r-1)}, & \omega_{kb}^{(r)} &= \omega_{kb}^{(r-1)}, & \omega_{kab}^{(r)} &= \frac{N_a}{\hat{N}^{A(r-1)}} \omega_{kab}^{(r-1)} \quad \text{para } r = 1, 3, 5, \dots \\ \omega_{ka}^{(r)} &= \omega_{ka}^{(r-1)}, & \omega_{kb}^{(r)} &= \frac{N_b}{\hat{N}_v^{B(r-1)}} \omega_{kb}^{(r-1)}, & \omega_{kab}^{(r)} &= \frac{N_b}{\hat{N}_v^{B(r-1)}} \omega_{kab}^{(r-1)} \quad \text{para } r = 2, 4, 6, \dots \\ \omega_{ka}^{(0)} &= \frac{1}{\pi_k^A}, & \omega_{kb}^{(0)} &= \frac{1}{\pi_k^B}, & \omega_{kab}^{(0)} &= \frac{1}{(\pi_k^A + \pi_k^B)}, \\ \hat{N}_{A(r)} &= \omega_{ka}^{(r)} N_a + \omega_{kab}^{(r)} N_{ab} & \text{e} & & \hat{N}_{B(r)} &= \omega_{kb}^{(r)} N_b + \omega_{kab}^{(r)} n_{ab}. \end{aligned}$$

O estimador $\bar{y}_{BLG(r)}$ é obtido fazendo $\bar{y}_{BLG(0)} = \bar{y}_{BLG}$ e aplicando r ajustes com respeito aos domínios, alternadamente. Como para um determinado valor de r as expressões para a variância do estimador $\bar{y}_{BLG(r)}$ são bastante complexas (Bankier, Lepkowski & Groves (1986) ; Brackstone & Rao (1979)), Skinner (1991) provou que quando o número de ajustes é suficientemente grande, ou seja, quando $r \rightarrow \infty$, tem-se que

$$\bar{y}_{BLG(\infty)} = N^{-1} \{ (N_A - \tilde{N}_{ab}) \bar{y}_a + \tilde{N}_{ab} \bar{y}_{ab} + (N_B - \tilde{N}_{ab}) \bar{y}_b \}, \quad (2.16)$$

Logo, a variância aproximada do estimador raking é dada por

$$\begin{aligned} \text{AVar}(\bar{y}_{BLG(\infty)}) &= \left(\frac{1}{\pi_k^A} \right)^{-1} N_a \sigma_a^2 + \left(\frac{1}{\pi_k^B} \right)^{-1} N_b \sigma_b^2 + (\pi_k^A + \pi_k^B)^{-1} N_{ab} \sigma_{ab}^2 \\ &+ (\mu_a + \mu_b - \mu_{ab})^2 \frac{N_{ab} N_a N_b}{n_A N_b + n_B N_a} (1 + \lambda^2), \end{aligned} \quad (2.17)$$

onde μ_a , μ_b e μ_{ab} são as médias populacionais de cada domínio e σ_a^2 , σ_b^2 e σ_{ab}^2 são as variâncias populacionais de cada domínio, e

$$\lambda^2 = \frac{N_a N_b (N_{ab})^2 \left[(\pi_k^A)^2 N_A - (\pi_k^B)^2 N_B \right]^2}{N_A N_B \pi_k^A \pi_k^B (\pi_k^A + \pi_k^B)^2 ((N_{ab})^2 - N_A N_B)^2}.$$

4 Estratégia de Máxima Pseudo-Verossimilhança

Skinner & Rao (1996) apresentaram o estimador de máxima pseudo-verossimilhança (*pseudo maximum likelihood* - PML), considerando a situação descrita pelo cenário 3. Esta estratégia é considerada pois o estimador de Fuller & Burmeister não pode ser aplicado diretamente em planos amostrais complexos, por não ser consistente para o parâmetro de interesse. Dessa forma, a estratégia de máxima pseudo-verossimilhança considera uma modificação no estimador de Fuller & Burmeister, de modo a validar a propriedade de consistência. De um modo geral, utilizando novamente estimadores do tipo Horvitz-Thompson para obter estimadores nos domínios, tem-se que:

$$\begin{aligned}\bar{y}_{PML} &= N^{-1} \left[(N_A - \hat{N}_{ab,PML}) \bar{y}_{HT(a)} + \hat{N}_{ab,PML} \bar{y}_{ab(*)} + (N_B - \hat{N}_{ab,PML}) \bar{y}_{HT(b)} \right] \\ &= w_{a(*)} \bar{y}_{HT(a)} + w_{ab(*)} \bar{y}_{ab(*)} + w_{b(*)} \bar{y}_{HT(b)}\end{aligned}\quad (2.18)$$

onde

$$\begin{aligned}\bar{y}_{ab(*)} &= \frac{[\pi^A \hat{N}_{ab}^A \bar{y}_{ab(A)} + \pi^B \hat{N}_{ab}^B \bar{y}_{ab(B)}]}{[\pi^A \hat{N}_{ab}^A + \pi^B \hat{N}_{ab}^B]} \\ w_{a(*)} &= (N_A - \hat{N}_{ab,PML}) / N & w_{b(*)} &= (N_B - \hat{N}_{ab,PML}) / N \\ w_{ab(*)} &= \hat{N}_{ab,PML} / N \\ \hat{N}_{ab}^A &= \sum_{k \in ab} \frac{1}{\pi^A} & \hat{N}_{ab}^B &= \sum_{k \in ab} \frac{1}{\pi^B}.\end{aligned}$$

Com $\pi^A = \pi_k^A$ e $\pi^B = \pi_k^B$, para todo $k \in \mathcal{A}$ e $k \in \mathcal{B}$ respectivamente. $\hat{N}_{ab,PML}$ é a menor raiz da equação quadrática

$$ax^2 + bx + c = 0, \quad (2.19)$$

com

$$\begin{aligned}a &= n_A + n_B, \\ b &= n_A n_B + n_B n_A + n_A \hat{N}_{ab}^A + n_B \hat{N}_{ab}^B;\end{aligned}$$

$$c = n_A \hat{N}_{ab}^A N_B + n_B \hat{N}_{ab}^B N_A.$$

Como $b^2 - 4ac = (n_A N_B + n_B N_A + n_A \hat{N}_{ab}^A + n_B \hat{N}_{ab}^B)^2 + 4n_A n_B (N_A - \hat{N}_{ab}^A) (N_B - \hat{N}_{ab}^B)$, (2.19) não possui raízes complexas, pois $\hat{N}_{ab}^A \leq N_A$ e $\hat{N}_{ab}^B \leq N_B$. Além disso, Skinner & Rao (1996) mostraram que

$$\text{AVar}(\hat{N}_{ab,PML}) = \left(\frac{N^2}{n}\right) \left[\alpha^2 \sigma_{ab(A)}^2 + (1 - \alpha)^2 \sigma_{ab(B)}^2 \right],$$

onde $N = N_A + N_B - N_{ab}$, $N = n_A + n_B - n_{ab}$, $n_{ab} = n_{ab(A)} + n_{ab(B)}$ e α é uma constante tal que $0 \leq \alpha \leq 1$, cujo valor que minimiza $\text{AVar}(\hat{N}_{ab,PML})$ é dado por

$$\alpha = \frac{\sigma_{ab(B)}^2}{\sigma_{ab(A)}^2 + \sigma_{ab(B)}^2}.$$

Para cálculo da variância do estimador PML utilizando o método de linearização de Taylor, são necessários desenvolvimentos apresentados anteriormente para o estimador de Fuller & Burmeister. Além destes desenvolvimentos é necessário utilizar o método de linearização de Taylor em $\bar{y}_{ab(*)}$. Logo,

$$\begin{aligned} \bar{y}_{ab*} &\approx \frac{\pi^A N_{ab} \mu_{ab} + \pi^B N_{ab} \mu_{ab}}{\pi^A \hat{N}_{ab}^A + \pi^B \hat{N}_{ab}^B} + \frac{1}{\pi^A N_{ab} + \pi^B N_{ab}} \left\{ \pi^A \hat{N}_{ab}^A \bar{y}_{ab(A)} + \pi^B \hat{N}_{ab}^B \bar{y}_{ab(B)} \right. \\ &\quad \left. - \left(\frac{\pi^A N_{ab} \mu_{ab} + \pi^B N_{ab} \mu_{ab}}{\pi^A N_{ab} + \pi^B N_{ab}} \right) (\pi^A \hat{N}_{ab}^A + \pi^B \hat{N}_{ab}^B) \right\} \\ &= \mu_{ab} + \frac{(N_{ab})^{-1}}{\pi^A + \pi^B} \left\{ \pi^A \hat{N}_{ab}^A \bar{y}_{ab(A)} + \pi^B \hat{N}_{ab}^B \bar{y}_{ab(B)} - \mu_{ab} (\pi^A \hat{N}_{ab}^A + \pi^B \hat{N}_{ab}^B) \right\} \\ &\approx \mu_{ab} + \frac{(N_{ab})^{-1}}{\pi^A + \pi^B} \left\{ \pi^A (\hat{N}_{ab}^A \mu_{ab} + N_{ab} \bar{y}_{ab(A)} - N_{ab} \mu_{ab}) + \pi^B (\hat{N}_{ab}^B \mu_{ab} + N_{ab} \bar{y}_{ab(B)} - N_{ab} \mu_{ab}) \right. \\ &\quad \left. - \mu_{ab} \pi^A \hat{N}_{ab}^A - \mu_{ab} \hat{N}_{ab}^B \pi^B \right\} \\ &= \mu_{ab} + \frac{(N_{ab})^{-1}}{\pi^A + \pi^B} \left\{ N_{ab} \pi^A \bar{y}_{ab(A)} - \pi^A N_{ab} \mu_{ab} + N_{ab} \pi^B \bar{y}_{ab(B)} - \pi^B N_{ab} \mu_{ab} \right\} \\ &= \mu_{ab} + \left(\frac{\pi^A}{\pi^A + \pi^B} \right) \bar{y}_{ab(A)} + \left(\frac{\pi^B}{\pi^A + \pi^B} \right) \bar{y}_{ab(B)} - \left(\frac{\pi^A + \pi^B}{\pi^A + \pi^B} \right) \mu_{ab} \\ &= \left(\frac{\pi^A}{\pi^A + \pi^B} \right) \bar{y}_{ab(A)} + \left(\frac{\pi^B}{\pi^A + \pi^B} \right) \bar{y}_{ab(B)} \end{aligned}$$

Logo, um pseudo-estimador para o estimador PML é dado por

$$\begin{aligned} \bar{y}_{PML(L)} \approx & w_a \bar{y}_a + w_{a*} \mu_a - w_a \mu_a + w_b \bar{y}_b + w_{b*} \mu_b - w_b \mu_b + \\ & + w_{ab} \left\{ \left(\frac{\pi^A}{\pi^A + \pi^B} \right) \bar{y}_{ab(A)} + \left(\frac{\pi^B}{\pi^A + \pi^B} \right) \bar{y}_{ab(B)} \right\} \\ & + w_{ab*} \mu_{ab} - w_{ab} \mu_{ab}. \end{aligned} \quad (2.20)$$

Dessa forma, baseando-se na expressão (2.7), a variância aproximada do estimador PML é dada em função das quantidades

$$\mathbf{d}^T = [w_a \ w_b \ w_{ab} \ w_{ab} \ \mu_a/N \ \mu_b/N \ \mu_{ab}/N] \text{ e}$$

$$\Sigma_{PML} = \begin{pmatrix} \sigma_{yA}^2 & \sigma_{yA;yab(A)} & 0 & 0 & 0 & 0 & 0 & 0 \\ \sigma_{yA;yab(A)} & \sigma_{yab(A)}^2 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \sigma_{yB}^2 & \sigma_{yB;yab(B)} & 0 & 0 & 0 & 0 \\ 0 & 0 & \sigma_{yB;yab(B)} & \sigma_{yab(B)}^2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \sigma_{N_a,PML}^2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & \sigma_{N_b,PML}^2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \sigma_{N_{ab},PML}^2 & 0 \end{pmatrix} \quad (2.21)$$

Foram revisadas quatro estratégias de estimação sob a abordagem de cadastro duplo, sob os cenários 2 e 3. A partir da forma dos estimadores e de suas respectivas variâncias, é possível tirar as seguintes conclusões:

1. A partir da forma da variância do estimador \bar{y}_{FB} , quando $\beta_1 = \gamma$ e $\beta_2 = (\alpha - \gamma) \mu_{ab} + (1 - \alpha) \mu^a + \alpha \mu^b$, onde $\alpha = N_a N_{ab} / (N_a N_b + N_b N_a)$ e $\gamma = N_a N^B / (N_a N_b + N_b N_a)$, \mathbf{d}^T será equivalente à forma apresentada para o estimador PML;
2. Sob a condição comentada no item 1, tem-se que $\text{Var}(\bar{y}_{PML}) \geq \text{Var}(\bar{y}_{FB})$;
3. Considerando as condições 1 e 2 e o estimador \bar{y}_{FB} , quando $\beta_2 = 0$, tem-se que $\bar{y}_{FB} \approx \bar{y}_H$, e portanto $\text{Var}(\bar{y}_{FB}) \geq \text{Var}(\bar{y}_H)$. Dessa forma, tem-se que $\text{Var}(\bar{y}_{PML}) \geq \text{Var}(\bar{y}_H)$;
4. Caso os valores de β_1 e β_2 sejam escolhidos de outra forma, o estimador sob a estratégia PML não será necessariamente mais eficiente que o estimador sob a estratégia de Hartley.

As estratégias revisadas não utilizam informações de possíveis variáveis auxiliares no processo de estimação. A estratégia apresentada a seguir representa uma forma de fazer uso de informações auxiliares.

5 Estratégia de Singh & Wu (2003)

No contexto de estimadores do tipo regressão, Singh & Wu (2003) apresentaram uma estratégia de estimação sob abordagem de cadastro duplo baseada em um estimador intuitivo, dado por

$$\bar{y} = w_a \bar{y}_{Ra} + w_b \bar{y}_{Rb} + \frac{1}{2} \left(w_{ab} \bar{y}_{Rab(A)} + w_{ab} \bar{y}_{Rab(B)} \right), \quad (2.22)$$

onde \bar{y}_{Ra} , \bar{y}_{Rb} , $\bar{y}_{Rab(A)}$ e $\bar{y}_{Rab(B)}$ são os estimadores regressão para as médias dos domínios a , b e ab e que serão apresentados no próximo capítulo, $w_a = N_a/N$, $w_b = N_b/N$, $w_{ab} = N_{ab}/N$ e $N = N_A + N_B - N_{ab}$. Este estimador foi proposto para situações em que pode ocorrer $n_{ab(A)} = n_{ab(B)} = 0$, não sendo necessário o cálculo de valores ótimos para minimização da sua variância. Esta estratégia não envolve obtenção de estimativas por mínimos quadrados, pois calcula os coeficientes de regressão através de métodos de calibração, e portanto não será explorada no desenvolvimento desta tese. A partir deste estimador, foi proposto o estimador regressão de calibração (RC), dado por

$$\bar{y}_{RC} = w_a \bar{y}_{RCa} + w_b \bar{y}_{RCb} + \zeta_A \left(w_{ab} \bar{y}_{RCab(A)} \right) + \zeta_B \left(w_{ab} \bar{y}_{RCab(B)} \right), \quad (2.23)$$

onde ζ_A e ζ_B são coeficientes de ajuste.

CAPÍTULO 3

Estimadores do tipo regressão em uma abordagem de cadastro duplo

1 Introdução

O estimador regressão para o caso em que q variáveis auxiliares estão disponíveis é visto em várias referências. Para um maior aprofundamento, o leitor pode consultar Särndal, Swensson & Wretman (2003). Neste tipo de situação, admite-se que o vetor de informações auxiliares, associado ao elemento k da população, é conhecido e denotado

$$\mathbf{x}_k = (x_{1k}, x_{2k}, \dots, x_{qk}).$$

Considere, por exemplo, a estimação de um total populacional

$$t_y = \sum_{k \in U} y_k. \quad (3.1)$$

Defina as seguintes quantidades:

$$\check{y}_k = \frac{y_k}{\pi_k}, \quad \hat{t}_{y\pi} = \sum_{k \in S} \check{y}_k \quad \hat{\mathbf{T}} = \sum_{k \in S} \frac{\mathbf{x}_k \mathbf{x}_k^T}{\sigma_k^2 \pi_k} \quad \hat{\mathbf{t}} = \sum_{k \in S} \frac{\mathbf{x}_k \check{y}_k}{\sigma_k^2},$$

$$\check{x}_{jk} = \frac{x_{jk}}{\pi_{jk}} \quad (j = 1, 2, \dots, q), \quad \hat{\mathbf{t}}_{x\pi} = (\hat{t}_{x_1\pi}, \hat{t}_{x_2\pi}, \dots, \hat{t}_{x_q\pi}), \quad \text{onde } \hat{t}_{x_j\pi} = \sum_{k \in S} \check{x}_{jk}$$

O estimador regressão generalizado para t_y , utilizando q variáveis auxiliares pode ser expresso da seguinte forma:

$$\begin{aligned}
 \hat{t}_{yR} &= \hat{t}_{y\pi} + \sum_{j=1}^q \hat{\beta}_{Vj} (t_{x_j} - \hat{t}_{x_j\pi}) = \hat{t}_{y\pi} + (\mathbf{t}_x - \hat{\mathbf{t}}_x)^T \hat{\beta} \\
 &= \sum_{k \in S} \check{y}_k + (\mathbf{t}_x - \hat{\mathbf{t}}_{x\pi})^T (\hat{\mathbf{T}})^{-1} \sum_{k \in S} \sigma_k^{-2} \mathbf{x}_k \check{y}_k \\
 &= \sum_{k \in S} \left[1 + (\mathbf{t}_x - \hat{\mathbf{t}}_{x\pi})^T (\hat{\mathbf{T}})^{-1} \sigma_k^{-2} \mathbf{x}_k \right] \check{y}_k = \sum_{k \in S} g_{ks} \check{y}_k \\
 &= f(\hat{t}_{y\pi}, \hat{\mathbf{t}}_{x\pi}, \hat{\mathbf{T}}, \hat{\mathbf{t}}), \tag{3.2}
 \end{aligned}$$

onde $f(\cdot)$ é uma função não-linear dos estimadores $\hat{t}_{y\pi}$, $\hat{\mathbf{t}}_{x\pi}$, $\hat{\mathbf{T}}$ e $\hat{\mathbf{t}}$. A variância aproximada deste estimador (Särndal, 2003) é dada por

$$\text{AVar}(\hat{t}_{yR}) = \sum_{k \in U} \sum_{l \in U} \Delta_k \frac{E_k E_l}{\pi_k \pi_l}, \quad \text{onde } E_k = y_k - \mathbf{x}_k^T \hat{\beta} \tag{3.3}$$

A seguir, serão apresentadas estratégias de estimação para a abordagem de cadastro duplo que fazem uso de estimadores do tipo regressão, e que foram inspiradas nas estratégias apresentadas no capítulo 2. Tais estratégias constituem a contribuição original desta tese.

2 Estimador regressão sob a estratégia de Hartley

Considere a situação de uso de uma abordagem de cadastro duplo sob o cenário 2, e que as seguintes condições são satisfeitas:

- O vetor de q informações auxiliares, denotado por $\mathbf{x}_{kA} = (x_{1k}, x_{2k}, \dots, x_{qk})$ é conhecido para todo $k \in \mathcal{A}$;
- O vetor de m informações auxiliares, denotado por $\mathbf{x}_{kB} = (x_{1k}, x_{2k}, \dots, x_{mk})$ é conhecido para todo $k \in \mathcal{B}$. É importante notar que as informações de variáveis auxiliares disponíveis em um cadastro não precisam necessariamente coincidir com as informações auxiliares disponíveis no outro.

Nos capítulos anteriores foi apresentada a seguinte partição da média populacional

$$t_y = \sum_{k \in \mathcal{A}} y_k + \sum_{k \in \mathcal{B}} y_k - \sum_{k \in \mathcal{A} \cap \mathcal{B}} y_k.$$

Dessa forma, o seguinte estimador regressão para a média populacional é proposto e denominado estimador regressão sob a estratégia de Hartley.

Estimador regressão sob a estratégia de Hartley

$$\begin{aligned}
 \bar{y}_{RH} &= N^{-1} \left\{ \hat{t}_{yR}^A + \hat{t}_{yR}^B - p \hat{t}_{yR}^{ab(A)} - (1-p) \hat{t}_{yR}^{ab(B)} \right\} \\
 &= N^{-1} \left\{ N_A \bar{y}_{RA} + N_B \bar{y}_{RB} - N_{ab} p \bar{y}_{Rab(A)} - (1-p) N_{ab} \bar{y}_{Rab(B)} \right\} \\
 &= w_A \bar{y}_{RA} + w_B \bar{y}_{RB} - w_{ab} \left(p \bar{y}_{Rab(A)} + (1-p) \bar{y}_{Rab(B)} \right) \\
 &= w_A \frac{f(\hat{t}_{y\pi}^A, \hat{\mathbf{t}}_{x\pi}^A, \hat{\mathbf{T}}^A, \hat{\mathbf{t}}^A)}{N_A} + w_B \frac{f(\hat{t}_{y\pi}^B, \hat{\mathbf{t}}_{x\pi}^B, \hat{\mathbf{T}}^B, \hat{\mathbf{t}}^B)}{N_B} \\
 &\quad - w_{ab} \left\{ p \frac{f(\hat{t}_{y\pi}^{ab(A)}, \hat{\mathbf{t}}_{x\pi}^{ab(A)}, \hat{\mathbf{T}}^{ab(A)}, \hat{\mathbf{t}}^{ab(A)})}{N_{ab}} + (1-p) \frac{f(\hat{t}_{y\pi}^{ab(B)}, \hat{\mathbf{t}}_{x\pi}^{ab(B)}, \hat{\mathbf{T}}^{ab(B)}, \hat{\mathbf{t}}^{ab(B)})}{N_{ab}} \right\}, \quad (3.4)
 \end{aligned}$$

onde \hat{t}_{yR}^A e \hat{t}_{yR}^B são os estimadores do tipo regressão para o total dos cadastros A e B e $\hat{t}_{yR}^{ab(A)}$ e $\hat{t}_{yR}^{ab(B)}$ são os estimadores do tipo regressão para o total do domínio ab . \bar{y}_{RA} e \bar{y}_{RB} são estimadores do tipo regressão para a média dos cadastros A e B e $\bar{y}_{Rab(A)}$ e $\bar{y}_{Rab(B)}$ são estimadores do tipo regressão para a média do domínio ab . $\hat{t}_{y\pi}^{ab(A)}$, $\hat{\mathbf{t}}_{x\pi}^{ab(A)}$, $\hat{\mathbf{T}}^{ab(A)}$ e $\hat{\mathbf{t}}^{ab(A)}$ são definidas de forma análoga às quantidades apresentadas em (3.2). Da mesma forma que a vista para a estratégia de Hartley sem considerar variáveis auxiliares, $0 \leq p \leq 1$.

Resultado 1 A variância aproximada do estimador \bar{y}_{RH} é função de variâncias aproximadas e covariâncias aproximadas para cada estimador regressão obtido nos cadastros e domínios:

$$\begin{aligned}
 AVar(\bar{y}_{RH}) &= (w_A)^2 AVar(\bar{y}_{RA}) + (w_B)^2 AVar(\bar{y}_{RB}) + (w_{ab})^2 p^2 AVar(\bar{y}_{Rab(A)}) \\
 &\quad + (w_{ab})^2 (1-p)^2 AVar(\bar{y}_{Rab(B)}) \\
 &\quad - 2p w_A w_{ab} ACov(\bar{y}_{RA}; \bar{y}_{Rab(A)}) - 2(1-p) w_B w_{ab} ACov(\bar{y}_{RB}; \bar{y}_{Rab(B)}) \quad (3.5)
 \end{aligned}$$

A partir da expressão (2.7), a variância aproximada de \bar{y}_{RH} é função de

$$\mathbf{d} = \begin{bmatrix} w_A \\ p w_{ab} \\ w_B \\ (1-p) w_{ab} \end{bmatrix} \quad \text{e} \quad \Sigma_{RH} = \begin{pmatrix} \sigma_A^2 & \sigma_{A;ab(A)} & 0 & 0 \\ \sigma_{A;ab(A)} & \sigma_{ab(A)}^2 & 0 & 0 \\ 0 & 0 & \sigma_B^2 & \sigma_{B;ab(B)} \\ 0 & 0 & \sigma_{B;ab(B)} & \sigma_{ab(B)}^2 \end{pmatrix}.$$

Os elementos de Σ_{RH} representam respectivamente as variâncias e covariâncias com base em estimadores regressão para cadastros e domínios respectivamente, como apresentado por Särndal, Swensson & Wretman (2003). Por exemplo,

$$\sigma_A^2 = \text{AVar}(\bar{y}_{RA}) = \sum_{k \in \mathcal{A}} \sum_{l \in \mathcal{A}} \Delta_{kl}^A \frac{E_k^A}{\pi_k^A} \frac{E_l^A}{\pi_l^A}, \quad \text{onde } E_k^A = y_k - \mathbf{x}_{kA}^T \tilde{\beta}^A$$

$$\sigma_B^2 = \text{AVar}(\bar{y}_{RB}) = \sum_{k \in \mathcal{B}} \sum_{l \in \mathcal{B}} \Delta_{kl}^B \frac{E_k^B}{\pi_k^B} \frac{E_l^B}{\pi_l^B}, \quad \text{onde } E_k^B = y_k - \mathbf{x}_{kB}^T \tilde{\beta}^B$$

3 Estimador regressão sob a estratégia de Fuller & Burmeister

Na situação em que os tamanhos populacionais dos domínios são desconhecidos (cenário 3), a estratégia de Fuller & Burmeister (1972) é considerada como alternativa à abordagem de Hartley (1962). Neste caso, o estimador regressão para a abordagem de cadastro duplo é denominado estimador regressão sob a estratégia de Fuller & Burmeister, e dado por

$$\bar{y}_{RFB} = w_a \bar{y}_{Ra} + w_b \bar{y}_{Rb} + \beta_1 w_{ab} \bar{y}_{Rab(A)} + w_{ab} (1 - \beta_1) \bar{y}_{Rab(B)} + N^{-1} \beta_2 (\hat{N}_{ab}^A - \hat{N}_{ab}^B), \quad (3.6)$$

onde β_1 e β_2 são constantes de ponderação escolhidas de modo a minimizar a variância aproximada do estimador \bar{y}_{RFB} .

Resultado 2 Admitindo a expressão (2.7), a expressão em termos matriciais da variância do estimador regressão sob a estratégia de Fuller & Burmeister é dada em função das quantidades

$$\mathbf{d} = \begin{bmatrix} w_a \\ w_a p_1 \\ w_b \\ w_b (1 - p_1) \\ N^{-1} p_2 \\ N^{-1} p_2 \end{bmatrix} \quad \text{e} \quad \Sigma_{RFB} = \begin{pmatrix} \sigma_{A(a)}^2 & \sigma_{A(a);A(ab)} & 0 & 0 & 0 & 0 \\ \sigma_{A(a);A(ab)} & \sigma_{A(ab)}^2 & 0 & 0 & 0 & 0 \\ 0 & 0 & \sigma_{B(b)}^2 & \sigma_{B(b);B(ab)} & 0 & 0 \\ 0 & 0 & \sigma_{B(b);B(ab)} & \sigma_{B(ab)}^2 & 0 & 0 \\ 0 & 0 & 0 & 0 & \sigma_{N_{A(ab)}}^2 & 0 \\ 0 & 0 & 0 & 0 & 0 & \sigma_{N_{B(ab)}}^2 \end{pmatrix}.$$

Os elementos de Σ_{RFB} representam as variâncias aproximadas dos estimadores do tipo regressão para quantidades nos cadastros e domínios, os quais estão presentes no estimador sob a estratégia. Da mesma forma que a vista para o estimador regressão sob a estratégia de Hartley (1962), um estimador para a variância deste estimador pode ser obtido substituindo-se as respectivas quantidades populacionais pelos seus respectivos estimadores de Horvitz-Thompson.

4 Estimador regressão sob a estratégia BLG

Esta estratégia de estimação no contexto de estimadores regressão é inspirada em Ban- kier, Lepkowski & Groves (1986), e também utiliza a mesma ponderação nos elementos dos cadastros apresentada anteriormente, e dada por

$$\omega_k = \begin{cases} \frac{1}{\pi_k^A}, & k \in a \\ \frac{1}{\pi_k^B}, & k \in b \\ \frac{1}{\pi_k^{A \cup B}} = \frac{1}{\pi_k^A + \pi_k^B}, & k \in ab \end{cases} \quad (3.7)$$

Dessa forma, é natural conceber um estimador regressão generalizado a partir desta estraté- gia. Considerando o problema de estimar, por exemplo, uma média populacional dada por

$$\mu = N^{-1}(t_{yA} + t_{yB}), \quad \text{onde } N = N_A + N_B,$$

tem-se que o estimador regressão BLG é dado por

$$\begin{aligned} \bar{y}_{RBLG} &= N^{-1} \left(\sum_{k \in S_A} g_{ks(A)}^* \omega_k y_k + \sum_{k \in S_B} g_{ks(B)}^* \omega_k y_k \right) \\ &= N^{-1} \left(\sum_{k \in S_a} g_{ksa} \omega_k y_k + \sum_{k \in S_b} g_{ksb} \omega_k y_k + \sum_{k \in S_{ab(A)}} g_{ksab(A)} \omega_k y_k + \sum_{k \in S_{ab(B)}} g_{ksab(B)} \omega_k y_k \right) \\ &= N^{-1} \left(\sum_{k \in S_a} g_{ksa} \frac{y_k}{\pi_k^A} + \sum_{k \in S_b} g_{ksb} \frac{y_k}{\pi_k^B} + \sum_{k \in S_{ab(A)}} g_{ksab(A)} \frac{y_k}{\pi_k^{A \cup B}} + \sum_{k \in S_{ab(B)}} g_{ksab(B)} \frac{y_k}{\pi_k^{A \cup B}} \right) \\ &= N^{-1} \left(\sum_{k \in S_a} g_{ksa} \frac{y_k}{\pi_k^A} + \sum_{k \in S_b} g_{ksb} \frac{y_k}{\pi_k^B} + \sum_{k \in S_{ab}} g_{ksab} \frac{y_k}{\pi_k^{A \cup B}} \right) \\ &= w_a \bar{y}_{Ra} + w_b \bar{y}_{Rb} + w_{ab} \bar{y}_{Rab}, \end{aligned} \quad (3.8)$$

onde $w_a = \frac{N_a}{N}$, $w_b = \frac{N_b}{N}$ e $w_{ab} = \frac{N_{ab}}{N}$. S_a , S_b e S_{ab} representam os elementos da amostra nos domínios a , b e ab . Dessa forma, a variância aproximada do estimador \bar{y}_{RBLG} é dada de modo mais simples que a apresentada para a estratégia de Hartley:

$$\begin{aligned} AVar(\bar{y}_{RBLG}) &= (w_a)^2 AVar(\bar{y}_{Ra}) + (w_b)^2 AVar(\bar{y}_{Rb}) + (w_{ab(A)})^2 AVar(\bar{y}_{Rab(A)}) \\ &+ (w_{ab(B)})^2 AVar(\bar{y}_{Rab(B)}) + 2w_a w_{ab(A)} ACov(\bar{y}_{Ra}; \bar{y}_{Rab(A)}) \\ &+ 2w_b w_{ab(B)} ACov(\bar{y}_{Rb}; \bar{y}_{Rab(B)}) \end{aligned} \quad (3.9)$$

Resultado 3 Admitindo a expressão (2.7), a variância do estimador BLG expressa em termos matriciais é dada em função de

$$\mathbf{d} = \begin{bmatrix} w_a \\ w_b \\ w_{ab(A)} \\ w_{ab(B)} \end{bmatrix} \quad \text{e} \quad \Sigma_{RBLG} = \begin{pmatrix} \sigma_a^2 & \sigma_{a;ab(A)} & 0 & 0 \\ \sigma_{a;ab(A)} & \sigma_{ab(A)}^2 & 0 & 0 \\ 0 & 0 & \sigma_b^2 & \sigma_{b;ab(B)} \\ 0 & 0 & \sigma_{b;ab(B)} & \sigma_{ab(B)}^2 \end{pmatrix},$$

Na expressão da variância, $\sigma_a^2 = AVar(\bar{y}_{Ra})$, $\sigma_{a;ab(A)} = ACov(\bar{y}_{Ra}; \bar{y}_{Rab(A)})$ são estimadores regressão para os domínios apresentadas na expressão (3.9). Os elementos de Σ_{RBLG} representam as variâncias e covariâncias aproximadas dos estimadores regressão para os domínios, respectivamente. É possível obter um estimador para esta matriz. Basta substituir as respectivas quantidades populacionais pelos respectivos estimadores de Horvitz-Thompson. Ainda, com base na estratégia, é possível também propor uma versão *raking* do estimador regressão. Neste caso, considerando novamente a mesma estrutura proposta por Skinner (1991), com

$$\begin{aligned} \omega_{ka}^{(r)} &= \frac{N^A}{\hat{N}^{A(r-1)}} \omega_{ka}^{(r-1)}, & \omega_{kb}^{(r)} &= \omega_{kb}^{(r-1)} & \omega_{kab}^{(r)} &= \frac{N^A}{\hat{N}^{A(r-1)}} \omega_{kab}^{(r-1)} \quad \text{para } r = 1, 3, 5, \dots \\ \omega_{ka}^{(r)} &= \omega_{ka}^{(r-1)}, & \omega_{kb}^{(r)} &= \frac{N^B}{\hat{N}^{B(r-1)}} \omega_{kb}^{(r-1)} & \omega_{kab}^{(r)} &= \frac{N^B}{\hat{N}^{B(r-1)}} \omega_{kab}^{(r-1)} \quad \text{para } r = 2, 4, 6, \dots \end{aligned}$$

$$\omega_{ka}^{(0)} = \frac{1}{\pi_k^A}, \quad \omega_{kb}^{(0)} = \frac{1}{\pi_k^B}, \quad \omega_{kab}^{(0)} = \frac{1}{(\pi_k^A + \pi_k^B)}$$

$$\hat{N}^{A(r)} = \omega_{ka}^{(r)} n^a + \omega_{kab}^{(r)} n_{ab}$$

$$\hat{N}^{B(r)} = \omega_{kb}^{(r)} n^b + \omega_{kab}^{(r)} n_{ab}$$

tem-se que

$$\bar{y}_{RBLG(r)} = N^{-1} \left(\sum_{k \in S_a} g_{ksa} \omega_{ka}^{(r)} y_k + \sum_{k \in S_b} g_{ksa} \omega_{kb}^{(r)} y_k + \sum_{k \in S_{ab}} g_{ksa} \omega_{kab}^{(r)} y_k \right), \quad (3.10)$$

No caso em que $r \rightarrow \infty$, tem-se que

$$\bar{y}_{RBLG(\infty)} = \hat{w}_{a*} \bar{y}_{Ra*} + \hat{w}_{b*} \bar{y}_{Rb*} + \hat{w}_{ab*} \bar{y}_{Rab*}, \quad (3.11)$$

onde

$$\begin{aligned} \hat{w}_{a*} &= N_A - \tilde{N}_{ab}, & \hat{w}_{b*} &= N_B - \tilde{N}_{ab}, & \hat{w}_{ab*} &= \tilde{N}_{ab} \\ \bar{y}_{Ra*} &= \frac{\sum_{k \in S_a} g_{ksa} y_k}{n_a}, & \bar{y}_{Rb*} &= \frac{\sum_{k \in S_b} g_{ksb} y_k}{n_b}, & \bar{y}_{Rab*} &= \frac{\sum_{k \in S_{ab}} g_{ksab} y_k}{n_{ab}} \quad e \end{aligned}$$

\tilde{N}_{ab} é a menor raiz da equação quadrática em x ,

$$(n_{ab(A)} + n_{ab(B)}) x^2 - \left\{ (n_{ab(A)} + n_{ab(B)}) (N_A + N_B) + \left[\left(\frac{1}{\pi_k^A} \right) + \left(\frac{1}{\pi_k^B} \right) \right] n_a n_b \right\} x + (n_{ab(A)} + n_{ab(B)}) N_A N_B = 0.$$

Resultado 4 A variância aproximada do estimador raking-regressão para a abordagem de cadastro duplo é dada por:

$$\begin{aligned} AVar(\bar{y}_{RBLG(\infty)}) &= (w_{a*})^2 AVar(\bar{y}_{Ra*}) + (w_{b*})^2 AVar(\bar{y}_{Rb*}) + (w_{ab*})^2 AVar(\bar{y}_{Rab*}) \\ &+ (\mu_a + \mu_b - \mu_{ab})^2 \frac{N_{ab} N_a N_b}{n_A N_b + n_B N_a} (1 + \lambda^2), \end{aligned} \quad (3.12)$$

com

$$\lambda^2 = \frac{N_a N_b N_{ab}^2 \left[(\pi_k^A)^2 N_A - (\pi_k^B)^2 N_B \right]^2}{N_A N_B \pi_k^A \pi_k^B (\pi_k^A + \pi_k^B)^2 (N_{ab}^2 - N_A N_B)^2}.$$

A partir de (2.7), esta variância é dada em função de

$$\mathbf{d} = \begin{bmatrix} w_a^* \\ w_b^* \\ w_{ab}^* \\ 1 \end{bmatrix} \quad \text{e} \quad \Sigma_{RBLG(\infty)} = \begin{bmatrix} \sigma_a^2 & 0 & 0 & 0 \\ 0 & \sigma_b^2 & 0 & 0 \\ 0 & 0 & \sigma_{ab}^2 & 0 \\ 0 & 0 & 0 & \sigma_{N_{ab}}^2 \end{bmatrix}$$

Na matriz $\Sigma_{RBLG(\infty)}$, os valores de σ_a^2 , σ_b^2 e σ_{ab}^2 são obtidos da expressão (3.12). Por exemplo, $\sigma_a^2 = \text{AVar}(\bar{y}_{Ra^*})$.

5 Estimador regressão sob a estratégia de Máxima Pseudo-Verossimilhança

Considere a situação de uso de uma abordagem de cadastro duplo sob o cenário 3, e que as informações sobre as variáveis auxiliares são conhecidas, conforme descrito anteriormente. Tem-se que o estimador regressão para a estratégia de máxima pseudo-verossimilhança é dado por

$$\begin{aligned} \bar{y}_{RPML} &= \frac{(N_A - \hat{N}_{ab,PML}) \bar{y}_{Ra} + \hat{N}_{ab,PML} \bar{y}_{ab^*} + (N_B - \hat{N}_{ab,PML}) \bar{y}_{Rb}}{N} \\ \bar{y}_{RPML} &= \hat{w}_{a^*} \bar{y}_{Ra} + \hat{w}_{ab^*} \bar{y}_{ab^*} + \hat{w}_{b^*} \bar{y}_{Rb} \end{aligned} \quad (3.13)$$

$$\bar{y}_{ab^*} = \frac{[\pi^A \hat{N}_{ab}^A \bar{y}_{Rab(A)} + \pi^B \hat{N}_{ab}^B \bar{y}_R^{ab(B)}]}{[\pi^A \hat{N}_{ab}^A + \pi^B \hat{N}_{ab}^B]}$$

As demais quantidades deste estimador são as mesmas apresentadas para a expressão (2.18).

Resultado 5 Admitindo a expressão (2.7), a variância aproximada do estimador regressão sob a estratégia PML é dada em função de

$$\mathbf{d}^T = \left[w_a \quad w_{ab} \left(\frac{\pi^A}{\pi^A + \pi^B} \right) \quad w_b \quad w_{ab} \left(\frac{\pi^B}{\pi^A + \pi^B} \right) \quad \frac{\mu_a}{N} \quad \frac{\mu_b}{N} \quad \frac{\mu_{ab}}{N} \right] \quad \text{e}$$

$$\Sigma_{RPML} = \begin{pmatrix} \sigma_a^2 & \sigma_{(a;ab(A))} & 0 & 0 & 0 & 0 & 0 \\ \sigma_{(a;ab(A))} & \sigma_{ab(A)}^2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \sigma_b^2 & \sigma_{(b;ab(B))} & 0 & 0 & 0 \\ 0 & 0 & \sigma_{(b;ab(B))} & \sigma_{ab(B)}^2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \sigma_{\hat{N}_{PML}^a}^2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & \sigma_{\hat{N}_{PML}^b}^2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \sigma_{\hat{N}_{PML}^{ab}}^2 \end{pmatrix}$$

Prova: Através de procedimentos já descritos anteriormente, tem-se que um pseudo-estimador para o estimador PML é dado por

$$\begin{aligned} \bar{y}_{RPML} &\approx w_a \bar{y}_{Ra} + \hat{w}_{a*} \mu_a - w_a \mu_a + w_b \bar{y}_{Rb} + \hat{w}_{b*} \mu_b - w_b \mu_b \\ &+ w_{ab} \left\{ \left(\frac{\pi^A}{\pi^A + \pi^B} \right) \bar{y}_{Rab(A)} + \left(\frac{\pi^B}{\pi^A + \pi^B} \right) \bar{y}_{Rab(B)} \right\} \\ &+ \hat{w}_{ab*} \mu_{ab} - w_{ab} \mu_{ab} \end{aligned} \quad (3.14)$$

Além dos procedimentos anteriores, é preciso linearizar o estimador \bar{y}_{ab*} . Tem-se que

$$\begin{aligned} \bar{y}_{ab*} &\approx \frac{\pi^A N_{ab} \mu^{ab} + \pi^B N_{(B)} \mu_{ab}}{\pi^A \hat{N}_{ab(A)} + \pi^B \hat{N}_{ab(B)}} + \frac{1}{\pi^A N_{ab} + \pi^B N_{ab}} \left\{ \pi^A \hat{N}_{ab}^A \bar{y}_{Rab(A)} + \pi^B \hat{N}_{ab}^B \bar{y}_{Rab(B)} \right. \\ &- \left. \left(\frac{\pi^A N_{ab} \mu_{ab} + \pi^B N_{ab} \mu_{ab}}{\pi^A N_{ab} + \pi^B N_{ab}} \right) (\pi^A \hat{N}_{ab(A)} + \pi^B \hat{N}_{ab(B)}) \right\} \\ &= \mu_{ab} + \frac{(N_{ab})^{-1}}{\pi^A + \pi^B} \left\{ \pi^A \hat{N}_{ab(A)} \bar{y}_{Rab(A)} + \pi^B \hat{N}_{ab(B)} \bar{y}_{Rab(B)} - \mu_{ab} (\pi^A \hat{N}_{ab(A)} + \pi^B \hat{N}_{ab(B)}) \right\} \\ &\approx \mu_{ab} + \frac{(N_{ab})^{-1}}{\pi^A + \pi^B} \left\{ \pi^A (\hat{N}_{ab(A)} \mu^{ab} + N_{ab} \bar{y}_{Rab(A)} - N_{ab} \mu_{ab}) + \pi^B (\hat{N}_{ab(B)} \mu_{ab} + N_{ab} \bar{y}_{Rab(B)} - N_{ab} \mu_{ab}) \right. \\ &- \left. \mu_{ab} \pi^A \hat{N}_{ab(A)} - \mu_{ab} \hat{N}_{ab(B)} \pi^B \right\} \\ &= \mu_{ab} + \frac{(N_{ab})^{-1}}{\pi^A + \pi^B} \left\{ N_{ab} \pi^A \bar{y}_{Rab(A)} - \pi^A N_{ab} \mu_{ab} + N_{ab} \pi^B \bar{y}_{Rab(B)} - \pi^B N_{ab} \mu_{ab} \right\} \\ &= \mu_{ab} + \left(\frac{\pi^A}{\pi^A + \pi^B} \right) \bar{y}_{Rab(A)} + \left(\frac{\pi^B}{\pi^A + \pi^B} \right) \bar{y}_{Rab(B)} - \left(\frac{\pi^A + \pi^B}{\pi^A + \pi^B} \right) \mu_{ab} \end{aligned}$$

$$= \left(\frac{\pi^A}{\pi^A + \pi^B} \right) \bar{y}_{Rab(A)} + \left(\frac{\pi^B}{\pi^A + \pi^B} \right) \bar{y}_{Rab(B)}.$$

Logo, obtendo a variância no estimador \bar{y}_{PML_L} , temos o resultado 5. ■

Sobre os estimadores propostos originalmente nesta tese é possível observar o seguinte:

1. Na forma da variância do estimador regressão sob a estratégia PML, se $p = \pi^A / (\pi^A + \pi^B)$ e $1 - p = \pi^B / (\pi^A + \pi^B)$, a variância do estimador regressão sob a estratégia de Hartley tende a se aproximar da variância do estimador regressão sob a estratégia PML, desde que $\sigma_{\hat{N}_{a,PML}}^2$, $\sigma_{\hat{N}_{b,PML}}^2$ e $\sigma_{\hat{N}_{ab,PML}}^2$ sejam negligíveis. Se esta condição não ocorrer, e $w_{a*} \approx w_a$, $w_{b*} \approx w_b$ e $w_{ab*} \approx w_{ab}$, a variância do estimador regressão sob a estratégia PML tende a se aproximar da variância do estimador regressão sob a estratégia BLG;
2. As variâncias dos estimadores regressão generalizados sob as estratégias PML e Fuller & Burmeister terão o mesmo vetor \mathbf{d}^T quando $\beta_1 = \theta$ e $\beta_2 = (\phi - \theta)\mu_{ab} + (1 - \phi)\mu_a + \phi\mu_b$, onde $\phi = n_A N_{ab} / (n_A N_b + n_B N_a)$ e $\theta = n_A N_B / (n_A N_B + n_B N_A)$;
3. A partir da condição 2, $AVar(\bar{y}_{RPML}) \geq AVar(\bar{y}_{FB})$;
4. Dadas as condições 2 e 3, quando $\beta_2 = 0$, tem-se que $\bar{y}_{RFB} = \bar{y}_{RH}$, e portanto, analogamente aos estimadores já conhecidos, conclui-se que $AVar(\bar{y}_{RFB}) \geq AVar(\bar{y}_{RH})$;
5. É possível notar que para as estratégias propostas, $AVar(\bar{y}_{RPML}) \geq AVar(\bar{y}_{RH})$.

CAPÍTULO 4

Propriedades Estatísticas dos Estimadores

1 Introdução

Conceitos sobre sequências de populações finitas são importantes para investigar o comportamento assintótico dos estimadores propostos nesta tese. Neste capítulo, tais conceitos são apresentados no contexto da abordagem de cadastro duplo. Para uma revisão da teoria considerando um cadastro apenas, o leitor pode consultar Särndal, Swensson & Wretman (2003), Park (2002) e Fuller (2009).

Considere uma sequência infinita de elementos identificáveis por $k = 1, 2, 3, \dots$, associada a uma sequência infinita de valores y_1, y_2, y_3, \dots , onde y_k é o valor associado ao elemento k . Uma sequência de populações U_1, U_2, U_3, \dots é definida de modo que U_ν contém os N_ν primeiros elementos da sequência infinita, isto é, $U_\nu = \{1, 2, \dots, N_\nu\}$. Dessa forma,

$$U_1 \subset U_2 \subset U_3 \subset \dots \quad \text{e} \quad N_1 < N_2 < N_3 < \dots$$

θ_ν é o valor de um parâmetro de interesse na população U_ν , que é função do vetor $\mathbf{Y}_\nu = \{y_1, y_2, \dots, y_{N_\nu}\}$. Considere a situação em que $U_\nu = \mathcal{A}_\nu \cup \mathcal{B}_\nu$, onde \mathcal{A}_ν e \mathcal{B}_ν representam os conjuntos de elementos da sequência pertencentes aos cadastros A e B respectivamente. Neste caso, $N_\nu = N_{A\nu} + N_{B\nu} - N_{ab\nu}$. O cenário que define a relação entre as sequências \mathcal{A}_ν e \mathcal{B}_ν é descrito pela figura (4.1) a seguir para $\nu = 2$. Outros cenários podem ser verificados no apêndice D.

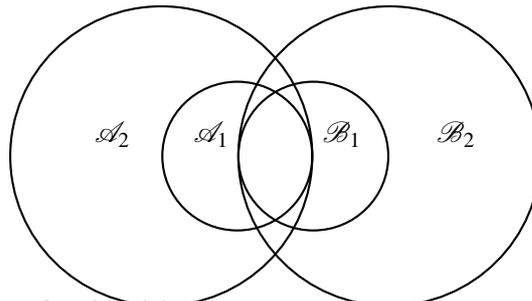


Figura 4.1: Cenário (a) $A_1 \subset A_2$, $B_1 \subset B_2$ e $A_1 \cap B_1 \subset A_2 \cap B_2$

A figura ilustra populações com seqüências para cada cadastro encaixadas. Isto também acontece com a união e a interseção entre os cadastros. Todos os estimadores apresentados nesta tese são baseados neste cenário.

São considerados planos amostrais probabilísticos denotados por $p_{A_v}(\cdot)$ e $p_{B_v}(\cdot)$ para A_v e B_v , que associam probabilidades $p_{A_v}(S_{A_v})$ e $p_{B_v}(S_{B_v})$ para as amostras S_{A_v} e S_{B_v} em cada cadastro. Os tamanhos das amostras obtidas de A_v e B_v são denotados por n_{A_v} e n_{B_v} , e será admitido que $n_{A_1} < n_{A_2} < \dots$ e $n_{B_1} < n_{B_2} < \dots$. As probabilidades de inclusão de primeira ordem dos elementos de A_v e B_v são denotadas por $\pi_{v_k}^A$ e $\pi_{v_k}^B$. As probabilidades de inclusão de segunda ordem dos elementos de A_v e B_v são denotadas por $\pi_{v_{kl}}^A$ e $\pi_{v_{kl}}^B$. Dessa forma, tem-se uma seqüência de populações finitas sob a abordagem de cadastro duplo indexadas por v , tal que, quando $v \rightarrow \infty$, observa-se que $n_{A_v} \rightarrow \infty$, $N_{A_v} \rightarrow \infty$, $n_{B_v} \rightarrow \infty$ e $N_{B_v} \rightarrow \infty$ simultaneamente.

Considere θ_v como um parâmetro de interesse. Sob a abordagem de cadastro duplo, é possível escrever $\theta_v = \theta_{A_v} + \theta_{B_v} - \theta_{abv}$, onde

- θ_{A_v} é um parâmetro de interesse em A_v ;
- θ_{B_v} é um parâmetro de interesse em B_v ;
- θ_{abv} é um parâmetro de interesse em em $A_v \cap B_v$.

Por exemplo, quando θ_v é uma média populacional,

$$\begin{aligned}
 \theta_v &= \mu_v = N_v^{-1} \sum_{k \in U_v} y_k = N_v^{-1} \left(\sum_{k \in A_v} y_k + \sum_{k \in B_v} y_k - \sum_{k \in abv} y_k \right) \\
 &= N_v^{-1} (N_{A_v} \mu_{A_v} + N_{B_v} \mu_{B_v} - N_{abv} \mu_{abv}) \\
 &= w_{A_v} \mu_{A_v} + w_{B_v} \mu_{B_v} - w_{abv} \mu_{abv} \\
 &= \theta_{A_v} + \theta_{B_v} - \theta_{abv}
 \end{aligned} \tag{4.1}$$

onde $w_{Av} = N_{Av}/N_v$, $w_{Bv} = N_{Bv}/N_v$ e $w_{abv} = N_{abv}/N_v$. Um estimador centrado para θ_v é dado por $\hat{\theta}_v = \hat{\theta}_{Av} + \hat{\theta}_{Bv} - \hat{\theta}_{abv}$. Para o caso da média populacional,

$$\begin{aligned}\hat{\theta}_{Av} &= w_A \sum_{k \in S_{Av}} (y_k / \pi_{vk}^A) \\ \hat{\theta}_{Bv} &= w_B \sum_{k \in S_{Bv}} (y_k / \pi_{vk}^B) \\ \hat{\theta}_{abv} &= w_{ab} \sum_{k \in S_{Av} \cap ab} (y_k / \pi_{vk}^A) + w_{ab} \sum_{k \in S_{Bv} \cap ab} (y_k / \pi_{vk}^B)\end{aligned}$$

2 Centralidade assintótica e consistência do estimador sob a abordagem de cadastro duplo

A partir da sequência de populações descrita anteriormente, é possível apresentar as seguintes definições:

Definição 1 Um estimador $\hat{\theta}_v$ é **assintoticamente centrado** para θ_v se

$$\lim_{v \rightarrow \infty} [E(\hat{\theta}_v - \theta_v)] = 0,$$

Sob a abordagem de cadastro duplo, tem-se a condição quando

$$\begin{aligned}\lim_{v \rightarrow \infty} [E(\hat{\theta}_{Av} - \theta_{Av})] &= 0, \\ \lim_{v \rightarrow \infty} [E(\hat{\theta}_{Bv} - \theta_{Bv})] &= 0, \\ \lim_{v \rightarrow \infty} [E(\hat{\theta}_{abv} - \theta_{abv})] &= 0.\end{aligned}$$

Definição 2 Um estimador $\hat{\theta}_v$ é **consistente de acordo com o plano** para θ_v se, para qualquer constante $\varepsilon > 0$,

$$\lim_{v \rightarrow \infty} P(|\hat{\theta}_v - \theta_v| > \varepsilon) = 0$$

Sob a abordagem de cadastro duplo, tem-se a condição quando

$$\begin{aligned}\lim_{v \rightarrow \infty} P(|\hat{\theta}_{Av} - \theta_{Av}| > \varepsilon) &= 0, \\ \lim_{v \rightarrow \infty} P(|\hat{\theta}_{Bv} - \theta_{Bv}| > \varepsilon) &= 0, \\ \lim_{v \rightarrow \infty} P(|\hat{\theta}_{abv} - \theta_{abv}| > \varepsilon) &= 0.\end{aligned}$$

Condições necessárias para verificação de consistência do estimador de Horvitz-Thompson

Para o caso em que é considerado apenas um cadastro, Isaki & Fuller (1982) provaram que as condições a seguir são suficientes para garantir a propriedade de consistência do estimador de Horvitz-Thompson. Estas condições serão necessárias para o desenvolvimento dos resultados nesta tese.

C1. A forma do estimador de Horvitz-Thompson é dada por

$$\bar{y}_{HT} = \frac{1}{N_V} \sum_{k \in S_V} \frac{y_k}{\pi_{Vk}},$$

C2. As probabilidades de inclusão de primeira e segunda ordem são dadas respectivamente por:

$$\begin{aligned} \pi_{Vk} &= P(k \in S_V) = n_V \left(\sum_{k \in S_V} \omega_k \right)^{-1} \omega_k, & k = 1, \dots, N_V \\ \pi_{Vkl} &= P(k \in S_V, l \in S_V), & k, l = 1, \dots, N_V \end{aligned}$$

C3. Seja

$$g_{Vkl} = \begin{cases} \pi_{Vk}\pi_{Vl} - \pi_{Vkl}, & \text{se } \pi_{Vk}\pi_{Vl} - \pi_{Vkl} > 0 \\ 0, & \text{c.c.} \end{cases}$$

C4. A sequência $\{y_k, \omega_k\}$ é tal que

$$\begin{aligned} \text{i.} & \frac{1}{N_V^2} \sum_{k \neq l=1} g_{Vkl}^r = \mathbf{O}(n_V^{-2r\delta}) \\ \text{ii.} & \frac{1}{N_V} \sum_{i \in U_V} \left[\frac{N_V}{n_V} (y_k - \bar{y}) \right]^{2s} < M < \infty \end{aligned}$$

C5. Para $\delta > 0$ e para todo N_V e μ_V é a média populacional da variável de interesse Y , as probabilidades de inclusão satisfazem:

$$N_V^{-2} \sum_{k \neq l \in U_V} g_{Vkl}^r = \mathbf{O}(n_V^{-2r\delta}),$$

onde para $r = 1$, $\left| \frac{y_k}{\pi_{Vk}} \right|$ é limitado. Se $s = 1$, então $n_V^2 g_{Vkl}$ é limitado para todo k, l .

Para o caso em que o parâmetro de interesse é uma média populacional, estas condições garantem que

$$\bar{y}_{HT} - \mu_v = \mathbf{O}_p \left(n_v^{-1/2} \right).$$

onde $\delta = 1/2$ e $\mathbf{O}_p(\cdot)$ indica ordem em probabilidade, como usado em teoria da probabilidade². Sobre as amostras obtidas de cada cadastro e seus estimadores, ainda é necessário estabelecer as seguintes condições:

A1. Denotando por S_{Av} e S_{Bv} as amostras obtidas de A e B , respectivamente, tem-se que, para qualquer variável de interesse y ,

1. $\max_{k \in S_{Av}} y_k = \mathbf{O}_p \left\{ (n_{Av})^{1/2} \right\}$ e
2. $\max_{k \in S_{Bv}} y_k = \mathbf{O}_p \left\{ (n_{Bv})^{1/2} \right\}$,

Esta condição é satisfeita para qualquer plano amostral em que $\pi_k^A \neq 0$ e $\pi_k^B \neq 0$. O uso simultâneo dos dois cadastros induz à obtenção do vetor que representa a união entre as amostras obtidas de cada cadastro, dado por

$$\mathbf{S}_v = S_{Av} \cup S_{Bv} = S_{av} \cup S_{bv} \cup S_{abv}, \quad \text{onde } S_{abv} = S_{Av} \cap S_{Bv}$$

Tem-se que $n_{Av} = n_{av} + n_{abv}$, $n_{Bv} = n_{bv} + n_{abv}$. O conjunto S_v é apenas uma das possíveis amostras obtidas através do uso simultâneo de dois cadastros. De modo geral, \mathbf{S}_v é um conjunto aleatório definido da seguinte forma:

$$\mathbf{S}_v = \mathbf{S}_v^A \cup \mathbf{S}_v^B = \{k \in \mathcal{A}_v, l \in \mathcal{B}_v \mid I_{vk}^A = 1, I_{vl}^B = 1, 1 \leq k \leq N_v^A; 1 \leq l \leq N_v^B\}.$$

\mathbf{S}_v é uma das $\binom{N_v^A}{n_{Av}} \times \binom{N_v^B}{n_{Bv}}$ amostras resultantes do processo de seleção utilizando cadastro duplo.

A2. Assintoticamente, o estimador de Horvitz-Thompson possui distribuição normal para qualquer variável de interesse y e satisfaz a condição **A1**, conforme descrito por Wu & Rao (2006).

A3. Sobre o tamanho populacional dos domínios, supondo que $N_a > 0$ e $N_b > 0$ e admitindo que $v \rightarrow \infty$, tem-se que

²Para um maior aprofundamento, o leitor pode consultar Sen, Singer & Lima (2010) e Fuller (2009). Uma breve revisão sobre o tema é apresentado no apêndice A. A prova das propriedades para o estimador de Hartley podem ser consultadas no apêndice B, página 106.

$$\begin{aligned} \frac{n_{Av}}{n_{Av} + n_V^B} &\rightarrow c_1 \in (0, 1) & \frac{N_{av}}{N_{Av}} &\rightarrow c_2 \in (0, 1) & \frac{N_V^b}{N_V^B} &\rightarrow c_3 \in (0, 1) \\ \frac{N_{av}}{N_V} &\rightarrow w_a \in (0, 1) & \frac{N_V^b}{N_V} &\rightarrow w_b \in (0, 1) \end{aligned}$$

É possível estender os resultados para um contexto mais amplo, quando F cadastros estão disponíveis. Neste caso, a sequência U_1, U_2, U_3, \dots é considerada de tal forma que

$$U_V = U_V^1 \cup U_V^2 \cup U_V^3 \dots \cup U_V^F = \bigcup_{f=1}^F U_V^f, \quad \text{onde } f = 1, \dots, F.$$

Lohr (2007, 2010) apresenta resultados para $F = 3$. Tem-se ainda que

$$N_V = \sum_{i=1}^F N_V^{U^i} - \sum_{i < j} N_V^{U^i \cap U^j} + \sum_{i < j < k} N_V^{U^i \cap U^j \cap U^k} + \dots + (-1)^{F+1} N_V^{\bigcap_{i=1}^F U^i}.$$

Logo,

$$\left(\bigcup_{f=1}^F U_V^f \right) \subset \left(\bigcup_{f=1}^F U_V^f \right) \subset \left(\bigcup_{f=1}^F U_V^f \right) \subset \dots \quad (4.2)$$

3 Centralidade assintótica e consistência do estimador regressão generalizado

No contexto de populações finitas, ao considerar a situação em que é de interesse estimar uma média populacional $\mu_V = N_V^{-1} \sum_{k \in U_V} y_k$, o estimador regressão generalizado é dado por

$$\hat{\mu}_{yRV} = N_V^{-1} \hat{t}_{yR} = N_V^{-1} \left[\hat{t}_{y\pi V} + \sum_{j=1}^q \hat{\beta}_{Vj} (t_{x_j V} - \hat{t}_{x_j \pi V}) \right],$$

onde \hat{t}_{yR} é o estimador regressão para o total populacional. Tem-se ainda que o estimador original de Horvitz-Thompson é obtido na situação particular em que $\hat{\beta}_{Vj} = 0, \forall j$. Seguindo a descrição de Robinson & Särndal (1983), sob a situação em que é considerado um modelo de superpopulação, $\hat{\beta}_{Vj}$ é função dos vetores aleatórios \mathbf{Y}_V e \mathbf{I}_V , onde $\mathbf{Y}_V = \{Y_1, Y_2, \dots, Y_{N_V}\}$ e $\mathbf{I}_V = \{I_{k1}, I_{k2}, \dots, I_{kN_V}\}$, onde I_{Vk} é uma variável indicadora de inclusão do elemento k na amostra, onde

$$I_{Vk} = \begin{cases} 1, & \text{se } k \in S_V \\ 0, & \text{c.c.} \end{cases}$$

Robinson & Särndal (1983) mostraram que o estimador regressão generalizado é consistente de acordo com o plano amostral e assintoticamente centrado, considerando válidas as condições a seguir:

$$\mathbf{A1.} \limsup_{v \rightarrow \infty} \left\{ \frac{\sum_{k \in U_v} x_{kj}^2}{N_v} \right\} < \infty \quad \text{para } j = 0, 1, \dots, q$$

$$\mathbf{A2.} \limsup_{v \rightarrow \infty} \left\{ \frac{\sum_{k \in U_v} Y_k^2}{N_v} \right\} < \infty, \text{ com probabilidade 1.}$$

$$\mathbf{A3.} \limsup_{v \rightarrow \infty} \left\{ E \left(\sum_{j=0}^q \hat{\beta}_{vj}^2 \mid \mathbf{Y}_v \right) \right\} < \infty, \text{ com probabilidade 1.}$$

$$\mathbf{A4.} \liminf_{v \rightarrow \infty} \left\{ N_v \left[\min_{1 \leq k \leq N_v} \pi_{vk} \right] \right\} = \infty.$$

$$\mathbf{A5.} \lim_{v \rightarrow \infty} \left\{ \max_{1 \leq k \neq l \leq N_v} \left[\frac{\pi_{vkl}}{\pi_{vk}\pi_{vl}} - 1 \right] \right\} = 0.$$

A partir destas considerações, são apresentados os seguintes resultados para o estimadores regressão propostos no capítulo 3:

Resultado 6 O estimador \bar{y}_{RH} é assintoticamente centrado e consistente de acordo com o plano amostral sob as condições **A1-A5**.

Resultado 7 O estimador \bar{y}_{RFB} é assintoticamente centrado e consistente de acordo com o plano amostral sob as condições **A1-A5**.

Resultado 8 Os estimadores \bar{y}_{RBLG} e $\bar{y}_{RBLG(\infty)}$ são assintoticamente centrados e consistentes de acordo com o plano amostral sob as condições **A1-A5**.

Resultado 9 O estimador \bar{y}_{RPML} é assintoticamente centrado e consistente de acordo com o plano amostral sob as condições **A1-A5**.

Prova: Apêndice **B**.

4 Erro quadrático médio assintótico dos estimadores regressão sob a abordagem de cadastro duplo

A seguir serão apresentadas condições necessárias para a minimização do erro quadrático médio, definido como $E\left\{(\bar{y}_{Rv} - \bar{y}_v)^2 | \mathbf{Y}\right\}$, dos estimadores regressão propostos. Para o caso em que apenas um cadastro é considerado,

$$\mu_k = E(Y_k) \quad \sigma_k^2 = \text{Var}(Y_k)$$

Dessa forma, em cada cadastro, adicionalmente são assumidas as seguintes condições sobre a estrutura das observações da variável de interesse:

A6. Os Y_k 's são não-correlacionados e

$$\limsup_{v \rightarrow \infty} \frac{\sum_{k \in U_v} \mu_k^2}{N} < \infty, \quad \text{e} \quad \limsup_{v \rightarrow \infty} \frac{\sum_{k \in U_v^d} \mu_k^2}{N_v} < \infty; \quad \sigma_k^2 < \infty$$

A7. Para um conjunto de q variáveis auxiliares no cadastro, existem constantes β_j tal que para todo v suficientemente grande,

$$n_v \sum_{j=0}^q E \left[\left(\hat{\beta}_{vj} - \beta_j \right)^2 \right] < \infty,$$

$$\text{A8.} \quad \liminf_{v \rightarrow \infty} \frac{N_v}{n_v} \left\{ \min_{1 \leq k \leq N_v} \pi_{vk} \right\} > 0$$

$$\text{A9.} \quad \limsup_{v \rightarrow \infty} \frac{(N_v)^2}{n_v} \left\{ \max_{1 \leq k \neq l \leq N_v} |\pi_{vkl} - \pi_{vk}\pi_{vl}| \right\} < \infty$$

A10. Existe uma constante γ tal que

$$\mu_k = \gamma \pi_{vk} + \sum_{j=0}^q \beta_j x_{jk}, \quad 1 \leq k \leq N_v$$

A partir destas condições para cadastros, as quais também podem ser admitidas para os domínios, serão apresentadas expressões para os EQM's dos estimadores propostos para cada uma das estratégias de estimadores regressão propostas nesta tese.

Resultado 10 Dada uma sequência de de populações finitas U_v , tem-se que para cada um dos cadastros, sob as condições **A1**, **A2**, **A5**, **A6–A9** o erro quadrático médio assintótico do estimador regressão é dado por

$$n_v E(\bar{y}_{Rv} - \bar{y}_v)^2 = D_v + W_v + Q_v < \infty, \text{ quando } v \rightarrow \infty$$

onde

$$\begin{aligned} D_v &= \frac{n_v}{(N_v)^2} \sum_{k \in U_v} \sigma_k^2 \left(\frac{1}{\pi_{vk} - 1} \right), \\ W_v &= \frac{n_v}{(N_v)^2} \sum_{k \in U_v} \sum_{\substack{l \in U_v \\ k \neq l}} \sigma_k^2 (\pi_{kl} - \pi_{vk} \pi_{vl}) \left(\frac{g_k}{\pi_{vk}} - \frac{g_l}{\pi_{vl}} \right)^2, \text{ com } g_k = \mu_k - \sum_{j=0}^q \beta_j x_{kj} \\ Q_v &= \frac{n_v}{(N_v)^2} E \left\{ \left[\sum_{j=0}^q (\hat{\beta}_{vj} - \beta_j) \right] \frac{1}{N_v} \sum_{k \in U_v} x_{kj} \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right) \right\}^2 + \\ &+ 2E \left\{ \left[\frac{\sqrt{n_v}}{N_v} \sum_{k \in U_v} (y_k - \mu_k) \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right) \right] \left[\frac{\sqrt{n_v}}{N_v} \sum_{j=0}^q (\hat{\beta}_{vj} - \beta_j) \frac{1}{N_v} \sum_{k \in U_v} x_{kj} \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right) \right] \right\} \\ &+ 2E \left\{ \left[\frac{\sqrt{n_v}}{N_v} \sum_{k \in U_v} (y_k - \mu_k) \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right) \right] \frac{\sqrt{n_v}}{N_v} \sum_{k \in U_v} g_k \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right) \right\} \\ &+ 2E \left\{ \left[\frac{\sqrt{n_v}}{N_v} \sum_{j=0}^q (\hat{\beta}_{vj} - \beta_j) \frac{1}{N_v} \sum_{k \in U_v} x_{kj} \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right) \right] \frac{\sqrt{n_v}}{N_v} \sum_{k \in U_v} g_k \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right) \right\} \end{aligned}$$

com $Q_v \rightarrow 0$, quando $v \rightarrow \infty$.

Prova: Seja

$$\sqrt{n}(\bar{y}_{Rv} - \bar{y}_v) = m_1 + m_2 + m_3,$$

onde

$$m_1 = \frac{\sqrt{n_v}}{N_v} \sum_{k \in U} (y_k - \mu_k) \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right)$$

$$m_2 = \frac{\sqrt{n_v}}{N_v} \sum_{j=0}^q (\hat{\beta}_{vj} - \beta_j) \frac{1}{N} \sum_{k \in U} x_{kj} \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right)$$

$$m_3 = \frac{\sqrt{n_v}}{N_v} \sum_{k \in U} g_k \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right)$$

Logo,

$$n_v E(\bar{y}_{RV} - \bar{y}_v)^2 = E(m_1^2) + E(m_2^2) + E(m_3^2) + E(m_1 m_2) + E(m_1 m_3) + E(m_2 m_3)$$

É possível notar que,

$$E(m_2) = E\{E(m_2^2 | \mathbf{I})\} = E\left\{ \frac{n_v}{N_v} \sum_{k \in U} \sigma_k^2 \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right)^2 \right\} = D_v.$$

Sob **A6** e **A8** tem-se que

$$\limsup_{v \rightarrow \infty} D_v \leq \limsup_{v \rightarrow \infty} \left(\frac{n_v}{N_v \min \pi_{vk}} \right) \frac{1}{N_v} \sum_{k \in U_v} \sigma_k^2 < \infty \quad (4.3)$$

Ainda, sob **A1**, **A4**, **A5** e **A7**,

$$E(m_2^2) \leq E\left\{ \frac{n_v}{(N_v)^2} \sum_{j=0}^q E\left[(\hat{\beta}_{vj} - \beta_j)^2 | \mathbf{I} \right] \sum_{j=0}^q (\hat{\beta}_{vj} - \beta_j) \frac{1}{N_v} \sum_{k \in U_v} x_{kj} \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right) \right\} \rightarrow 0, \quad (4.4)$$

quando $v \rightarrow \infty$.

$$E(m_3^2) = \frac{n_v}{(N_v)^2} \sum_{k \in U_v} g_k^2 \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right) + \frac{n_v}{(N_v)^2} \sum_{k \in U_v} \sum_{\substack{l \in U_v \\ k \neq l}} g_k g_l \left(\frac{I_{vk}}{\pi_{vk} \pi_{vl}} - 1 \right) = W_v, \quad (4.5)$$

onde

$$\limsup_{v \rightarrow \infty} W_v \leq \limsup_{v \rightarrow \infty} \left\{ \left(\frac{N_v^2}{n_v} \max_v |\pi_{vkl} - \pi_{vk} \pi_{vl}| \right) \left(\frac{n_v^2}{N_v^2 \max \{\pi_{vk}^2\}} \right) \frac{8}{N_v} \left(\sum_{k \in U_v} \mu_k^2 + 2 \sum_{j=0}^q \beta_j^2 \sum_{k \in U_v} x_{kj}^2 \right) \right\} < \infty. \quad (4.6)$$

Ainda, sob **A1**, **A7**, **A9** e **A10**:

$$E(m_1 m_2) \leq [E(m_1)^2 E(m_2)^2]^{1/2} \rightarrow 0,$$

$$E(m_1 m_3) = E[E(m_1 | \mathbf{I}_v) m_2]^{1/2} = 0,$$

$$E(m_2 m_3) \leq [E(m_2)^2 E(m_3)^2]^{1/2} \rightarrow 0,$$

o que leva ao resultado 2. ■

Todas as condições apresentadas são utilizadas sob a suposição de que n_v cresça com mesma taxa de N_v . Se isso não acontecer, tem-se que $D_v \rightarrow 0$ e $W_v \rightarrow 0$, o que ilustra a consistência do estimador, mas que não leva a um resultado interessante por não apresentar uma forma geral para o EQM.

Considerando **C15**, π_{vk} precisa ser constante para todo v , se $\gamma \neq 0$. Sob a perspectiva de um plano amostral de amostragem aleatória simples, por exemplo, D_v é mínimo quando

$$\pi_{vk} = \frac{n\sigma_k}{\sum_{k \in U_v} \sigma_k}.$$

Dessa forma, tem-se que

$$EQM(\bar{y}_{Rv}) = \frac{\left[(N_v)^{-1} \sum_{k \in U_v} \sigma_k \right]^2 - \left(\frac{n_v}{(N_v)^2} \right) \sum_{k \in U_v} \sigma_k^2}{n_v}$$

A seguir, serão apresentadas expressões para o erro quadrático médio assintótico de cada um dos estimadores regressão propostos.

4.1 Estratégia de Hartley

Sob a estratégia de Hartley, o erro quadrático médio do estimador \bar{y}_{RH} é dado pela expressão:

$$\begin{aligned} EQM(\bar{y}_{RH}) = E(\bar{y}_{RH} - \mu)^2 &= (w_A)^2 E(\bar{y}_{RA} - \mu_A)^2 + (w_B)^2 E(\bar{y}_{RB} - \mu_B)^2 \\ &+ 2w_A w_B E\{(\bar{y}_{Ra} - \mu_A)(\bar{y}_{Rb} - \mu_B)\} \\ &- 2w_A w_{ab} p E\left[(\bar{y}_{Ra} - \mu_A)(\bar{y}_{Rab(A)} - \mu_{ab})\right] - \end{aligned}$$

$$\begin{aligned}
 & - 2w_A w_{ab}(1-p)E \left[(\bar{y}_{Ra} - \mu_A) (\bar{y}_{Rab(B)} - \mu_{ab}) \right] \\
 & - 2w_B w_{ab} p E \left[(\bar{y}_{Rb} - \mu_B) (\bar{y}_{Rab(A)} - \mu_{ab}) \right] \\
 & - 2w_B w_{ab}(1-p)E \left[(\bar{y}_{Rb} - \mu_B) (\bar{y}_{Rab(B)} - \mu_{ab}) \right] \\
 & + p^2 (w_{ab})^2 E \left(\bar{y}_{Rab(A)} - \mu_{ab} \right)^2 \\
 & + 2p(1-p)(w_{ab})^2 E \left[\left(\bar{y}_{Rab(A)} - \mu_{ab} \right) \left(\bar{y}_{Rab(B)} - \mu_{ab} \right) \right] \\
 & + (w_{ab})^2 (1-p)^2 E \left[\left(\bar{y}_{Rab(B)} - \mu_{ab} \right)^2 \right] \tag{4.7}
 \end{aligned}$$

Quando $v \rightarrow \infty$, tem-se que

$$\begin{aligned}
 E \{ (\bar{y}_{RAv} - \mu_{Av}) (\bar{y}_{RBv} - \mu_{Bv}) \} & = 0; \\
 E \left\{ \left(\bar{y}_{Rab(A)v} - \mu_{abv} \right) \left(\bar{y}_{Rab(B)v} - \mu_{abv} \right) \right\} & = 0; \\
 E \left\{ (\bar{y}_{RAv} - \mu_{Av}) \left(\bar{y}_{Rab(B)v} - \mu_{abv} \right) \right\} & = 0; \\
 E \left\{ \left(\bar{y}_{Rab(A)v} - \mu_{abv} \right) \left(\bar{y}_{Rab(B)v} - \mu_{abv} \right) \right\} & = 0,
 \end{aligned}$$

pois planos independentes são aplicados em A e B . Além disso, tem-se que,

$$\begin{aligned}
 E \left\{ (\bar{y}_{RAv} - \mu_{Av}) \left(\bar{y}_{Rab(A)v} - \mu_{vab} \right) \right\} & \longrightarrow C_{(A,ab(A))}; \\
 E \left\{ (\bar{y}_{RBv} - \mu_{Bv}) \left(\bar{y}_{Rab(B)v} - \mu_{ab} \right) \right\} & \longrightarrow C_{(B,ab(B))}; \\
 E \left\{ (\bar{y}_{RAv} - \mu_{Av})^2 \right\} & = D_v^A + W_v^A < \infty; \\
 E \left\{ (\bar{y}_{RBv} - \mu_{Bv})^2 \right\} & = D_v^B + W_v^B < \infty; \\
 E \left\{ \left(\bar{y}_{Rab(A)} - \mu_{ab} \right)^2 \right\} & = D_v^{ab(A)} + W_v^{ab(A)} < \infty; \\
 E \left\{ \left(\bar{y}_{Rab(B)v} - \mu_{ab} \right)^2 \right\} & = D_v^{ab(B)} + W_v^{ab(B)} < \infty,
 \end{aligned}$$

$C_{(s_1, s_2)}$ indica a covariância entre as quantidades provenientes dos conjuntos s_1 e s_2 , conforme descrito por Särndal (1992). Logo,

$$\begin{aligned}
 EQM(\bar{y}_{RHv}) = E(\bar{y}_{RHv} - \mu_v)^2 &= (w_{vA})^2 (D_v^A + W_v^A) + (w_{vB})^2 (D_v^B + W_v^B) \\
 &- 2w_{vA}w_{vab}pC_{(A,ab(A))} - 2w_{vB}w_{vab}(1-p)C_{(A,ab(A))} \\
 &+ p^2 (D_v^{ab(A)} + W_v^{ab(A)}) \\
 &+ (w_{vab})^2 (1-p)^2 E(D_v^{ab(B)} + W_v^{ab(B)}) \\
 &\leq (w_{vA})^2 (D_v^A + W_v^A) + (w_{vB})^2 (D_v^B + W_v^B) \\
 &+ p^2 (D_v^{ab(A)} + W_v^{ab(A)}) \\
 &+ (w_{vab})^2 (1-p)^2 E(D_v^{ab(B)} + W_v^{ab(B)}) \tag{4.8}
 \end{aligned}$$

Resultado 11 Admitindo que $w_A = \tau_A$, $w_B = \tau_B$ e $w_{vab} = \tau_{ab}$, tem-se então que o erro quadrático médio assintótico do estimador regressão sob a estratégia de Hartley é dado por

$$\begin{aligned}
 EQM(\bar{y}_{RHv}) &\leq \tau_A^2 EQM(\bar{y}_{RAv}) + \tau_B^2 EQM(\bar{y}_{RBv}) \\
 &+ p^2 EQM(\bar{y}_{Rab(A)v}) + (1-p)^2 \tau_{ab}^2 EQM(\bar{y}_{Rab(B)v}) < \infty.
 \end{aligned}$$

4.2 Estratégia de Fuller & Burmeister

Sob a estratégia em que a magnitude dos domínios é desconhecida tem-se (apêndice B),

$$\begin{aligned}
 \bar{y}_{RFB} - \mu_v &\approx w_a(\bar{y}_{RAv} - \mu_{va}) + w_{vb}(\bar{y}_{RBv} - \mu_b) + p_1 w_{vab}(\bar{y}_{vRab(A)} - \mu_{vab}) + (1-p_1)w_{vab}(\bar{y}_{Rab(B)v} - \mu_{vab}) \\
 &+ N^{-1}(\hat{N}_a \mu_{va} + \hat{N}_b \mu_{vb} + p_1 \hat{N}_{ab(A)} \mu_{vab} + (1-p_1) \hat{N}_{vab(B)} \mu_{vab} + p_2 \hat{N}_{ab(A)} - p_2 \hat{N}_{ab(B)})
 \end{aligned}$$

$$\begin{aligned}
 &= w_a (\bar{y}_{Rav} - \mu_{va}) + p_1 w_{vab} (\bar{y}_{Rab(A)v} - \mu_{vab}) \\
 &+ w_b (\bar{y}_{Rbv} - \mu_{vb}) + (1 - p_1) w_{vab} (\bar{y}_{Rab(B)v} - \mu_{vab}) + \mathbf{O}(1).
 \end{aligned}$$

Dessa forma,

$$\begin{aligned}
 (\bar{y}_{RFBv} - \mu_v)^2 &\approx (w_{va})^2 (\bar{y}_{Rav} - \mu_{va})^2 + 2p_1 w_{va} w_{vab} (\bar{y}_{Rav} - \mu_{va}) (\bar{y}_{Rab(A)v} - \mu_{vab}) \\
 &+ p_1^2 (w_{vab})^2 (\bar{y}_{Rab(A)v} - \mu_{vab})^2 + 2w_{va} w_{vb} (\bar{y}_{Rav} - \mu_{va}) (\bar{y}_{Rbv} - \mu_{vb}) \\
 &+ 2(1 - p_1) w_{va} w_{vab} (\bar{y}_{Rav} - \mu_{va}) (\bar{y}_{Rab(B)v} - \mu_{vab}) \\
 &+ 2p_1 w_{vb} w_{vab} (\bar{y}_{Rbv} - \mu_{vb}) (\bar{y}_{Rab(A)v} - \mu_{vab}) \\
 &+ 2p_1 (1 - p_1) (w_{vab})^2 (\bar{y}_{Rab(A)v} - \mu_{vab}) (\bar{y}_{Rab(B)v} - \mu_{vab}) \\
 &+ (w_{vb})^2 (\bar{y}_{Rbv} - \mu_{vb})^2 + 2(1 - p_1) w_{vb} w_{vab} (\bar{y}_{Rbv} - \mu_{vb}) (\bar{y}_{Rab(B)v} - \mu_{vab}) \\
 &+ (1 - p_1)^2 (w_{vab})^2 (\bar{y}_{Rab(B)v} - \mu_{vab})^2 + \mathbf{O}(N_v^{-1})
 \end{aligned}$$

Calculando o valor esperado, ou seja, $E(\bar{y}_{RFBv} - \mu_v)^2$ e admitindo $v \rightarrow \infty$, tem-se

$$\begin{aligned}
 E(\bar{y}_{Rav} - \mu_{va})^2 &= D_v^a + W_v^a < \infty \\
 E(\bar{y}_{Rbv} - \mu_{vb})^2 &= D_v^b + W_v^b < \infty \\
 E(\bar{y}_{Rab(A)v} - \mu_{vab})^2 &= D_v^{ab(A)} + W_v^{ab(A)} < \infty \\
 E(\bar{y}_{Rab(B)v} - \mu_{vab})^2 &= D_v^{ab(B)} + W_v^{ab(B)} < \infty \\
 E\left\{(\bar{y}_{Rav} - \mu_{va})(\bar{y}_{Rbv} - \mu_{vb})\right\} &= 0 \\
 E\left\{(\bar{y}_{Rab(A)v} - \mu_{vab})(\bar{y}_{Rab(B)v} - \mu_{vab})\right\} &= 0 \\
 E\left\{(\bar{y}_{Rav} - \mu_{va})(\bar{y}_{Rab(B)v} - \mu_{vab})\right\} &= 0
 \end{aligned}$$

$$E \left\{ (\bar{y}_{Rbv} - \mu_{vb}) (\bar{y}_{Rab(A)v} - \mu_{vab}) \right\} = 0$$

$$E \left\{ (\bar{y}_{Rav} - \mu_{va}) (\bar{y}_{Rab(A)v} - \mu_{vab}) \right\} = C_{(a;ab(A))} < \infty$$

$$E \left\{ (\bar{y}_{Rbv} - \mu_{bv}) (\bar{y}_{Rab(B)v} - \mu_{vab}) \right\} = C_{(b;ab(A))} < \infty$$

Resultado 12 O erro quadrático médio assintótico do estimador regressão sob a estratégia de Fuller & Burmeister é dado por

$$\begin{aligned} EQM(\bar{y}_{RFB}) &\approx (w_{va})^2 (D_v^a + W_v^a)^2 + p_1^2 (w_{vab})^2 (D_v^{ab(A)} + W_v^{ab(A)}) \\ &+ (w_{vb})^2 (D_v^b + W_v^b) + (1 - p_1)^2 (w_{vab})^2 (D_v^{ab(B)} + W_v^{ab(B)}) \\ &+ 2p_1 w_{va} w_{vab} C_{(a;ab(A))} + 2(1 - p_1) w_{vb} w_{vab} C_{(b;ab(B))} \end{aligned}$$

4.3 Estratégia BLG

Sob a estratégia adotada por Bankier, Lepkowski & Groves (1986), é possível notar que a forma do estimador apresentada em (3.8) é inicialmente escrita da seguinte forma:

$$\bar{y}_{RBLGv} = \frac{\sum_{k \in S_{Av}} g_{ks(A)}^* \omega_k y_k + \sum_{k \in S_{Bv}} g_{ks(B)}^* \omega_k y_k}{N} = w_{Av} \bar{y}_{RAv} + w_{Bv} \bar{y}_{Rv},$$

onde

$$\omega_k = \begin{cases} \frac{1}{\pi_{yk}^A}, & k \in a \\ \frac{1}{\pi_{vk}^B}, & k \in b \\ \frac{1}{\pi_{vk}^{A \cup B}} = \frac{1}{\pi_{vk}^A + \pi_{vk}^B}, & k \in ab \end{cases}$$

Dessa forma, tem-se que

$$\bar{y}_{RBLGv} - \mu = w_v^A (\bar{y}_{Rv}^A - \mu_v^A) + w_v^B (\bar{y}_{Rv}^B - \mu_v^B)$$

$$\begin{aligned} (\bar{y}_{RBLGv} - \mu_v)^2 &= (w_v^A)^2 (\bar{y}_{Rv}^A - \mu_v^A)^2 + (w_v^B)^2 (\bar{y}_{Rv}^B - \mu_v^B)^2 \\ &+ 2w_v^A w_v^B (\bar{y}_{Rv}^A - \mu_v^A) (\bar{y}_{Rv}^B - \mu_v^B) \end{aligned}$$

$$E(\bar{y}_{RBLGv} - \mu_v)^2 = (w_v^A)^2 E(\bar{y}_{Rv}^A - \mu_v^A)^2 + (w_v^B)^2 E(\bar{y}_{Rv}^B - \mu_v^B)^2 \quad (4.9)$$

Resultado 13 Para $w_v^A \rightarrow \tau_A$ e $w_v^B \rightarrow \tau_B$, tem-se que o erro quadrático médio do estimador regressão sob a estratégia BLG é dado por

$$EQM(\bar{y}_{RBLG}) = \tau_A^2 EQM(\bar{y}_{Rv}^A) + \tau_B^2 EQM(\bar{y}_{Rv}^B)$$

A partir da estratégia de Hartley, é possível notar que $EQM(\bar{y}_{RBLG}) < EQM(\bar{y}_{RHv})$. Os EQM's serão aproximadamente iguais quando os EQM's dos estimadores referentes aos domínios são aproximadamente nulos.

Para a versão *raking* o estimador regressão apresentado em (3.11) também temos resultado similar, basta notar que existem constantes α^A, α^B que satisfazem:

$$\begin{aligned} \omega_{ka}^{(\infty)} &= \frac{1}{\pi_{vk}^A} \left[\prod_{r \text{ par}} \left(\frac{N_{Av}}{\hat{N}^{A(r-1)}} \right) \right] = \frac{1}{\pi_{vk}^A} \alpha^A \\ \omega_{kb}^{(\infty)} &= \frac{1}{\pi_{vk}^B} \left[\prod_{r \text{ impar}} \left(\frac{N_v^B}{\hat{N}^{B(r-1)}} \right) \right] = \frac{1}{\pi_{vk}^B} \alpha^B \\ \omega_{kab}^{(\infty)} &= \frac{1}{\pi_{vk}^A + \pi_{vk}^B} \left[\prod_{r \text{ par}} \left(\frac{N_{Av}}{\hat{N}^{A(r-1)}} \right) \prod_{r \text{ impar}} \left(\frac{N_v^B}{\hat{N}^{B(r-1)}} \right) \right] = \left(\frac{1}{\pi_{vk}^A + \pi_{vk}^B} \right) \alpha^A \alpha^B \quad (4.10) \end{aligned}$$

4.4 Estratégia de Máxima Pseudo-Verossimilhança

Com base na estratégia de Máxima Pseudo-Verossimilhança, é possível notar que

$$\begin{aligned} \bar{y}_{RPML} &= w_{va} \bar{y}_{Rv}^a + w_{vb} \bar{y}_{Rv}^b + w_{vab} \left(\frac{\pi_{vk}^A}{\pi_{vk}^A + \pi_{vk}^B} \right) \bar{y}_{Rab(A)} + w_{vab} \left(\frac{\pi_{vk}^B}{\pi_{vk}^A + \pi_{vk}^B} \right) \bar{y}_{Rv}^{ab(B)} \\ &+ \hat{w}^a \mu_{va} + \hat{w}^b \mu_{vb} + \hat{w}^{ab} \mu_{vab} - w_{va} \mu_{va} - w_{vb} \mu_{vb} - w_{vab} \mu_{vab} \end{aligned}$$

Logo,

$$\begin{aligned} \bar{y}_{RPMLV} - \mu_v &\approx w_{va} (\bar{y}_{Rv}^a - \mu_{va}) + w_{vb} (\bar{y}_{Rv}^b - \mu_{vb}) \\ &+ w_{vab} \left(\frac{\pi_{vk}^A}{\pi_{vk}^A + \pi_{vk}^B} \right) \bar{y}_{Rab(A)v} + w_{vab} \left(\frac{\pi_{vk}^B}{\pi_{vk}^A + \pi_{vk}^B} \right) \bar{y}_{Rab(B)v} + O(1) \end{aligned}$$

No caso em que $\left(\frac{\pi_{vk}^A}{\pi_{vk}^A + \pi_{vk}^B} \right) = p_1$ e $\left(\frac{\pi_{vk}^B}{\pi_{vk}^A + \pi_{vk}^B} \right) = (1 - p_1)$, onde p_1 é uma constante de ponderação, tal que $0 \leq p_1 \leq 1$, o estimador \bar{y}_{RPML} assume forma similar à do estimador \bar{y}_{RFB} . Dessa forma, $EQM(\bar{y}_{RPMLV})$ terá forma similar à $EQM(\bar{y}_{RFBV})$.

5 Distribuições assintóticas dos estimadores propostos

Skinner & Rao (1996), mostraram que é possível representar os estimadores \bar{y}_H , \bar{y}_{FB} e \bar{y}_{PML} como funções do vetor

$$\hat{\eta} = \left(\bar{y}_{HTa} \quad \bar{y}_{HTab(A)} \quad \hat{N}^{ab(A)}/N \quad \bar{y}_{HTb} \quad \bar{y}_{HTab(B)} \quad \hat{N}^{ab(B)}/N \right),$$

que é estimador de

$$\eta = \left(\mu_a \quad \mu_{ab} \quad N_{ab}/N \quad \mu_b \quad \mu_{ab} \quad N_{ab}/N \right).$$

O estimador \bar{y}_{BLG} também pode ser representado de forma análoga, como função de

$$\hat{\eta} = \left(\bar{y}_{HT(A)} \quad \bar{y}_{HT(B)} \right)$$

sendo mais simples que a forma apresentada para os outros estimadores. Para todas as estratégias, $\hat{\eta}$ é consistente para η . Logo,

$$n^{1/2} (\hat{\eta} - \eta) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \Sigma)$$

onde $n = n_A + n_B - n_{ab}$. Σ é uma matriz bloco diagonal da forma

$$\Sigma = \begin{pmatrix} \Sigma_A & \mathbf{0} \\ \mathbf{0} & \Sigma_B \end{pmatrix}$$

onde

$$\Sigma_A = \begin{pmatrix} \sigma_a^2 & \sigma_{a,ab(A)} & 0 \\ \sigma_{a,ab(A)} & \sigma_{ab(A)}^2 & 0 \\ 0 & 0 & \sigma_{N_{ab}^A} \end{pmatrix} \quad \Sigma_B = \begin{pmatrix} \sigma_b^2 & \sigma_{b,ab(B)} & 0 \\ \sigma_{b,ab(B)} & \sigma_{ab(B)}^2 & 0 \\ 0 & 0 & \sigma_{N_{ab}^B} \end{pmatrix}$$

Os estimadores do tipo regressão propostos também podem ser escritos como funções de um vetor denotado por

$$\hat{\phi} = \left(\bar{y}_{Ra} \quad \bar{y}_{Rab(A)} \quad \hat{N}^{ab(A)}/N \quad \bar{y}_{Rb} \quad \bar{y}_{Rv}^{ab(B)} \quad \hat{N}^{ab(B)}/N \right),$$

Foram apresentadas condições no capítulo 3 que tornam $\hat{\phi}$ consistente para ϕ . Logo,

$$n^{1/2}(\hat{\phi} - \phi) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \Sigma_*).$$

onde, de acordo com a estratégia utilizada, Σ_* é dada por Σ_{RH} , Σ_{RFB} , Σ_{RBLG} ou Σ_{RPML} . Todos os estimadores propostos podem ser dados na forma $\hat{\theta} = f(\hat{\phi})$, onde $f(\cdot)$ é uma função contínua, não linear e diferenciável em relação aos componentes de $\hat{\phi}$. Por exemplo, sob a estratégia de Hartley, $f(\hat{\phi}) = \bar{y}_{RH}$. Logo, pelo método método delta, tem-se que

$$n^{1/2}(f(\hat{\phi}) - f(\phi)) \xrightarrow{d} \mathcal{N}(0, [f'(\phi)]^T \Sigma_* [f'(\phi)]). \quad (4.11)$$

onde $f'(\phi)$ é vetor de derivadas parciais. Ao substituir quantidades populacionais pelas suas respectivas quantidades amostrais, é razoável considerar que um intervalo de confiança aproximado para $f(\phi)$ é dado por

$$f(\hat{\phi}) \pm z_{1-\alpha/2} \sqrt{[f'(\phi)]^T \hat{\Sigma}_* [f'(\phi)]/n^{1/2}} \quad (4.12)$$

CAPÍTULO 5

Inferência para casos particulares do estimador regressão generalizado

Nos capítulos anteriores foram apresentados vários estimadores do tipo regressão para cadastro duplo. Neste capítulo são ilustrados casos particulares destes estimadores, considerando que um plano de amostragem aleatória simples é aplicado em cada um dos cadastros.

1 Estimadores do tipo razão sob a estratégia de Hartley

Considerando-se o estimador regressão sob a estratégia de Hartley apresentado pela expressão (3.4), e considerando o cenário em que o estimador regressão reduz-se ao caso de estimadores do tipo razão, é possível obter os estimadores apresentados a seguir.

1.1 Estimador razão RH1

Sob a situação particular em que uma variável auxiliar está disponível em ambos os cadastros, e considerando que um plano de amostragem aleatória simples é empregado em cada cadastro, o estimador RH1 do tipo razão para μ , sob a estratégia de Hartley, é considerado para situações em que o total das variáveis auxiliares é conhecido para os domínios. Tem-se que

$$\bar{y}_{RH1} = w_a \frac{\tilde{y}_a}{\tilde{x}_a} \mu_{X_a} + w_b \frac{\tilde{y}_b}{\tilde{x}_b} \mu_{X_b} + w_{ab} p \frac{\tilde{y}_{ab(A)}}{\tilde{x}_{ab(A)}} \mu_{X_{ab}} + w_{ab} (1-p) \frac{\tilde{y}_{ab(B)}}{\tilde{x}_{ab(B)}} \mu_{X_{ab}}, \quad (5.1)$$

onde

- $w_a = N_a/N$;

- $w_b = N_b/N$;
- $w_{ab} = N_{ab}/N$;
- μ_{X_b} é a média populacional da variável auxiliar no domínio b ;
- μ_{X_a} é a média populacional da variável auxiliar no domínio a ;
- μ_{X_b} é a média populacional da variável auxiliar no domínio b ;
- \tilde{y}_a é a média amostral da variável de interesse no domínio a ;
- \tilde{y}_b é a média amostral da variável de interesse no domínio b ;
- $\tilde{y}_{ab(A)}$ é a média amostral das informações da variável de interesse no cadastro A e que pertencem ao domínio ab ;
- $\tilde{y}_{ab(B)}$ é a média amostral das informações da variável de interesse no cadastro B e que pertencem ao domínio ab .
- \tilde{x}_a é a média amostral da variável auxiliar no domínio a ;
- \tilde{x}_b é a média amostral da variável auxiliar no domínio b ;
- $\tilde{x}_{ab(A)}$ é a média amostral das informações da variável auxiliar no cadastro A e que pertencem ao domínio ab ;
- $\tilde{x}_{ab(B)}$ é a média amostral das informações da variável auxiliar no cadastro B e que pertencem ao domínio ab .

\tilde{y}_a , \tilde{y}_b , e $\tilde{y}_{ab(B)}$ representam as médias amostrais dos valores pertencentes aos domínios a , b e ab e identificados nas amostras de tamanho n_A e n_B obtidas de cada um dos cadastros. A variância aproximada deste estimador, expressa em termos matriciais, é dada por $A\text{Var}(\bar{y}_{RH1}) = \mathbf{d}^T \Sigma_{RH1} \mathbf{d}$, onde

$$\mathbf{d} = \begin{bmatrix} 1 \\ p \\ -R_a \\ -pR_{ab(A)} \\ 1 \\ 1-p \\ -R_b \\ -(1-p)R_{ab(B)} \end{bmatrix} \quad \text{e} \quad \Sigma_{RH1} = \begin{pmatrix} \Sigma_A & \mathbf{0} \\ \mathbf{0} & \Sigma_B \end{pmatrix},$$

onde

$$\Sigma_A = \begin{pmatrix} \Sigma_{y(a;ab(A))} & \Sigma_{xy(a;ab(A))} \\ \Sigma_{xy(a;ab(A))} & \Sigma_{x(a;ab(A))} \end{pmatrix} = \begin{pmatrix} \sigma_{ya}^2 & \sigma_{(ya;yab(A))} & \sigma_{(ya;xa)} & \sigma_{(ya;xab(A))} \\ \sigma_{(ya;yab(A))} & \sigma_{yab(A)}^2 & \sigma_{(xa;yab(A))} & \sigma_{(yab(A);xab(A))} \\ \sigma_{(ya;xa)} & \sigma_{(xa;yab(A))} & \sigma_{xa}^2 & \sigma_{(xa;xab(A))} \\ \sigma_{(ya;xab(A))} & \sigma_{(xab(A);yab(A))} & \sigma_{(xa;xab(A))} & \sigma_{xab(A)}^2 \end{pmatrix}$$

$$\Sigma_B = \begin{pmatrix} \Sigma_{y(b;ab(B))} & \Sigma_{xy(b;ab(B))} \\ \Sigma_{xy(b;ab(B))} & \Sigma_{x(b;ab(B))} \end{pmatrix} = \begin{pmatrix} \sigma_{yb}^2 & \sigma_{(yb;yab(B))} & \sigma_{(yb;xb)} & \sigma_{(yb;xab(B))} \\ \sigma_{(yb;yab(B))} & \sigma_{yab(B)}^2 & \sigma_{(xa;yab(B))} & \sigma_{(yab(B);xab(B))} \\ \sigma_{(yb;xb)} & \sigma_{(xb;yab(B))} & \sigma_{xb}^2 & \sigma_{(xb;xab(B))} \\ \sigma_{(yb;xab(B))} & \sigma_{(xab(B);yab(B))} & \sigma_{(xb;xab(B))} & \sigma_{xab(B)}^2 \end{pmatrix}$$

A forma da variância aproximada deste estimador é dada por

$$AVar(\bar{y}_{RH1}) = N^{-2} \sigma_{RH1}^2,$$

$$\begin{aligned} \text{em que } \sigma_{RH1}^2 &= (N_A)^2 \left(\frac{1-f_A}{n_A} \right) w_a \sigma_{ya}^2 + p^2 (N_A)^2 \left(\frac{1-f_A}{n_A} \right) w_{ab(A)} \sigma_{yab(A)}^2 \\ &+ R_a^2 (N_A)^2 \left(\frac{1-f_A}{n_A} \right) w_a \sigma_{xa}^2 + p^2 R_{ab(A)}^2 (N_A)^2 \left(\frac{1-f_A}{n_A} \right) w_{ab(A)} \sigma_{xab(A)}^2 \\ &- 2R_a^2 \left(\frac{N_a}{n_a} \right)^2 n_A (1-f_A) \rho_{xya} \sigma_{ya} \sigma_{xa} \\ &- 2p^2 R_{ab(A)}^2 \left(\frac{N_{ab}}{n_{ab(A)}} \right)^2 n_A (1-f_A) \rho_{xyab(A)} \sigma_{yab(A)} \sigma_{xab(A)} \\ &+ (N^B)^2 \left(\frac{1-f_B}{n_B} \right) w_b \sigma_{yb}^2 + (1-p)^2 (N^B)^2 \left(\frac{1-f_B}{n_B} \right) w_{ab(B)} \sigma_{yab(B)}^2 \\ &+ R_b^2 (N^B)^2 \left(\frac{1-f_B}{n_B} \right) w_b \sigma_{xb}^2 + (1-p)^2 R_{ab(B)}^2 (N_A)^2 \left(\frac{1-f_B}{n_B} \right) w_{ab(B)} \sigma_{xab(B)}^2 \\ &- 2R_b^2 \left(\frac{N_b}{n_b} \right)^2 n_B (1-f_B) \rho_{xyb} \sigma_{yb} \sigma_{xb} \\ &- 2(1-p)^2 R_{ab(B)}^2 \left(\frac{N_{ab}}{n_{ab(B)}} \right)^2 n_A (1-f_A) \rho_{xyab(B)} \sigma_{yab(B)} \sigma_{xab(B)}, \end{aligned}$$

onde

- $w_{ab(A)} = N_{ab}/N_A$, $w_{ab(B)} = N_{ab}/N_B$;
- $R_a = \frac{\mu_{ya}}{\mu_{xa}}$, $R_{ab(A)} = \frac{\mu_{yab(A)}}{\mu_{xab(A)}}$;
- $R_b = \frac{\mu_{yb}}{\mu_{xb}}$, $R_{ab(B)} = \frac{\mu_{yab(B)}}{\mu_{xab(B)}}$;
- σ_{ya}^2 , σ_{yb}^2 , $\sigma_{yab(A)}^2$, $\sigma_{yab(B)}^2$ são as variâncias populacionais referentes aos domínios;
- ρ_{xya} , $\rho_{xyab(A)}$, ρ_{xyb} e $\rho_{xyab(B)}$ representam os coeficientes de correlação linear entre as variáveis auxiliares e as variáveis de interesse nos domínios.

Ao substituir as quantidades populacionais pelas suas respectivas quantidades amostrais, tem-se um estimador para a variância do estimador $RH1$. Na situação em que apenas um cadastro contém informação auxiliar disponível, A por exemplo, este estimador reduz-se à forma apresentada e analisada por Coelho (2007), onde

$$\bar{y}_{RH1} = w_a \frac{\tilde{y}_a}{\tilde{x}_a} \mu_{X_a} + w_{ab} p \frac{\tilde{y}_{ab(A)}}{\tilde{x}_{ab(A)}} \mu_{X_{ab}} + w_b \tilde{y}_b + w_{ab} (1 - p) \tilde{y}_{ab(B)}. \quad (5.2)$$

1.2 Estimador razão RH2

O estimador $RH2$ do tipo razão para μ , sob a estratégia de Hartley, é considerado para a situação em que existe informação apenas para o total da variável auxiliar nos cadastros. Considerando novamente o caso em que um plano de amostragem aleatória simples é utilizado, tem-se que

$$\bar{y}_{RH2} = \left(\frac{w_a \tilde{y}_a + p w_{ab} \tilde{y}_{ab(A)}}{w_a \tilde{x}_a + p w_{ab} \tilde{x}_{ab(A)}} \right) \mu_{X_A} + \left(\frac{w_b \tilde{y}_b + (1 - p) w_{ab} \tilde{y}_{ab(B)}}{w_b \tilde{x}_b + (1 - p) w_{ab} \tilde{x}_{ab(B)}} \right) \mu_{X_B}, \quad (5.3)$$

onde w_a , w_b , w_{ab} , \tilde{y}_a , \tilde{y}_b , $\tilde{y}_{ab(A)}$, $\tilde{y}_{ab(B)}$, \tilde{x}_a , \tilde{x}_b , $\tilde{x}_{ab(A)}$ e $\tilde{x}_{ab(B)}$ são definidos de maneira análoga à apresentada para o estimador $RH1$. μ_{X_A} e μ_{X_B} são as médias populacionais da variável auxiliar nos cadastros A e B . A variância aproximada deste estimador, expressa em termos matriciais, é dada por

$$AVar(\bar{y}_{RH2}) = \mathbf{d}^T \Sigma_{RH2} \mathbf{d},$$

onde

$$\mathbf{d}^T = \left[\frac{\mu_{X_A}}{\mu_{x(A)}^*}, \frac{p \mu_{X_A}}{\mu_{x(A)}^*}, -\frac{G_A \mu_{X_A}}{\mu_{x(A)}^*}, -\frac{p G_A \mu_{X_A}}{\mu_{x(A)}^*}, \frac{\mu_{X_B}}{\mu_{x(B)}^*}, \frac{(1 - p) \mu_{X_B}}{\mu_{x(B)}^*}, -\frac{G_B \mu_{X_B}}{\mu_{x(B)}^*}, -\frac{G_B (1 - p) \mu_{X_B}}{\mu_{x(B)}^*} \right]$$

Na expressão da variância aproximada, $\Sigma_{RH2} = \Sigma_{RH1}$. Tem-se ainda que

$$\begin{aligned} \mu_{y(A)}^* &= \frac{N_a \mu_{ya} + p N_{ab} \mu_{yab(A)}}{N_A} & \mu_{x(A)}^* &= \frac{N_a \mu_{xa} + p N_{ab} \mu_{xab(A)}}{N_A} \\ \mu_{y(B)}^* &= \frac{N_b \mu_{yb} + p N_{ab} \mu_{yab(B)}}{N_B} & \mu_{x(B)}^* &= \frac{N_b \mu_{xb} + p N_{ab} \mu_{xab(B)}}{N_B} \\ G_A &= \frac{\mu_{y(A)}^*}{\mu_{x(A)}^*} & G_B &= \frac{\mu_{y(B)}^*}{\mu_{x(B)}^*} \end{aligned}$$

Ao substituir as quantidades populacionais pelas suas respectivas quantidades amostrais, tem-se um estimador para a variância do estimador $RH2$. Considerando novamente a situação em que apenas o cadastro A possui informação auxiliar disponível, este estimador reduz-se à outra forma de estimador razão apresentada por Coelho (2007), onde

$$\bar{y}_{RH2} = \left(\frac{w_a \tilde{y}_a + p w_{ab} \tilde{y}_{ab(A)}}{w_a \tilde{x}_a + p w_{ab} \tilde{x}_{ab(A)}} \right) \mu_{X_A} + w_b \tilde{y}_b + (1 - p) w_{ab} \tilde{y}_{ab(B)}, \quad (5.4)$$

2 Estimadores do tipo razão sob a estratégia de Fuller & Burmeister

Os estimadores do tipo razão sob a estratégia de Hartley são propostos para o caso em que N_a , N_b e N_{ab} são conhecidos. Quando isto não acontecer, e admitindo o uso de um plano de amostragem aleatória simples nos cadastros, e considerando o cenário em que o estimador regressão sob a abordagem de Fuller & Burmeister reduz-se ao estimador do tipo razão, é possível obter os estimadores apresentados a seguir.

2.1 Estimador razão FB1

O estimador $RFB1$ do tipo razão para μ , sob a estratégia de Fuller & Burmeister, é dado por

$$\bar{y}_{RFB1} = \hat{w}_a \frac{\tilde{y}_a}{\tilde{x}_a} \mu_{X_a} + \hat{w}_b \frac{\tilde{y}_b}{\tilde{x}_b} \mu_{X_b} + \hat{w}_{ab} \beta_1 \frac{\tilde{y}_{ab(A)}}{\tilde{x}_{ab(A)}} \mu_{X_{ab}} + \hat{w}_{ab} (1 - \beta_1) \frac{\tilde{y}_{ab(B)}}{\tilde{x}_{ab(B)}} \mu_{X_{ab}}, \quad (5.5)$$

onde

- $\hat{w}_a = (N_A - \hat{N}_{ab,s})/N = \hat{N}_a/N$, $\hat{w}_b = (N_B - \hat{N}_{ab,s})/N = \hat{N}_b/N$ e $\hat{w}_{ab} = \hat{N}_{ab,s}/N$ são estimadores de w_a , w_b e w_{ab} . $\hat{N}_{ab,s}$ é definido como a menor raiz de (2.13), e os demais termos são definidos de maneira análoga à apresentada para o estimador razão sob a estratégia de Hartley.

Aplicando o método de linearização de Taylor em cada um dos termos do estimador apresentado em (5.5), é possível obter um pseudo-estimador de \bar{y}_{RFB1} , para que seja possível a obtenção de uma variância aproximada do estimador. Tem-se

$$A\text{Var}(\bar{y}_{RFB1}) = N^{-2} \mathbf{d}^T \Sigma_{RFB1} \mathbf{d} + O(1), \tag{5.6}$$

onde

$$\mathbf{d} = \begin{bmatrix} 1 \\ \beta_1 \\ -R_a \\ -\beta_1 R_{ab(A)} \\ 1 \\ 1 - \beta_1 \\ -R_b \\ -(1 - \beta_1) R_{ab(B)} \\ 1 \\ 1 \end{bmatrix}$$

com $R_a = \mu_{ya} / \mu_{xa}$, $R_a = \mu_{yab(A)} / \mu_{xab(A)}$, $R_b = \mu_{yb} / \mu_{xb}$ e $R_{ab(B)} = \mu_{yab(B)} / \mu_{xab(B)}$. A matriz de variâncias e covariâncias do estimador é dada por

$$\Sigma_{RFB1} = \begin{pmatrix} \Sigma_{y(a;ab(A))} & \Sigma_{xy(a;ab(A))} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \Sigma_{xy(a;ab(A))} & \Sigma_{x(a;ab(A))} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \Sigma_{y(b;ab(B))} & \Sigma_{xy(b;ab(B))} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \Sigma_{xy(b;ab(B))} & \Sigma_{x(b;ab(B))} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \Sigma_{N_{(A,B)}^{ab}} \end{pmatrix}$$

as submatrizes de Σ_{RFB1} com exceção de $\Sigma_{N_{(A,B)}^{ab}}$ são iguais às submatrizes de Σ_{RH1} . Tem-se ainda que

$$\Sigma_{N_{(A,B)}^{ab}} = \begin{pmatrix} \text{Var}(\hat{N}_{ab}^A) & 0 \\ 0 & \text{Var}(\hat{N}_{ab}^B) \end{pmatrix}$$

Ao substituir as quantidades populacionais pelas suas respectivas quantidades amostrais, tem-se um estimador para a variância do estimador $RFB1$. Considerando novamente a situação em

que apenas o cadastro A possui informação auxiliar disponível, este estimador reduz-se à outra forma de estimador razão apresentada e analisada por Coelho (2007), onde

$$\bar{y}_{RFB1} = \hat{w}_a \frac{\tilde{y}_a}{\tilde{x}_a} \mu_{X_a} + \hat{w}_{ab} \beta_1 \frac{\tilde{y}_{ab(A)}}{\tilde{x}_{ab(A)}} \mu_{X_{ab}} + \hat{w}_b \tilde{y}_b + \hat{w}_{ab} (1 - \beta_1) \tilde{y}_{ab(B)}. \quad (5.7)$$

2.2 Estimador razão RFB2

O estimador $RFB2$ do tipo razão para μ , sob a estratégia de Fuller & Burmeister, é considerado para a situação em que existe informação apenas para o total da variável auxiliar nos cadastros. Considerando-se novamente o caso em que um plano de amostragem aleatória simples é utilizado em ambos os cadastros, tem-se que

$$\bar{y}_{RFB2} = \left(\frac{\hat{w}_a \tilde{y}_a + \beta_1 \hat{w}_{ab} \tilde{y}_{ab(A)}}{\hat{w}_a \tilde{x}_a + \beta_1 \hat{w}_{ab} \tilde{x}_{ab(A)}} \right) \mu_{X_A} + \left(\frac{\hat{w}_b \tilde{y}_b + (1 - \beta_1) \hat{w}_{ab} \tilde{y}_{ab(B)}}{\hat{w}_b \tilde{x}_b + (1 - \beta_1) \hat{w}_{ab} \tilde{x}_{ab(A)}} \right) \mu_{X_B}, \quad (5.8)$$

onde μ_{X_A} e μ_{X_B} são as médias populacionais para as únicas variáveis auxiliares nos cadastros A e B , e as demais quantidades são definidas de forma análoga às quantidades apresentadas para o estimador $RFB1$. Sob a situação em que apenas uma variável está disponível para um dos cadastro, A por exemplo, tem-se que

$$\bar{y}_{RFB2} = \left(\frac{\hat{w}_a \tilde{y}_a + \beta_1 \hat{w}_{ab} \tilde{y}_{ab(A)}}{\hat{w}_a \tilde{x}_a + \beta_1 \hat{w}_{ab} \tilde{x}_{ab(A)}} \right) \mu_{X_A} + \hat{w}_b \tilde{y}_b + \hat{w}_{ab} (1 - p) \tilde{y}_{ab(B)}, \quad (5.9)$$

com variância aproximada igual a

$$A\text{Var}(\bar{y}_{RFB2}) = N_v^{-2} [t_{X_A}^2 A\text{Var}(\hat{G}_A) + A\text{Var}(\hat{G}_B)],$$

onde t_{X_A} representa o total populacional da variável auxiliar no cadastro A . Tem-se que

$$\begin{aligned} A\text{Var}(\hat{G}_A) &= \left(\frac{1}{t_{X_A}^*} \right)^2 \left\{ \left(\mu_{y_a}^2 + \beta_1^2 \mu_{y_{ab(A)}}^2 + G_A^2 \mu_{x_a}^2 + p^2 G_A^2 \mu_{x_{ab(A)}}^2 \right) \frac{N_a N_b N_{ab} g_A g_B}{n_A N_b g_B + n_B N_a g_A} \right. \\ &+ N_a^2 \left(\frac{1 - f_A}{n_A P_a} \right) \sigma_{y_a}^2 + \beta_1^2 N_{ab}^2 \left(\frac{1 - f_{\mathcal{A}}}{n_{\mathcal{A}} P_{ab(A)}} \right) \sigma_{y_{ab(A)}}^2 + G_A^2 N_a^2 \left(\frac{1 - f_A}{n_A P_a} \right) \sigma_{x_a}^2 \\ &+ \beta_1^2 G_A^2 N_{ab}^2 \left(\frac{1 - f_A}{n_A P_{ab(A)}} \right) \sigma_{x_{ab(A)}}^2 - 2 G_A N_a^2 \frac{n_A}{n_a^2} \left(1 - \frac{n_A}{N_A} \right) \rho_{x_{y_a}} \sigma_{x_a} \sigma_{y_a} \end{aligned}$$

$$- 2\beta_1^2 G_A N_{ab}^2 \frac{n_{\mathcal{A}}}{n_{ab(A)}^2} \left(1 - \frac{n_{\mathcal{A}}}{N_{\mathcal{A}}}\right) \rho_{xyab(A)} \sigma_{xab(A)} \sigma_{yab(A)} \Big\},$$

onde $G_{\mathcal{A}} = \frac{t_{y,\mathcal{A}}^*}{t_{x,\mathcal{A}}^*}$, com $t_{y,\mathcal{A}}^* = N_a \mu_{y_a} + p N_{ab} \mu_{yab(A)}$, $t_{x,\mathcal{A}}^* = N_a \mu_{x_a} + p N_{ab} \mu_{xab(A)}$ e $\sigma_{x_a}^2$, $\sigma_{xab(A)}^2$, $\sigma_{y_a}^2$ e $\sigma_{yab(A)}^2$ representam as variâncias das variáveis nos domínios. Tem-se ainda que

$$\begin{aligned} A\text{Var}(\hat{G}_B) &= G_B + \left\{ (1 - \beta_1)^2 (\mu_{yab(B)}^2) + (\mu_{yb}^2) \right\} \frac{N_a N_b N_{ab} g_A g_B}{n_A N_b g_B + n_B N_a g_A} \\ &+ (1 - \beta_1)^2 N_{ab}^2 \left(\frac{1 - f_B}{n_B P_{ab(B)}} \right) \sigma_{yab(B)}^2 + N_b^2 \left(\frac{1 - f_B}{n_B P_b} \right) \sigma_{yb}^2 \end{aligned}$$

onde $G_B = N_b \mu_{y_b} + (1 - p) N_{ab} \mu_{yab(B)}$, $P_a = n_a / n_A$, $P_{ab(A)} = n_a b(A) / n_A$, $P_b = n_b / n_B$ e $P_{ab(B)} = n_a b(B) / n_B$. A derivação da variância do estimador *RFB2*, bem como a obtenção do valor de β_1 que minimiza sua variância pode ser consultada em Coelho (2007).

3 Estimadores do tipo razão sob a estratégia BLG

Sob a estratégia adotada por Bankier, Lepkowski & Groves (1986), e considerando que existe informação de uma variável auxiliar em ambos os cadastros, o estimador *RBLG* do tipo razão para μ , sob a estratégia BLG é dado por

$$\bar{y}_{RBLG} = \frac{\frac{\sum_{k \in S_A} \omega_k y_k}{\sum_{k \in S_A} \omega_k x_k} t_{X_A} + \frac{\sum_{k \in S_B} \omega_k y_k}{\sum_{k \in S_B} \omega_k x_k} t_{X_B}}{N} \quad (5.10)$$

em que \tilde{y}_a , \tilde{x}_a , \tilde{y}_b e \tilde{x}_b representam as médias amostrais das variáveis de interesse e das variáveis auxiliares em cada cadastro, respectivamente. Através do método de linearização de Taylor, tem-se que a variância aproximada do estimador de Bankier é dada por

$$\begin{aligned} A\text{Var}(\bar{y}_{RBLG}) &= (w_A)^2 \text{Var}(\tilde{y}_a) + \left(w_A \frac{\mu_{Y_A}}{\mu_{X_A}} \right)^2 \text{Var}(\tilde{x}_a) \\ &+ (w_B)^2 \text{Var}(\tilde{y}_b) + \left(w_B \frac{\mu_{Y_B}}{\mu_{X_B}} \right)^2 \text{Var}(\tilde{x}_b) \\ &- 2(w_A)^2 \frac{\mu_{Y_A}}{\mu_{X_A}} \text{Cov}(\tilde{y}_a, \tilde{x}_a) - 2(w_B)^2 \frac{\mu_{Y_B}}{\mu_{X_B}} \text{Cov}(\tilde{y}_b, \tilde{x}_b). \end{aligned} \quad (5.11)$$

Em termos matriciais, a variância aproximada do estimador RBLG é dada por

$$AVar(\bar{y}_{RBLG}) = \mathbf{d}^T \Sigma_{RBLG} \mathbf{d},$$

em que

$$\mathbf{d} = \begin{bmatrix} w_A \\ w_A \frac{\mu_{Y_A}}{\mu_{X_A}} \\ w_B \\ w_B \frac{\mu_{Y_B}}{\mu_{X_B}} \end{bmatrix} \quad \text{e} \quad \Sigma_{RBLG} = \begin{pmatrix} \sigma_{yA}^2 & \sigma_{xyA} & 0 & 0 \\ \sigma_{xyA} & \sigma_{xA}^2 & 0 & 0 \\ 0 & 0 & \sigma_{yB}^2 & \sigma_{xyB} \\ 0 & 0 & \sigma_{xyB} & \sigma_{xB}^2 \end{pmatrix}$$

Da mesma forma que a vista para o estimador regressão, é possível propor uma versão *raking* para este estimador, da seguinte forma:

$$\begin{aligned} \bar{y}_{RBLG} &= \frac{\frac{\sum_{k \in S_A} \omega_k^{(r)} y_k}{\sum_{k \in S_A} \omega_k^{(r)} x_k} t_{X_A} + \frac{\sum_{k \in S_B} \omega_k^{(r)} y_k}{\sum_{k \in S_B} \omega_k^{(r)} x_k} t_{X_B}}{N} \\ &= w_A \frac{\tilde{y}_{(r)A}}{\tilde{x}_{(r)A}} \mu_{X_A} + w_B \frac{\tilde{y}_{(r)B}}{\tilde{x}_{(r)B}} \mu_{X_B}, \end{aligned} \quad (5.12)$$

em que a variância aproximada é similar à apresentada em (5.11).

4 Estimadores do tipo razão sob a estratégia de Máxima-Pseudo Verossimilhança

Sob esta estratégia, a proposta de estimador razão se torna simples. Sob a situação em que é de interesse a estimação de uma média populacional sob a abordagem de cadastro duplo, tem-se que

$$\bar{y}_{RPML} = N^{-1} \left\{ \frac{(N_A - \hat{N}_{ab,PML}) \tilde{y}_a}{(N_A - \hat{N}_{ab,PML}) \tilde{x}_a} t_{X_a} + \frac{\hat{N}_{ab,PML} \tilde{y}_{ab(*)}}{\hat{N}_{ab,PML} \tilde{x}_{ab(*)}} t_{X_{ab}} + \frac{(N_B - \hat{N}_{ab,PML}) \tilde{y}_b}{(N_B - \hat{N}_{ab,PML}) \tilde{x}_b} t_{X_b} \right\}$$

$$= N^{-1} \left\{ \frac{\tilde{y}_a}{\tilde{x}_a} t_{X_a} + \frac{\tilde{y}_{ab(*)}}{\tilde{x}_{ab(*)}} t_{X_{ab}} + \frac{\tilde{y}_b}{\tilde{x}_b} t_{X_b} \right\}$$

Como visto anteriormente, através do método de linearização de Taylor, um pseudo-estimador para $\tilde{y}_{ab(*)}$ é dado por

$$\begin{aligned} \tilde{y}_{ab(*)} &\approx \left(\frac{\pi^A}{\pi^A + \pi^B} \right) \tilde{y}_{ab(A)} + \left(\frac{\pi^B}{\pi^A + \pi^B} \right) \tilde{y}_{ab(B)} \\ &= p_A \tilde{y}_{ab(A)} + p_B \tilde{y}_{ab(B)} = p_A \tilde{y}_{ab(A)} + (1 - p_A) \tilde{y}_{ab(B)} = p_A (\tilde{y}_{ab(A)} - \tilde{y}_{ab(B)}) + \tilde{y}_{ab(B)} = g(\tilde{y}_{ab(A)}, \tilde{y}_{ab(B)}) \end{aligned}$$

Dessa forma, tem-se então que

$$\frac{g(\tilde{y}_{ab(A)}, \tilde{y}_{ab(B)})}{g(\tilde{x}_{ab(A)}, \tilde{x}_{ab(B)})} \approx \frac{g(\mu_{yab}, \mu_{yab})}{g(\mu_{xab}, \mu_{xab})} + \frac{1}{g(\mu_{xab}, \mu_{xab})} \left(g(\tilde{y}_{ab(A)}, \tilde{y}_{ab(B)}) - \frac{g(\mu_{yab}, \mu_{yab})}{g(\mu_{xab}, \mu_{xab})} g(\tilde{x}_{ab(A)}, \tilde{x}_{ab(B)}) \right)$$

É possível notar que $g(\mu_{yab}, \mu_{yab}) = \mu_{yab}$ e $g(\mu_{xab}, \mu_{xab}) = \mu_{xab}$. Dessa forma,

$$\frac{\tilde{y}_{ab(*)}}{\tilde{x}_{ab(*)}} = \frac{g(\tilde{y}_{ab(A)}, \tilde{y}_{ab(B)})}{g(\tilde{x}_{ab(A)}, \tilde{x}_{ab(B)})} \approx \frac{\mu_{yab}}{\mu_{xab}} + \frac{1}{\mu_{xab}} \left(g(\tilde{y}_{ab(A)}, \tilde{y}_{ab(B)}) - \frac{\mu_{yab}}{\mu_{xab}} g(\tilde{x}_{ab(A)}, \tilde{x}_{ab(B)}) \right)$$

Logo,

$$\begin{aligned} \bar{y}_{RPML} &\approx N^{-1} \left\{ \left[\frac{\mu_{ya}}{\mu_{xa}} + \frac{1}{\mu_{xa}} \left(\tilde{y}_a - \frac{\mu_{ya}}{\mu_{xa}} \tilde{x}_a \right) \right] t_{X_a} + \left[\frac{\mu_{yb}}{\mu_{xb}} + \frac{1}{\mu_{xb}} \left(\tilde{y}_b - \frac{\mu_{yb}}{\mu_{xb}} \tilde{x}_b \right) \right] t_{X_b} \right. \\ &\quad \left. + \left[\frac{\mu_{yab}}{\mu_{xab}} + \frac{1}{\mu_{xab}} \left(g(\tilde{y}_{ab(A)}, \tilde{y}_{ab(B)}) - \frac{\mu_{yab}}{\mu_{xab}} g(\tilde{x}_{ab(A)}, \tilde{x}_{ab(B)}) \right) \right] t_{X_{ab}} \right\} \\ AVar(\bar{y}_{RPML}) &= \left(\frac{t_{X_a}}{\mu_{xa}} \right)^2 \text{Var}(\tilde{y}_a) + \left(\frac{\mu_{ya}}{\mu_{xa}} t_{X_a} \right)^2 \text{Var}(\tilde{x}_a) - 2\mu_{ya} \left(\frac{t_{X_a}}{\mu_{xa}} \right)^2 \text{Cov}(\tilde{y}_a; \tilde{x}_a) \\ &\quad + \left(\frac{t_{X_b}}{\mu_{xb}} \right)^2 \text{Var}(\tilde{y}_b) + \left(\frac{\mu_{yb}}{\mu_{xb}} t_{X_b} \right)^2 \text{Var}(\tilde{x}_b) - 2\mu_{yb} \left(\frac{t_{X_b}}{\mu_{xb}} \right)^2 \text{Cov}(\tilde{y}_b; \tilde{x}_b) \\ &\quad + p_A^2 \left(\frac{t_{X_{ab}}}{\mu_{xab}} \right)^2 \text{Var}(\tilde{y}_{ab(A)}) + p_A^2 \left(\frac{\mu_{yab}}{\mu_{xab}} \right)^2 \text{Var}(\tilde{x}_{ab(A)}) - 2p_A^2 \mu_{yab} \left(\frac{t_{X_{ab}}}{\mu_{xab}} \right)^2 \text{Cov}(\tilde{y}_{ab(A)}; \tilde{x}_{ab(A)}) \end{aligned}$$

$$\begin{aligned}
 &+ p_B^2 \left(\frac{t_{X_{ab}}}{\mu_{xab}} \right)^2 \text{Var}(\tilde{y}_{ab(B)}) + p_B^2 \left(\frac{\mu_{yab}}{\mu_{xab}} \right)^2 \text{Var}(\tilde{x}_{ab(B)}) - 2p_B^2 \mu_{yab} \left(\frac{t_{X_{ab}}}{\mu_{xab}} \right)^2 \text{Cov}(\tilde{y}_{ab(B)}; \tilde{x}_{ab(B)}) \\
 &= \mathbf{d}^T \Sigma_{RPML} \mathbf{d},
 \end{aligned}$$

em que \mathbf{d}^T e Σ_{RPML} são dados por

$$\mathbf{d}^T = [\mathbf{d}_1^T \quad \mathbf{d}_2^T]$$

$$\mathbf{d}_1^T = \left[\begin{array}{cccccc} \frac{t_{X_a}}{\mu_{xa}} & \frac{\mu_{ya}}{\mu_{xa}} t_{X_a} & \sqrt{\mu_{ya}} \frac{t_{X_a}}{\mu_{xa}} & p_A \frac{t_{X_{ab}}}{\mu_{xab}} & p_A \frac{\mu_{yab}}{\mu_{xab}} & p_A \sqrt{\mu_{yab}} \frac{t_{X_{ab}}}{\mu_{xab}} \end{array} \right]$$

$$\mathbf{d}_2^T = \left[\begin{array}{cccccc} \frac{t_{X_b}}{\mu_x^b} & \frac{\mu_{yb}}{\mu_{xb}} t_{X_b} & \sqrt{\mu_y^b} \frac{t_{X_b}}{\mu_x^b} & p_B \frac{t_{X_{ab}}}{\mu_{xab}} & p_B \frac{\mu_{yab}}{\mu_{xab}} & p_B \sqrt{\mu_{yab}} \frac{t_{X_{ab}}}{\mu_{xab}} \end{array} \right]$$

$$\Sigma_{RPML} = \begin{pmatrix} \Sigma_A & \mathbf{0} \\ \mathbf{0} & \Sigma_B \end{pmatrix}$$

em que

$$\Sigma_A = \begin{pmatrix} \sigma_{ya}^2 & \sigma_{y(a;ab(A))} & \sigma_{xya} & \sigma_{xy(a;ab(A))} \\ \sigma_{y(a;ab(A))} & \sigma_{yab}^2 & \sigma_{xy(a;ab(A))} & \sigma_{xy(a;ab(A))} \\ \sigma_{xya} & \sigma_{xy(a;ab(A))} & \sigma_{xa}^2 & \sigma_{x(a;ab(A))} \\ \sigma_{xy(a;ab(A))} & \sigma_{xa}^2 & \sigma_{x(a;ab(A))} & \sigma_{xab(A)}^2 \end{pmatrix}$$

$$\Sigma_B = \begin{pmatrix} \sigma_{yb}^2 & \sigma_{y(b;ab(B))} & \sigma_{xyb} & \sigma_{xy(b;ab(B))} \\ \sigma_{y(b;ab(B))} & \sigma_{yab}^2 & \sigma_{xy(b;ab(B))} & \sigma_{xy(b;ab(B))} \\ \sigma_{xyb} & \sigma_{xy(b;ab(B))} & \sigma_{xb}^2 & \sigma_{x(b;ab(B))} \\ \sigma_{xy(b;ab(B))} & \sigma_{xb}^2 & \sigma_{x(b;ab(B))} & \sigma_{xab(B)}^2 \end{pmatrix}$$

Sob um plano de amostragem aleatória simples, é possível observar que a variância aproximada da estratégia PML coincide com a variância do estimador razão $R1$ de Hartley. Como visto, a informação de uma variável auxiliar pode ser alocada no processo de estimação sob cada estratégia proposta. Qual das estratégias será a melhor vai depender do cenário em questão. Se os cadastros estiverem sob o cenário 2 (N_A , N_B , N_a , N_b e N_{ab} conhecidos), os estimadores de Hartley e Bankier podem ser utilizados, sendo a estratégia de Bankier uma estratégia mais simples de implementação, caso exista grande quantidade de informação no

domínio ab . Já se o interesse for obter uma maior precisão, a estratégia de Hartley deve ser considerada, uma vez que a constante p é escolhida de modo a minimizar a variância do estimador utilizado. Caso a situação seja a do cenário 3 (N_A , N_B , N_a , N_b e N_{ab} desconhecidos), os estimadores de Fuller & Burmeister e PML podem ser utilizados, sendo a estratégia de Fuller & Burmeister mais eficiente que a PML pela mesma justificativa apresentada para a estratégia de Hartley. Para ambos os cenários, as estratégias de Bankier e PML surgem como alternativas mais direcionadas à agilidade de obtenção de resultados do que eficiência.

1 Discussão metodológica

Os desempenhos dos estimadores propostos nesta tese foram avaliados através do método de simulação de Monte Carlo. Admitindo-se que em ambos os cadastros exista informação de duas variáveis auxiliares disponíveis, o cenário considerado foi que a relação entre a variável de interesse e as variáveis auxiliares é descrita pelo modelo

$$\xi : y_k = \beta_0 + \beta_1 x_{k1} + \cdots + \beta_q x_{kq} + \varepsilon_k = \beta_0 + \sum_{r=1}^q x_{kr} + \varepsilon_k = \mathbf{x}_k^T \underline{\beta} + \varepsilon_k, \quad (6.1)$$

onde foi fixado $q = 2$ em ambos os cadastros. Os estimadores apresentados nos capítulos 2 e os estimadores regressão propostos no capítulo 3 foram avaliados de acordo com o plano de amostragem aleatória simples aplicado em cada cadastro.

Foi admitido que $N_A + N_B - N_{ab} = 4500$, conforme descrição na tabela (6.1) a seguir. Para cada tamanho populacional, foram consideradas amostras de tamanho $n_A = n_B = 125, 250, 500, 750, 1000, 1250$.

Tabela 6.1: Tamanhos populacionais utilizados para a , b e ab

POPULAÇÃO	N_a	N_b	N_{ab}	$N_A = N_B$
1	2000	2000	500	2500
2	1750	1750	1000	2750
3	1500	1500	1500	3000

Para a avaliação dos estimadores foram considerados o viés relativo, desvio padrão e

erro quadrático médio dos estimadores para a média populacional. O viés de um estimador $\hat{\theta}$ é dado por

$$B(\hat{\theta}) = E(\hat{\theta} - \theta) = E(\hat{\theta}) - \theta$$

e o erro quadrático médio por

$$EQM(\hat{\theta}) = \text{Var}(\hat{\theta}) + \{B(\hat{\theta})\}^2.$$

O viés relativo é dado por

$$VR(\hat{\theta}) = \frac{E(\hat{\theta} - \theta)}{\theta}$$

onde θ é o parâmetro de interesse e $\hat{\theta}$ é estimador de θ . O esquema de geração de dados foi conduzido da seguinte forma:

Passo 1. Gerar pares (i, j) referentes às observações das variáveis auxiliares para o domínio a , onde $i, j = 1, \dots, N_a$, considerando uma distribuição $\mathcal{N}(\mu_{xaq}, \sigma_{xaq}^2)$, com $q = 1, 2$ e $\mu_{ax1} = 4$ e $\mu_{ax2} = 8$. Foram ainda considerados dois cenários para a variância das variáveis auxiliares:

1. $\sigma_{xa1}^2 = 1$ e $\sigma_{xa2}^2 = 1$;
2. $\sigma_{xa1}^2 = 10$ e $\sigma_{xa2}^2 = 4$.

Passo 2. Gerar pares (i, j) para as variáveis auxiliares no domínio ab , onde $i, j = 1, \dots, N_{ab}$, a partir de valores de uma distribuição $\mathcal{N}(\mu_{xabq}, \sigma_{xabq}^2)$, $q=1,2$. Analogamente ao item anterior, tem-se que $\mu_{xab1} = 2$, $\mu_{xab2} = 1$, e dois cenários para as variâncias das variáveis auxiliares neste domínio:

1. $\sigma_{xab1}^2 = 1$ e $\sigma_{xab2}^2 = 1$;
2. $\sigma_{xab1}^2 = 10$ e $\sigma_{xab2}^2 = 4$.

Passo 3. Gerar N_a observações da variável de interesse y para o domínio a , onde $i = 1, \dots, N_a$, seguindo distribuição $\mathcal{N}(\beta_0 + \beta_1^a x_{1kA} + \beta_2^a x_{2kA}, \sigma_k^2)$, onde $\beta_i^a = 1$ e $\sigma_k^2 = \sigma_{ya}^2 = 1$;

Passo 4. Gerar N_{ab} observações da variável de interesse para o domínio ab , a partir de uma distribuição $\mathcal{N}(\beta_0 + \beta_1^{ab} x_{1kA} + \beta_2^{ab} x_{2kA}, \sigma_k^2)$, onde $\beta_i^{ab} = 1$ e $\sigma_k^2 = \sigma_{yab}^2 = 1$.

Passo 5. Repetir os passos 1, 2, 3 e 4 para o domínio b e cadastro B . No passos 1, os cenários para as variâncias foram os seguintes:

1. $\sigma_{xb1}^2 = 1$, $\sigma_{xb2}^2 = 1$;

$$2. \sigma_{xb1}^2 = 8 \text{ e } \sigma_{xb2}^2 = 12.$$

Passo 6. Na estrutura de populações geradas nos passos de 1 a 5 para os cadastros A e B , coletar uma amostra de tamanho especificado através do plano de amostragem aleatória simples e obter estimativas para a média populacional utilizando os estimadores dos capítulos 2 e 3;

Passo 7. Obter estimativas de desvio padrão (DP), Viés Relativo (VR), Erro Quadrático Médio (EQM). Além disso, para cada estimador $\hat{\theta}$ é construído um intervalo do tipo

$$\hat{\theta} \pm 1.96\sqrt{\text{Var}(\hat{\theta})}.$$

Este intervalo foi considerado para cálculo da proporção de vezes em que o intervalo de confiança dos estimadores propostos cobre o verdadeiro valor do parâmetro para cada caso. Esta proporção será chamada de *Proporção de Cobertura do Intervalo* (PCI).

As populações geradas nos passos de 1 a 5 permaneceram fixas no estudo de simulação. Utilizando o método de Monte Carlo, os passos 6 e 7 foram repetidos $r = 1000$ vezes, para cada tamanho de amostra especificado. Os resultados são apresentados nas tabelas 6.2 a 6.13 e pelas figuras 6.1 a 6.15. Além dos estimadores para a variável de interesse, foram considerados também no estudo de simulação os estimadores para N_{ab} sob as estratégias de Fuller & Burmeister (1972) e Skinner & Rao (1996), denotados por $\hat{N}_{ab,s}$ e $\hat{N}_{ab,PML}$.

2 Resultados da Simulação

2.1 Desempenho dos estimadores propostos sob o plano de amostragem aleatória simples

O desempenho dos estimadores propostos sob o plano de amostragem aleatória simples, sob os cenários 2 e 3, é apresentado a seguir pelas tabelas (6.2) a (6.13) e figuras (6.1) a (6.6).

De modo geral, é possível observar que os estimadores já propostos na literatura tiveram desempenho similar tanto em relação à proporção de cobertura do intervalo de confiança quanto na redução da variância à medida em que o tamanho de amostra nos cadastros aumenta. Observando-se os cenários tem-se que, sob o cenário 2, o estimador de Hartley apresentou evidência de melhor desempenho em relação ao estimador BLG, como mostram os valores de desvio padrão e EQM. Sob o cenário 3, tem-se evidência de melhor desempenho do estimador PML em relação ao estimador de Fuller & Burmeister, quando observados também os valores de desvio padrão e EQM. Isto foi verificado para todos os valores populacionais de N_{ab} , o que fornece evidência de que em situações onde a variabilidade é idêntica para

todos os cadastros, basta que o cenário em questão seja determinado e utilizado o estimador apropriado para este cenário.

Comparados aos estimadores de Hartley, Fuller, Bankier e PML quando o tamanho de amostra nos cadastros aumenta, as versões propostas no contexto de estimadores regressão apresentaram considerável redução em relação à variância e erro quadrático médio, apresentando um viés relativo maior. Além disso, há evidência de que os estimadores regressão propostos são mais confiáveis para se estimar o parâmetro populacional de interesse, como mostra as PCI's para cada estimador, bem próximas dos níveis nominais de um intervalo de confiança calculado com $z_{1-\alpha/2} = 1.96$. Este desempenho foi o mesmo para todos os valores de N_{ab} apresentados.

Os resultados mostraram também que cada um dos estimadores regressão propostos apresentou viés relativo maior que sua respectiva versão simples apresentada no capítulo 2. Porém, os resultados forneceram evidência de que esse viés tende a diminuir com o aumento do tamanho da amostra em ambos os cadastros.

Na situação em que as variâncias em domínios e cadastros são diferentes, os estimadores regressão considerados sob as estratégias de Hartley (cenário 2) e Fuller & Burmeister (cenário 3) ainda continuaram a apresentar considerável redução da variância em relação aos estimadores já propostos quando o tamanho da amostra aumenta nos cadastros. Os estimadores de Bankier (cenário 2) e PML (cenário 3) apresentaram redução nas PCI's quando o tamanho de N_{ab} aumenta. Isto fornece evidência de que estes estimadores são sensíveis ao aumento da quantidade de informações do domínio ab , pois à medida em que N_{ab} aumenta, é maior a chance de selecionar mais elementos deste domínio, e constatou-se que os estimadores mencionados apresentaram problemas devido a esse fato. Ainda, é possível notar que se for considerada a quantidade $w_{ab} = N_{ab}/N$, onde $N = N_A + N_B - N_{ab}$ (total populacional), tem-se que os estimadores foram avaliados nas situações em que $w^{ab} = 500/4500 = 0.11$, $w_{ab} = 1000/4500 = 0.22$ e $w_{ab} = 1500/4500 = 0.33$. Como mostram as tabelas (6.8.) a (6.12), a PCI para estes estimadores diminuiu quando w_{ab} ficou acima de 0.10 (10%) do total populacional.

A distribuição empírica de cada um dos estimadores propostos foi obtida, como mostram as figuras (6.7) a (6.9), e fornece evidência de que a distribuição assintótica de todos os estimadores propostos é centrada no real valor do parâmetro, que no estudo de simulação foi a média populacional. Para todos os casos, há evidência de que os estimadores obtidos seguem distribuição normal, como apresentado no capítulo 4.

Os gráficos referentes aos desvios padrões dos estimadores revelam que sob o cenário 2, os estimadores regressão sob as estratégias de Hartley e Bankier apresentaram desempenho similar na redução da variância, com pequena vantagem para o estimador regressão sob a estratégia BLG. Sob o cenário 3, o estimador regressão sob a estratégia de Fuller & Burmeister apresentou melhor desempenho para todos os casos. Desconsiderando-se os cenários,

os estimadores regressão sob as estratégias PML e Fuller & Burmeister foram os que apresentaram menor variância em comparação aos demais estimadores, bem como apresentaram também maior redução com o aumento do tamanho da amostra em cada cadastro.

2.2 Estudo sobre a distribuição de $\hat{N}_{ab,s}$ e $\hat{N}_{ab,PML}$

Os estimadores de N_{ab} propostos por Fuller & Burmeister (1972) e Skinner & Rao (1996), e denotados por $\hat{N}_{ab,s}$ e $\hat{N}_{ab,PML}$, são definidos como a solução de equações quadráticas e utilizados nos casos em que não é possível identificar duplicatas nas amostras obtidas de A e B . Além disso, não estão em função dos valores das variáveis utilizadas no processo de estimação. As figuras (6.10) a (6.12) ilustram a distribuição empírica de $\hat{N}_{ab,s}$ e as figuras (6.13) a (6.15) ilustram a distribuição empírica de $\hat{N}_{ab,PML}$. Para cada uma das populações geradas na tabela (6.10), foram obtidas 1000 estimativas de N_{ab} de acordo com as estratégias mencionadas e construídos gráficos de densidade estimada.

É possível verificar que, para $N_{ab} = 500$ por exemplo, à medida em que os tamanhos de amostra aumentam, a distribuição empírica das estimativas fornece evidência de que as distribuições de $\hat{N}_{ab,s}$ e $\hat{N}_{ab,PML}$ convergem para uma distribuição centrada no real valor de N_{ab} , o que sugere distribuição normal. Além disso, verifica-se também que há evidência de redução da variância do estimador para tamanhos de amostra elevados, como o observado para $n_A = n_B = 1250$ em ambos os casos. Para os demais valores de referência para N_{ab} analisados no estudo de simulação, o comportamento observado foi o mesmo, o que evidencia que $\hat{N}_{ab,s}$ e $\hat{N}_{ab,PML}$ são estimadores que garantem uma boa aproximação para N_{ab} , e essa aproximação se torna ainda melhor à medida em que os tamanhos de amostra nos cadastros aumentam.

Tabela 6.2: Avaliação das estratégias de Hartley, Fuller, Bankier e PML para $N = 4500$, $N_{ab} = 500$ e desvios unitários.

Cenário	Estimador	$n_A = n_B$	Estatística			
			Desvio Padrão	$100 \times VR$	EQM	PCI
2	\bar{y}_H	125	0.1111	-0.0191	0.0123	0.953
		250	0.0768	0.0386	0.0059	0.940
		500	0.0493	0.0059	0.0024	0.955
		750	0.0384	0.0189	0.0014	0.952
		1000	0.0298	0.0201	0.0008	0.952
		1250	0.0249	0.0255	0.0006	0.957
	\bar{y}_{BLG}	125	0.3311	-0.0453	0.0042	0.948
		250	0.2223	0.0915	0.0496	0.946
		500	0.1477	0.0740	0.0219	0.947
		750	0.1139	-0.0208	0.0219	0.944
		1000	0.0976	0.0446	0.0095	0.959
		1250	0.0804	0.0034	0.0064	0.948
3	\bar{y}_{FB}	125	0.1108	-0.0224	0.0122	0.953
		250	0.0861	0.0374	0.0059	0.942
		500	0.0692	0.0058	0.0024	0.953
		750	0.0393	0.0019	0.0014	0.952
		1000	0.0303	0.0019	0.0008	0.951
		1250	0.0278	0.0025	0.0006	0.956
	\bar{y}_{PML}	125	0.1108	-0.0002	0.0117	0.950
		250	0.0767	0.0003	0.0045	0.948
		500	0.0491	0.0005	0.0020	0.951
		750	0.0384	0.0001	0.0010	0.951
		1000	0.0298	0.0001	0.0006	0.952
		1250	0.0248	0.0025	0.0005	0.953

Tabela 6.3: Avaliação das estratégias de estimador regressão para $N = 4500$, $N_{ab} = 500$ e desvios unitários.

Cenário	Estimador	$n_A = n_B$	Estatística			
			Desvio Padrão	$100 \times VR$	EQM	PCI
2	\bar{y}_{RH}	125	0.0649	0.0034	0.0042	0.948
		250	0.0454	0.0044	0.0021	0.939
		500	0.0297	0.0240	0.0009	0.956
		750	0.0221	0.0192	0.0004	0.949
		1000	0.0180	0.0204	0.0003	0.950
		1250	0.0143	0.0194	0.0002	0.947
	\bar{y}_{RBLG}	125	0.0687	0.0511	0.0047	0.956
		250	0.0488	0.0629	0.0024	0.939
		500	0.0315	0.0400	0.0010	0.954
		750	0.0235	0.0342	0.0005	0.940
		1000	0.0191	0.0375	0.0004	0.943
		1250	0.0151	0.0392	0.0002	0.942
3	\bar{y}_{RFB}	125	0.1108	-0.0224	0.0123	0.953
		250	0.0442	0.0420	0.0020	0.944
		500	0.0287	0.0230	0.0008	0.952
		750	0.0215	0.0189	0.0005	0.940
		1000	0.0176	0.0184	0.0003	0.956
		1250	0.0138	0.0203	0.0001	0.950
	\bar{y}_{RPML}	125	0.0821	0.0218	0.0067	0.952
		250	0.0544	0.0394	0.0029	0.951
		500	0.0358	0.0258	0.0013	0.949
		750	0.0280	0.0142	0.0008	0.955
		1000	0.0229	0.0116	0.0005	0.943
		1250	0.0184	0.0053	0.0003	0.947

Figura 6.1: Desvio Padrão dos estimadores propostos para $N = 4500$, $N_{ab} = 500$ e desvios unitários.

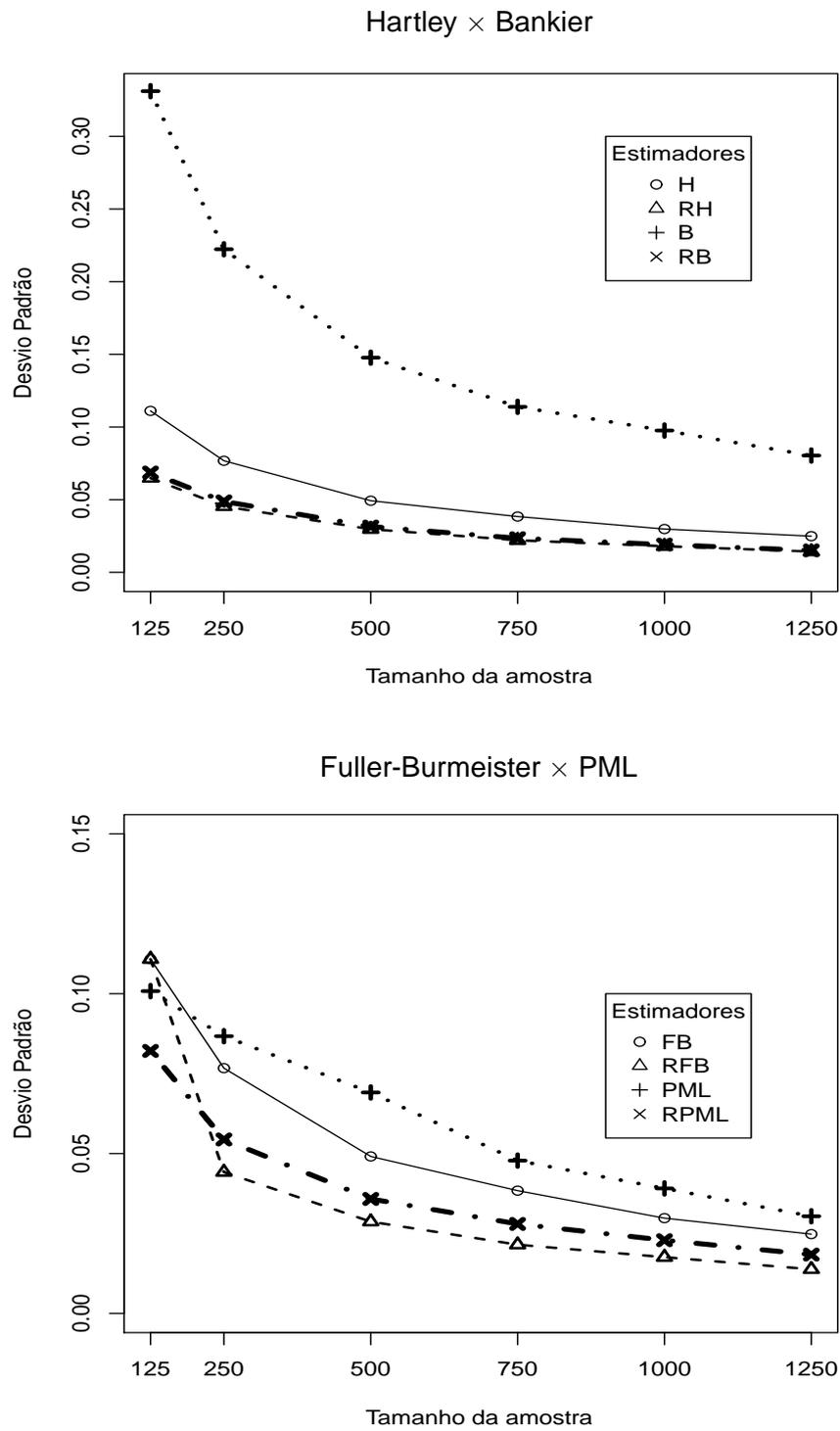


Tabela 6.4: Avaliação das estratégias de Hartley, Fuller, Bankier e PML para $N = 4500$, $N_{ab} = 1000$ e desvios unitários.

Cenário	Estimador	$n_A = n_B$	Estatística			
			Desvio Padrão	$100 \times VR$	EQM	PCI
2	\bar{y}_H	125	0.1151	0.0559	0.0133	0.945
		250	0.0793	0.0470	0.0063	0.947
		500	0.0524	0.0414	0.0028	0.949
		750	0.0409	0.0385	0.0017	0.950
		1000	0.0320	0.0219	0.0010	0.950
		1250	0.0274	0.0267	0.0007	0.949
	\bar{y}_{BLG}	125	0.4540	0.0938	0.2062	0.951
		250	0.2924	0.0149	0.0855	0.951
		500	0.2018	0.0967	0.0409	0.947
		750	0.1538	0.0647	0.0237	0.940
		1000	0.1278	0.0578	0.0163	0.942
		1250	0.1023	0.0679	0.0105	0.942
3	\bar{y}_{FB}	125	0.1154	0.0570	0.0133	0.943
		250	0.0792	0.0112	0.0063	0.946
		500	0.0525	0.0405	0.0028	0.950
		750	0.0410	0.0384	0.0017	0.950
		1000	0.0320	0.0219	0.0010	0.949
		1250	0.0274	0.0266	0.0008	0.947
	\bar{y}_{PML}	125	0.1125	0.0112	0.0126	0.940
		250	0.0761	0.0123	0.0058	0.949
		500	0.0518	0.0590	0.0026	0.953
		750	0.0399	0.0189	0.0016	0.955
		1000	0.0321	0.0002	0.0010	0.949
		1250	0.0274	0.0003	0.0008	0.947

Tabela 6.5: Avaliação das estratégias de estimador regressão para $N = 4500$, $N_{ab} = 1000$ e desvios unitários.

Cenário	Estimador	$n_A = n_B$	Estatística			
			Desvio Padrão	$100 \times VR$	EQM	PCI
2	\bar{y}_{RH}	125	0.0798	0.0138	0.0056	0.948
		250	0.0762	0.0118	0.0058	0.953
		500	0.0336	0.0656	0.0011	0.941
		750	0.0255	0.0121	0.0006	0.950
		1000	0.0212	0.0034	0.0004	0.945
		1250	0.0186	0.0032	0.0003	0.946
	\bar{y}_{RBLG}	125	0.0693	0.1228	0.0051	0.938
		250	0.0469	0.1304	0.0025	0.934
		500	0.0314	0.1258	0.0012	0.921
		750	0.0229	0.1197	0.0008	0.908
		1000	0.0188	0.1128	0.0006	0.910
		1250	0.0152	0.1143	0.0004	0.900
3	\bar{y}_{RFB}	125	0.0674	0.0188	0.0045	0.946
		250	0.0446	0.0155	0.0019	0.946
		500	0.0308	0.0045	0.0009	0.944
		750	0.0234	0.0059	0.0005	0.949
		1000	0.0189	0.0044	0.0004	0.939
		1250	0.0162	0.0029	0.0003	0.947
	\bar{y}_{RPML}	125	0.0817	0.1120	0.0069	0.950
		250	0.0539	0.0790	0.0030	0.946
		500	0.0357	0.0008	0.0014	0.942
		750	0.0267	0.1006	0.0008	0.919
		1000	0.0210	0.1033	0.0006	0.913
		1250	0.0177	0.0991	0.0004	0.915

Figura 6.2: Desvio Padrão das estratégias de estimador regressão para $N = 4500$, $N_{ab} = 1000$ e desvios unitários.

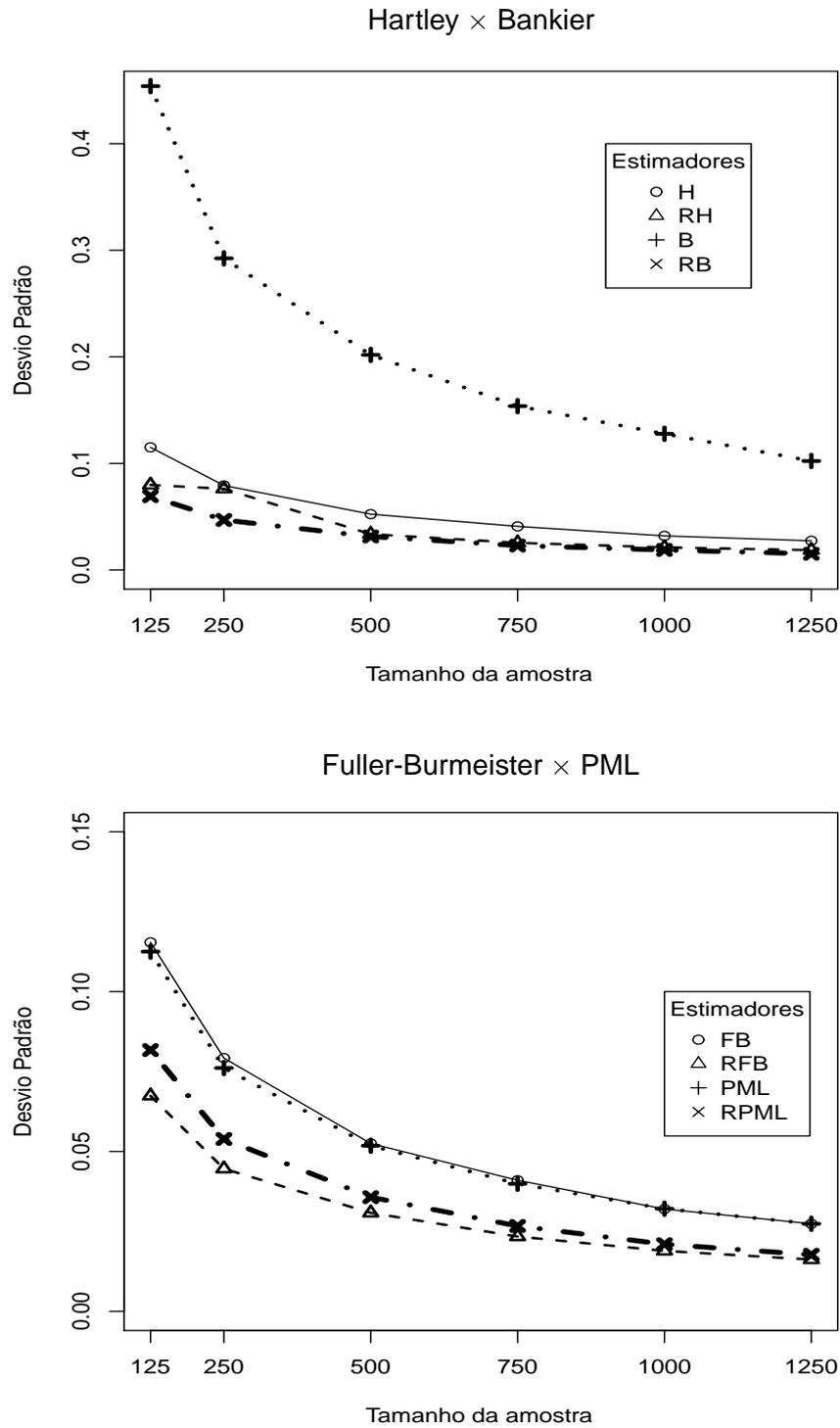


Tabela 6.6: Avaliação das estratégias de Hartley, Fuller, Bankier e PML para $N = 4500$, $N_{ab} = 1500$ e desvios unitários.

Cenário	Estimador	$n_A = n_B$	Estatística			
			Desvio Padrão	$100 \times VR$	EQM	PCI
2	\bar{y}_H	125	0.1169	0.0430	0.0137	0.952
		250	0.0792	0.0341	0.0063	0.952
		500	0.0541	0.0056	0.0029	0.950
		750	0.0428	0.0014	0.0018	0.953
		1000	0.0358	0.0012	0.0013	0.950
		1250	0.0283	0.0012	0.0008	0.950
	\bar{y}_{BLG}	125	0.5182	0.1671	0.2688	0.949
		250	0.3474	0.0975	0.1207	0.954
		500	0.2496	0.0448	0.0623	0.949
		750	0.1940	0.0676	0.0376	0.958
		1000	0.1628	0.0228	0.0265	0.956
		1250	0.1328	0.0149	0.0176	0.950
3	\bar{y}_{FB}	125	0.1164	0.0388	0.0135	0.946
		250	0.0792	0.0324	0.0062	0.953
		500	0.0541	0.0499	0.0029	0.947
		750	0.0427	0.0116	0.0018	0.953
		1000	0.0358	0.0133	0.0013	0.950
		1250	0.0284	0.0014	0.0008	0.950
	\bar{y}_{PML}	125	0.1194	0.0458	0.0195	0.946
		250	0.0892	0.0344	0.0082	0.951
		500	0.0642	0.0399	0.0059	0.949
		750	0.0527	0.0316	0.0019	0.950
		1000	0.0378	0.0233	0.0017	0.950
		1250	0.0289	0.0230	0.0009	0.939

Tabela 6.7: Avaliação das estratégias de estimador regressão para $N = 4500$, $N_{ab} = 1500$ e desvios unitários.

Cenário	Estimador	$n_A = n_B$	Estatística			
			Desvio Padrão	$100 \times VR$	EQM	PCI
2	\bar{y}_{RH}	125	0.0789	0.0261	0.0062	0.955
		250	0.0541	0.0153	0.0029	0.948
		500	0.0373	0.0158	0.0013	0.953
		750	0.0358	0.0023	0.0008	0.954
		1000	0.0285	0.0134	0.0013	0.950
		1250	0.0197	0.0008	0.0004	0.952
	\bar{y}_{RBLG}	125	0.0703	0.1547	0.0053	0.940
		250	0.0465	0.1221	0.0024	0.943
		500	0.0316	0.1296	0.0012	0.904
		750	0.0237	0.1289	0.0008	0.901
		1000	0.0177	0.1327	0.0006	0.931
		1250	0.0140	0.1410	0.0005	0.899
3	\bar{y}_{RFB}	125	0.0684	0.0189	0.0047	0.941
		250	0.0461	0.0083	0.0021	0.951
		500	0.0311	0.0150	0.0009	0.949
		750	0.0240	0.0035	0.0006	0.947
		1000	0.0235	0.0133	0.0013	0.950
		1250	0.0165	0.0046	0.0003	0.943
	\bar{y}_{RPML}	125	0.0843	0.1165	0.0073	0.946
		250	0.0535	0.1326	0.0031	0.938
		500	0.0372	0.1387	0.0016	0.943
		750	0.0282	0.1271	0.0010	0.919
		1000	0.0209	0.1325	0.0007	0.937
		1250	0.0160	0.1255	0.0005	0.899

Figura 6.3: Desvio Padrão das estratégias de estimador regressão para $N = 4500$, $N_{ab} = 1500$ e desvios unitários.

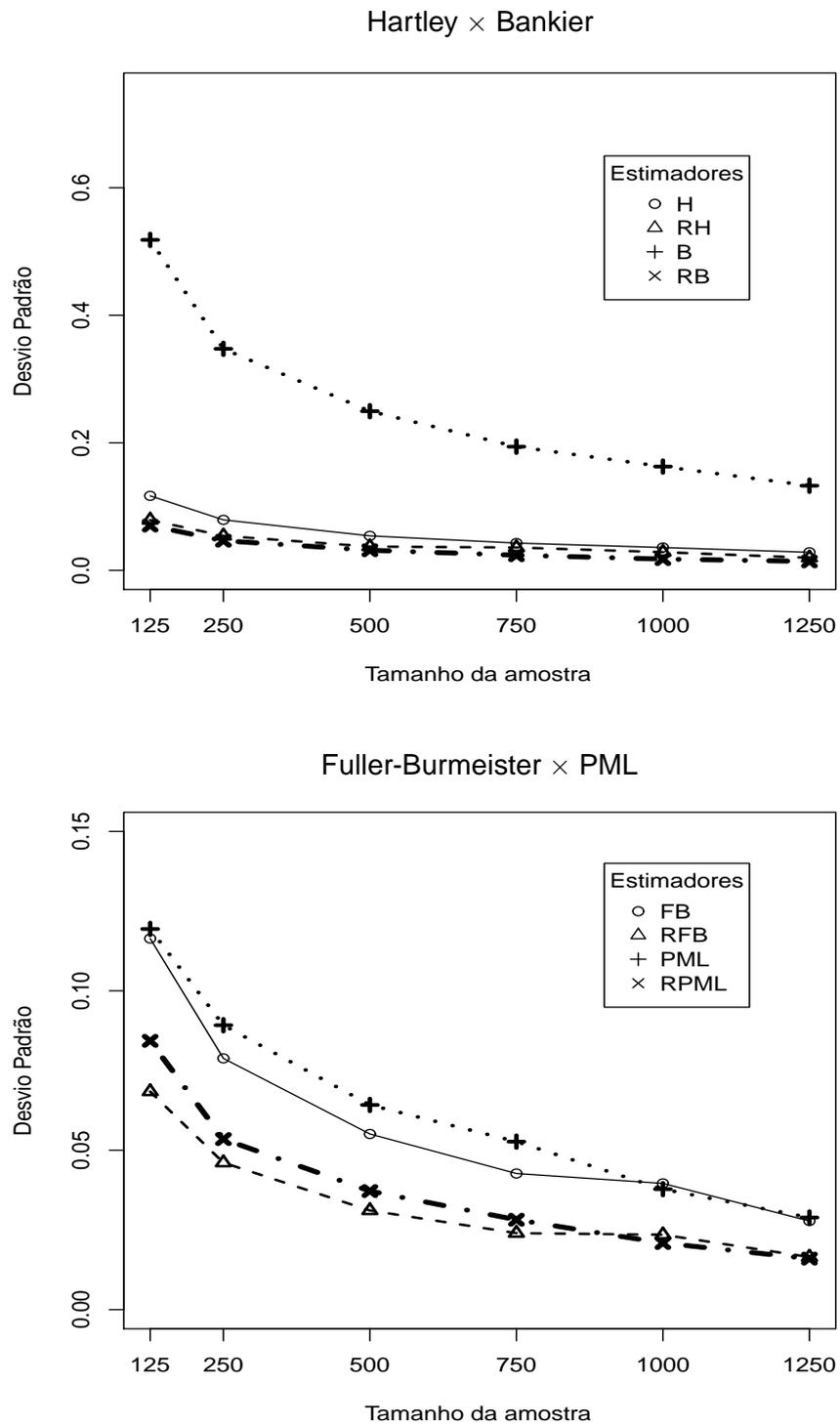


Tabela 6.8: Avaliação das estratégias de Hartley, Fuller, Bankier e PML para $N = 4500$, $N_{ab} = 500$ e desvios diferentes.

Cenário	Estimador	$n_A = n_B$	Estatística			
			Desvio Padrão	$100 \times VR$	EQM	PCI
2	\bar{y}_H	125	0.2719	0.1317	0.0742	0.948
		250	0.1863	0.0351	0.0347	0.946
		500	0.1212	0.0268	0.0147	0.953
		750	0.0928	0.0190	0.0100	0.937
		1000	0.0751	0.0236	0.0056	0.953
		1250	0.0614	0.0274	0.0037	0.961
	\bar{y}_{BLG}	125	0.4062	0.1653	0.1655	0.950
		250	0.2795	0.0856	0.0782	0.945
		500	0.1866	0.0409	0.0349	0.944
		750	0.1423	0.0215	0.0203	0.946
		1000	0.1220	0.0481	0.0149	0.957
		1250	0.1004	0.0191	0.0101	0.953
3	\bar{y}_{FB}	125	0.2714	0.1389	0.0740	0.940
		250	0.1859	0.0322	0.0346	0.947
		500	0.1209	0.0273	0.0146	0.955
		750	0.0928	0.0185	0.0086	0.939
		1000	0.0751	0.0232	0.0056	0.953
		1250	0.0614	0.0272	0.0037	0.961
	\bar{y}_{PML}	125	0.2818	0.1400	0.0741	0.947
		250	0.1944	0.0325	0.0350	0.947
		500	0.1218	0.0276	0.0147	0.960
		750	0.1032	0.0190	0.0090	0.940
		1000	0.0842	0.0233	0.0056	0.949
		1250	0.0745	0.0235	0.0038	0.959

Tabela 6.9: Avaliação das estratégias de estimador regressão propostos para $N = 4500$, $N_{ab} = 500$ e desvios diferentes.

Cenário	Estimador	$n_A = n_B$	Estatística			
			Desvio Padrão	$100 \times VR$	EQM	PCI
2	\bar{y}_{RH}	125	0.0649	0.029	0.0042	0.948
		250	0.0454	0.0396	0.0021	0.939
		500	0.0298	0.0193	0.0009	0.956
		750	0.0221	0.0144	0.0004	0.948
		1000	0.0181	0.0156	0.0003	0.953
		1250	0.0143	0.0146	0.0002	0.949
	\bar{y}_{RBLG}	125	0.0687	0.1056	0.0049	0.945
		250	0.0488	0.0398	0.0025	0.935
		500	0.0315	0.1166	0.0012	0.927
		750	0.0235	0.1224	0.0008	0.902
		1000	0.0191	0.1911	0.0006	0.880
		1250	0.0151	0.1518	0.0005	0.811
3	\bar{y}_{RFB}	125	0.0617	0.0275	0.0038	0.951
		250	0.0442	0.0371	0.0020	0.945
		500	0.0287	0.0182	0.0008	0.953
		750	0.0215	0.0135	0.0005	0.940
		1000	0.0176	0.0136	0.0003	0.956
		1250	0.0138	0.0155	0.0002	0.954
	\bar{y}_{RPML}	125	0.0815	0.1680	0.0071	0.938
		250	0.0550	0.1861	0.0036	0.936
		500	0.0360	0.1728	0.0018	0.911
		750	0.0281	0.1509	0.0012	0.881
		1000	0.0232	0.1586	0.0010	0.863
		1250	0.0185	0.1518	0.0007	0.818

Figura 6.4: Desvio Padrão das estratégias de estimador regressão para $N = 4500$, $N_{ab} = 500$ e desvios diferentes.

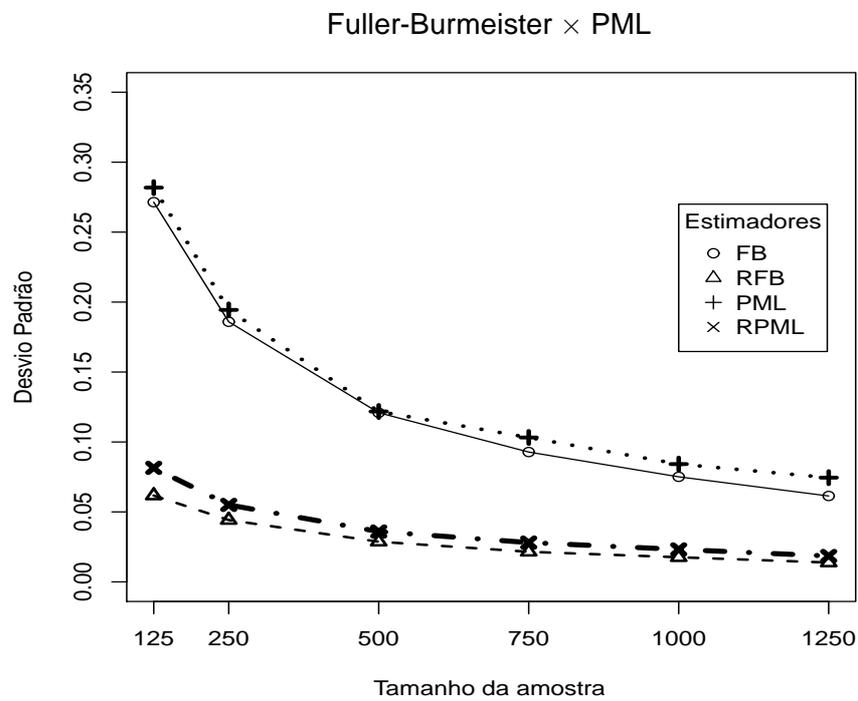
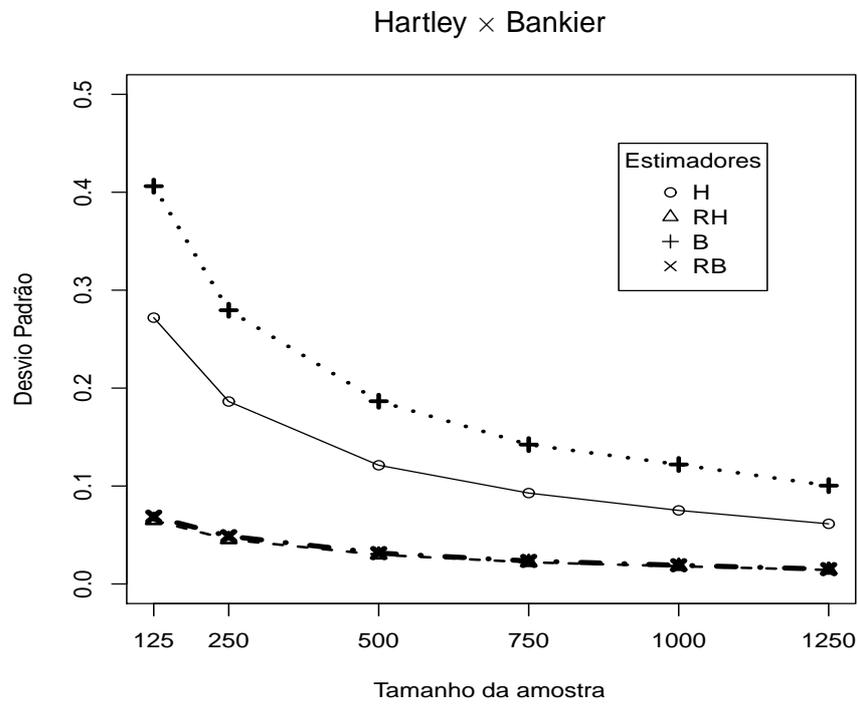


Tabela 6.10: Avaliação das estratégias de Hartley, Fuller, Bankier e PML para $N = 4500$, $N_{ab} = 1000$ e desvios diferentes.

Cenário	Estimador	$n_A = n_B$	Estatística			
			Desvio Padrão	$100 \times VR$	EQM	PCI
2	\bar{y}_H	125	0.2187	0.0714	0.0479	0.954
		250	0.1445	0.0417	0.0209	0.954
		500	0.0994	0.0919	0.0098	0.953
		750	0.0740	0.0853	0.5485	0.955
		1000	0.0638	0.0207	0.0041	0.951
		1250	0.0514	0.0667	0.0027	0.955
	\bar{y}_{BLG}	125	0.4822	0.2200	0.2332	0.955
		250	0.3281	0.1280	0.1079	0.948
		500	0.2137	0.0159	0.0456	0.957
		750	0.1693	0.0261	0.0287	0.950
		1000	0.1418	0.0726	0.0202	0.949
		1250	0.1116	0.0228	0.0124	0.954
3	\bar{y}_{FB}	125	0.2184	0.0676	0.0478	0.953
		250	0.1446	0.0430	0.0209	0.953
		500	0.0994	0.0097	0.0100	0.953
		750	0.0741	0.0086	0.0054	0.955
		1000	0.0638	0.0207	0.0041	0.951
		1250	0.0514	0.0668	0.0027	0.955
	\bar{y}_{PML}	125	0.2203	0.0678	0.0490	0.953
		250	0.1354	0.0425	0.0210	0.953
		500	0.1042	0.0101	0.0115	0.953
		750	0.0839	0.0088	0.0060	0.955
		1000	0.0745	0.0210	0.0043	0.951
		1250	0.0516	0.0670	0.0030	0.955

Tabela 6.11: Avaliação das estratégias de estimador regressão propostos para $N = 4500$, $N_{ab} = 1000$ e desvios diferentes.

Cenário	Estimador	$n_A = n_B$	Estatística			
			Desvio Padrão	$100 \times VR$	EQM	PCI
2	\bar{y}_{RH}	125	0.0690	0.0311	0.0047	0.945
		250	0.0478	0.0419	0.0023	0.947
		500	0.0315	0.0338	0.0010	0.956
		750	0.0240	0.0330	0.0005	0.940
		1000	0.0196	0.0350	0.0004	0.945
		1250	0.0163	0.0340	0.0003	0.939
	\bar{y}_{RBLG}	125	0.0759	0.0539	0.0057	0.887
		250	0.0544	0.0174	0.0029	0.852
		500	0.0349	0.0294	0.0012	0.856
		750	0.0267	0.0369	0.0007	0.906
		1000	0.0216	0.0292	0.0005	0.909
		1250	0.0183	0.0286	0.0003	0.912
3	\bar{y}_{RFB}	125	0.0647	0.0261	0.0042	0.955
		250	0.0461	0.0432	0.0022	0.952
		500	0.0296	0.0361	0.0010	0.949
		750	0.0227	0.0299	0.0005	0.945
		1000	0.0186	0.0341	0.0004	0.948
		1250	0.0155	0.0344	0.0002	0.940
	\bar{y}_{RPML}	125	0.0933	0.0923	0.0088	0.806
		250	0.0640	0.1247	0.0043	0.901
		500	0.0398	0.1113	0.0018	0.912
		750	0.0325	0.0993	0.0012	0.901
		1000	0.0259	0.1149	0.0010	0.909
		1250	0.0205	0.1000	0.0005	0.888

Figura 6.5: Desvio Padrão das estratégias de estimador regressão para $N = 4500$, $N_{ab} = 1000$ e desvios diferentes.

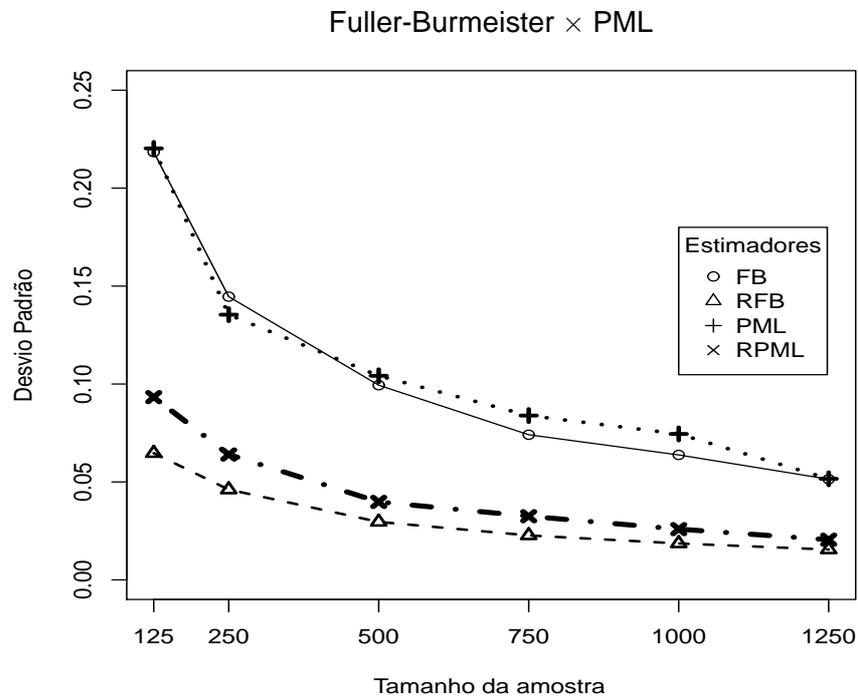
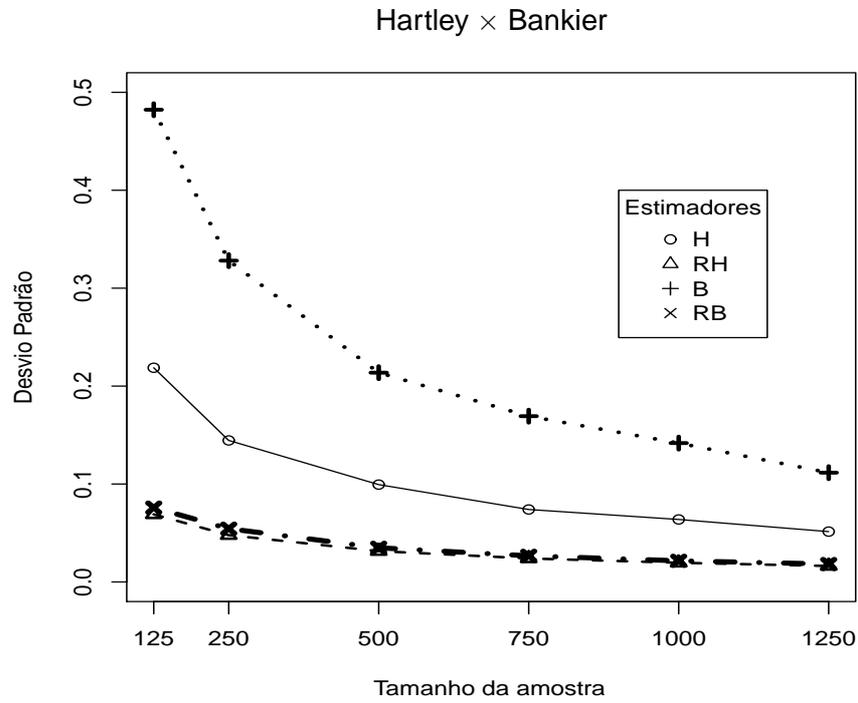


Tabela 6.12: Avaliação das estratégias de Hartley, Fuller, Bankier e PML para $N = 4500$, $N_{ab} = 1500$ e desvios diferentes.

Cenário	Estimador	$n_A = n_B$	Estatística			
			Desvio Padrão	$100 \times VR$	EQM	PCI
2	\bar{y}_H	125	0.1985	0.0724	0.0395	0.948
		250	0.1391	0.0109	0.0194	0.952
		500	0.0976	0.0023	0.0095	0.951
		750	0.0760	0.0031	0.0058	0.952
		1000	0.0593	0.0014	0.0035	0.955
		1250	0.0491	0.0085	0.0024	0.955
	\bar{y}_{BLG}	125	0.5062	0.0260	0.2562	0.945
		250	0.3405	0.0549	0.1159	0.943
		500	0.2384	0.0088	0.0568	0.944
		750	0.1869	0.0479	0.0349	0.942
		1000	0.1577	0.0174	0.0249	0.948
		1250	0.1259	0.0448	0.0159	0.941
3	\bar{y}_{FB}	125	0.1983	0.0694	0.0394	0.949
		250	0.1391	0.0117	0.0193	0.951
		500	0.0977	0.0029	0.0095	0.951
		750	0.0760	0.0031	0.0058	0.951
		1000	0.0593	0.0014	0.0035	0.954
		1250	0.0491	0.0083	0.0024	0.957
	\bar{y}_{PML}	125	0.1900	0.0702	0.0400	0.939
		250	0.1404	0.0120	0.0198	0.941
		500	0.0880	0.0032	0.0095	0.941
		750	0.0776	0.0035	0.0056	0.939
		1000	0.0590	0.0019	0.0037	0.950
		1250	0.0479	0.0081	0.0028	0.957

Tabela 6.13: Avaliação das estratégias de estimador regressão propostos para $N = 4500$, $N_{ab} = 1500$ e desvios diferentes.

Cenário	Estimador	$n_A = n_B$	Estatística			
			Desvio Padrão	$100 \times VR$	EQM	PCI
2	\bar{y}_{RH}	125	0.0748	0.0303	0.0056	0.954
		250	0.0503	0.0264	0.0025	0.954
		500	0.0331	0.0158	0.0011	0.956
		750	0.0254	0.0087	0.0006	0.952
		1000	0.0212	0.0375	0.0005	0.940
		1250	0.0183	0.0123	0.0003	0.944
	\bar{y}_{RBLG}	125	0.0886	0.5346	0.0113	0.788
		250	0.0587	0.5340	0.0069	0.721
		500	0.0376	0.5283	0.0048	0.670
		750	0.0286	0.5191	0.0041	0.763
		1000	0.0246	0.5499	0.0043	0.799
		1250	0.0211	0.5249	0.0038	0.709
3	\bar{y}_{RFB}	125	0.0692	0.0231	0.0048	0.952
		250	0.0458	0.0177	0.0021	0.951
		500	0.0298	0.0142	0.0010	0.945
		750	0.0226	0.0114	0.0005	0.951
		1000	0.0196	0.0332	0.0004	0.945
		1250	0.0167	0.0134	0.0003	0.954
	\bar{y}_{RPML}	125	0.1097	0.4907	0.0149	0.786
		250	0.0716	0.5027	0.0082	0.785
		500	0.0495	0.4918	0.0054	0.797
		750	0.0379	0.4948	0.0044	0.701
		1000	0.0324	0.4768	0.0038	0.731
		1250	0.0262	0.4950	0.0037	0.751

Figura 6.6: Desvio Padrão das estratégias de estimador regressão para $N = 4500$, $N_{ab} = 1500$ e desvios diferentes.

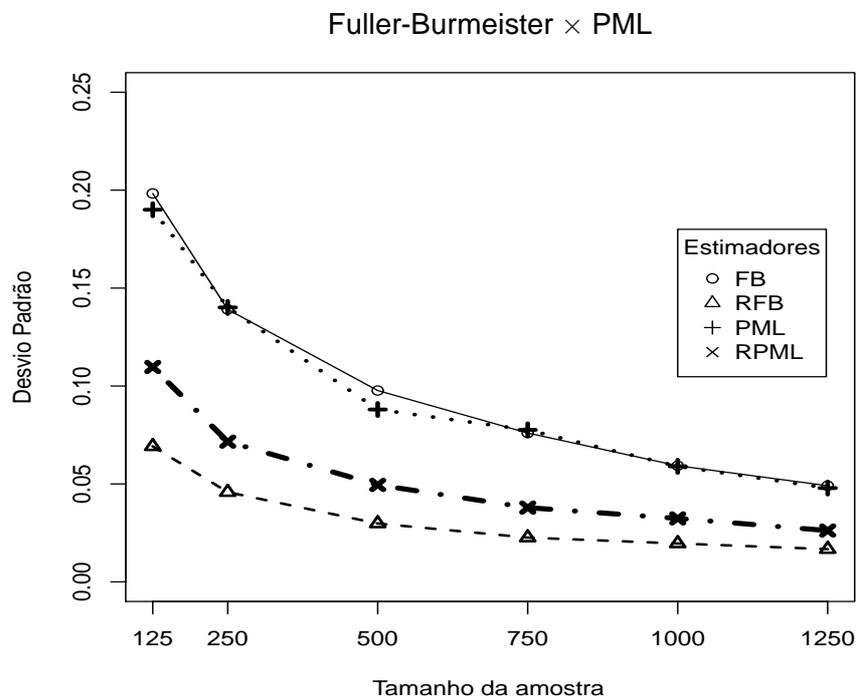
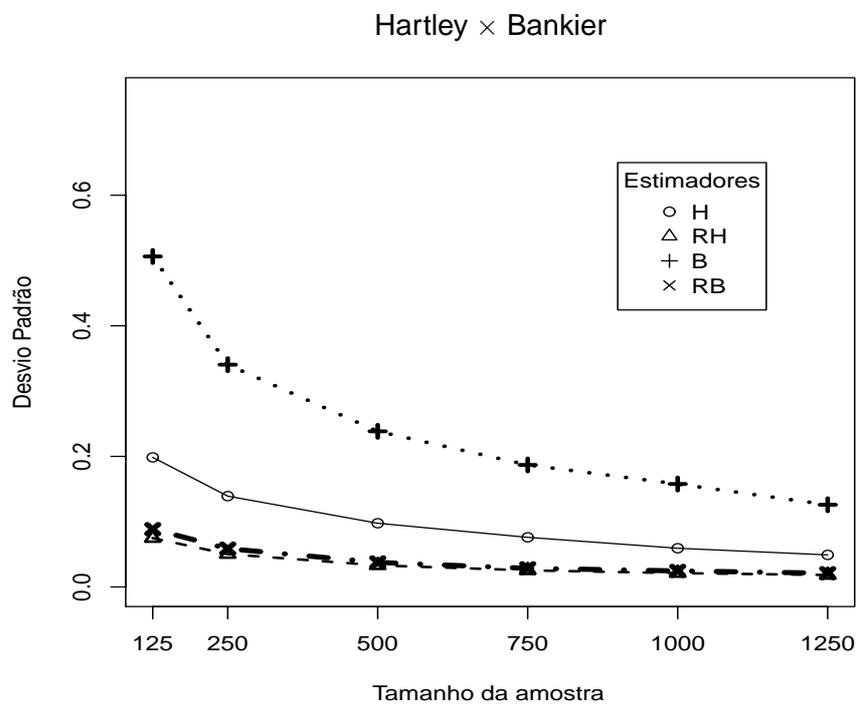


Figura 6.7: Densidade estimada dos estimadores propostos para $N = 4500$ e $N_{ab} = 500$

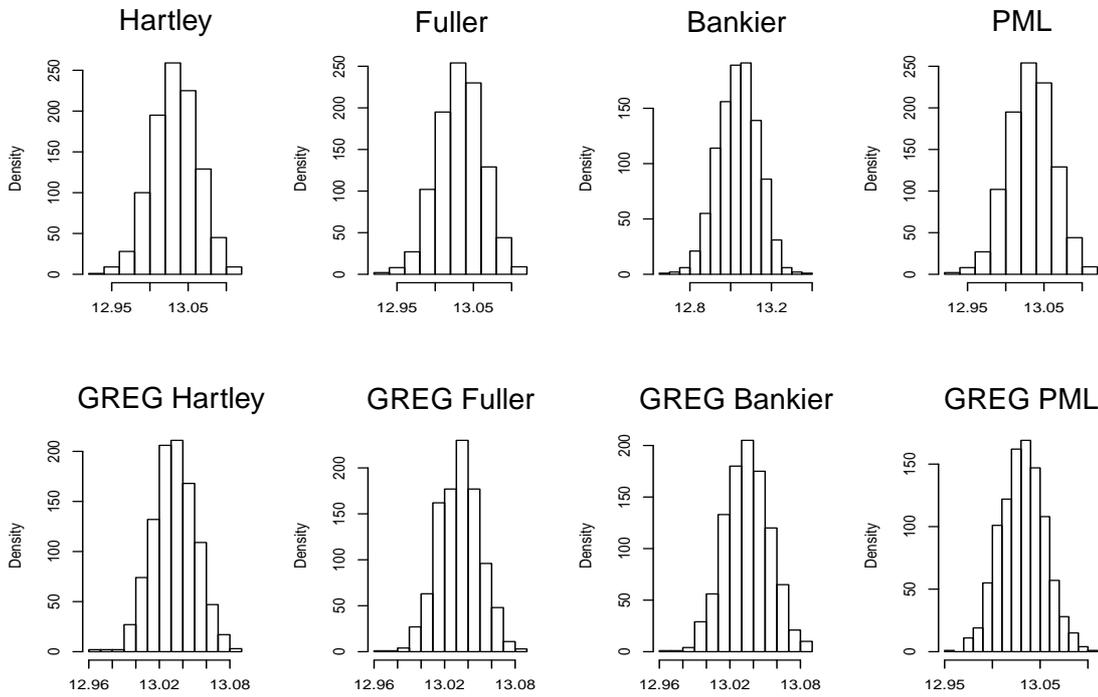


Figura 6.8: Densidade estimada dos estimadores propostos para $N = 4500$ e $N_{ab} = 1000$

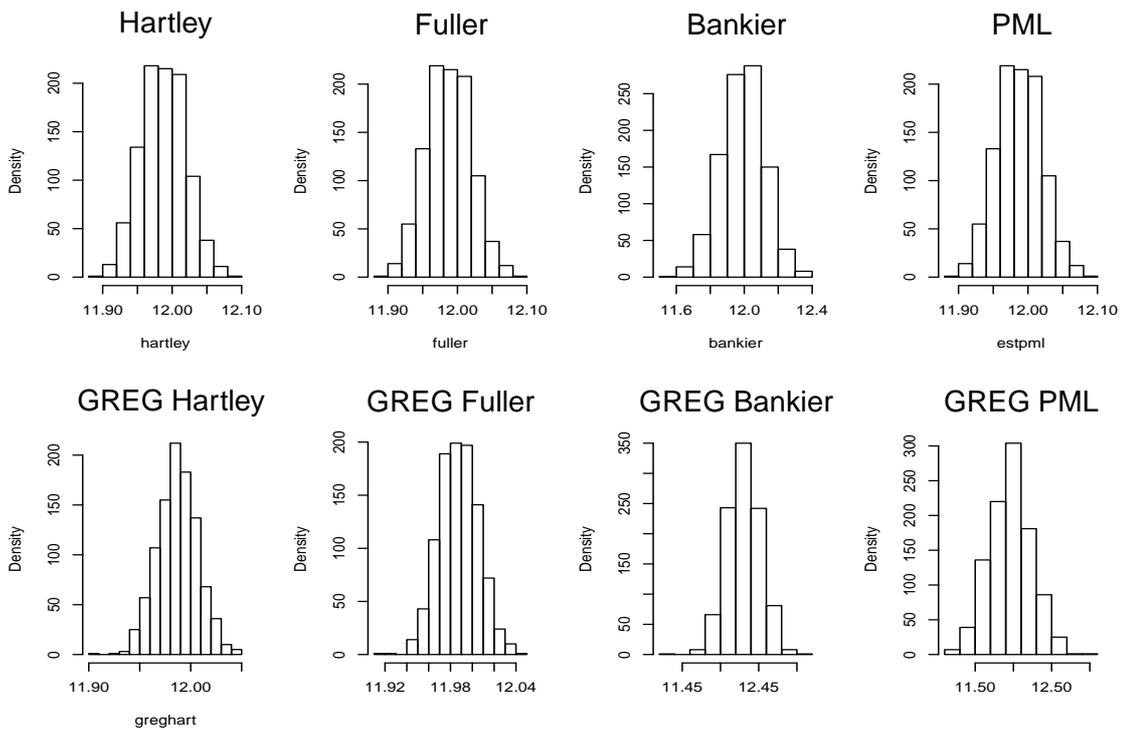


Figura 6.9: Densidade estimada dos estimadores propostos para $N = 4500$ e $N_{ab} = 1500$

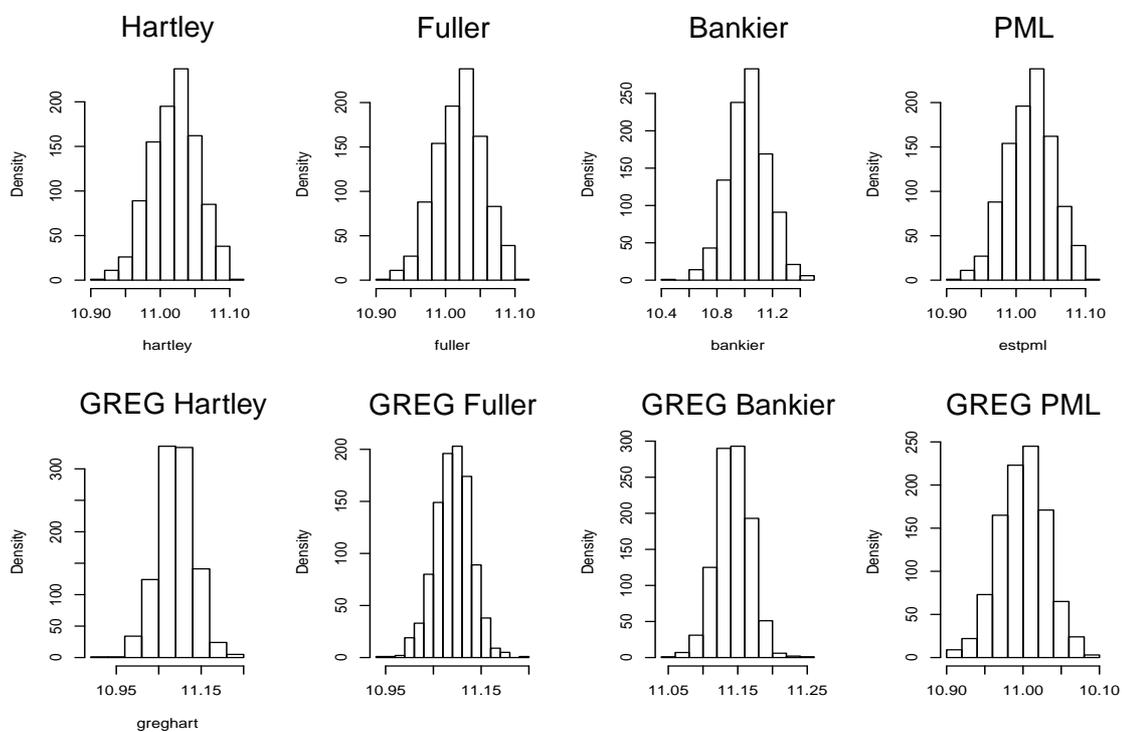


Figura 6.10: Densidade Estimada de $\hat{N}_{ab,s}$, com valor de referência $N_{ab} = 500$

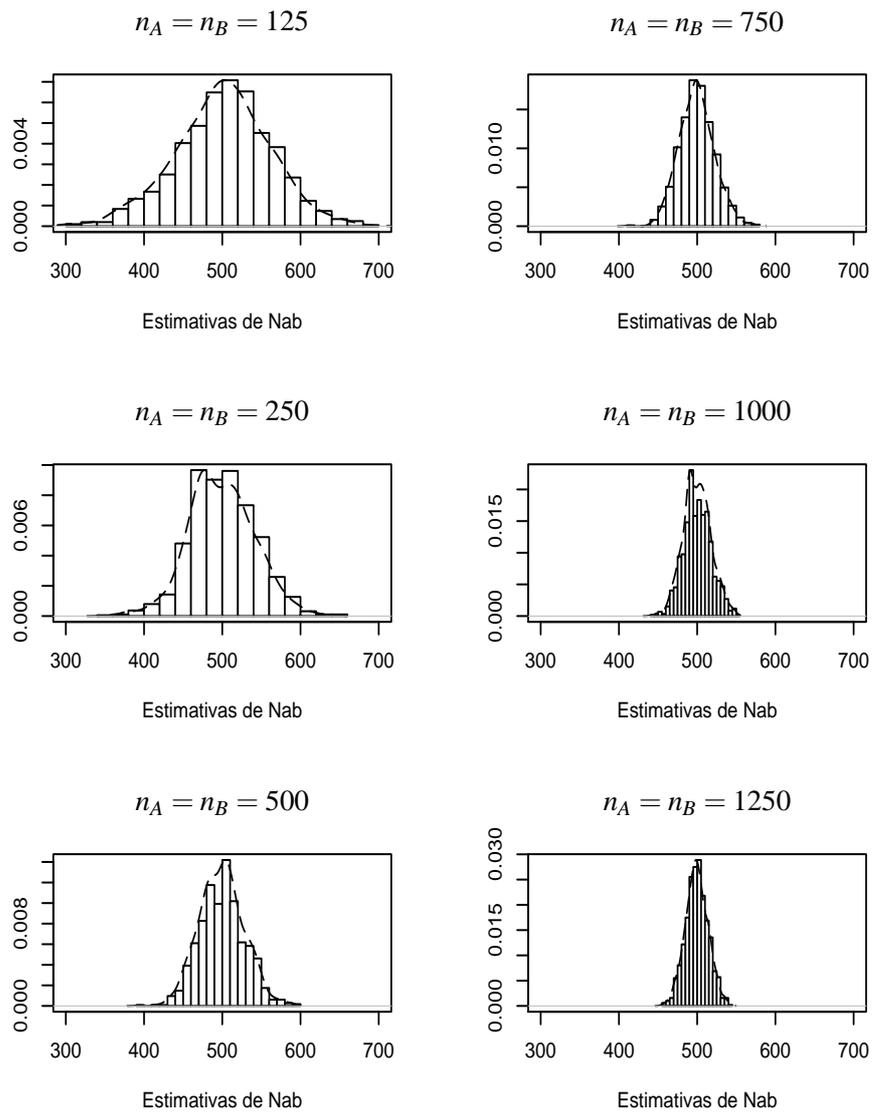


Figura 6.11: Densidade Estimada de $\hat{N}_{ab,s}$, com valor de referência $N_{ab} = 1000$

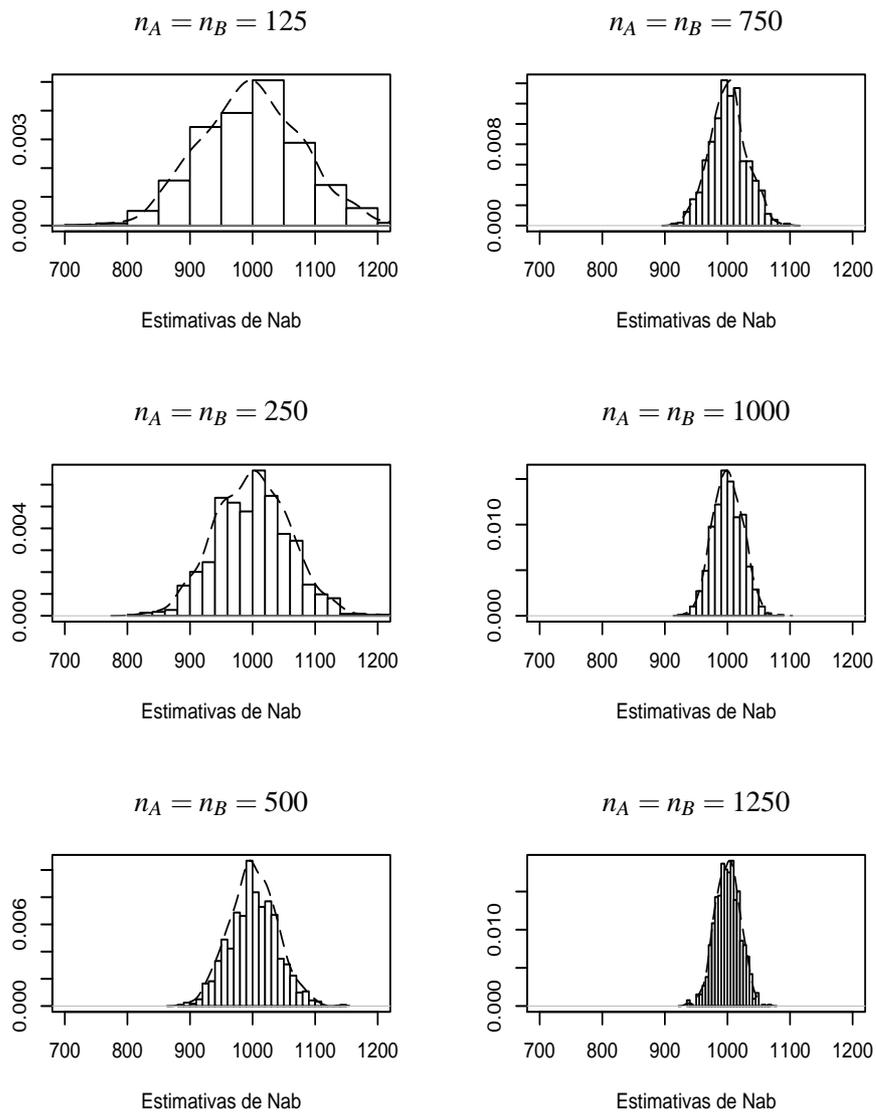


Figura 6.12: Densidade Estimada de $\hat{N}_{ab,s}$, com valor de referência $N_{ab} = 1500$

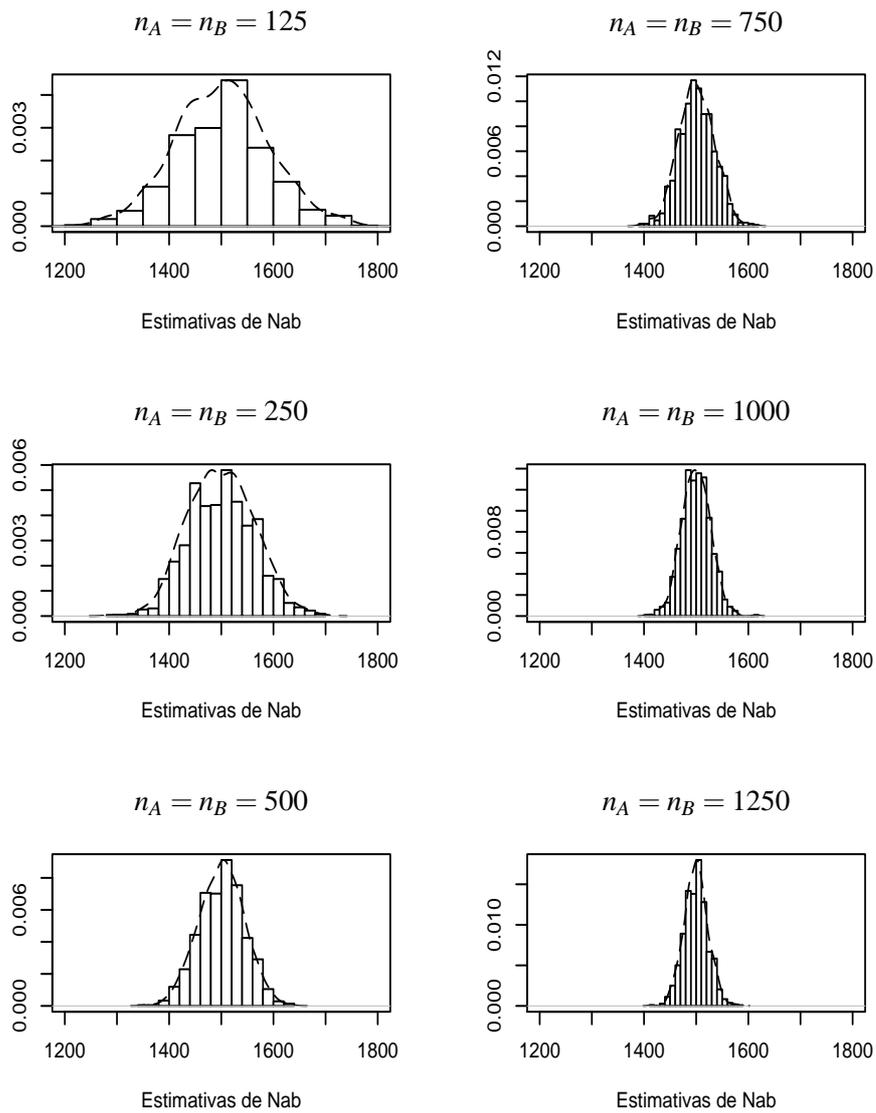


Figura 6.13: Densidade Estimada de $\hat{N}_{ab,PML}$, com valor de referência $N_{ab} = 500$

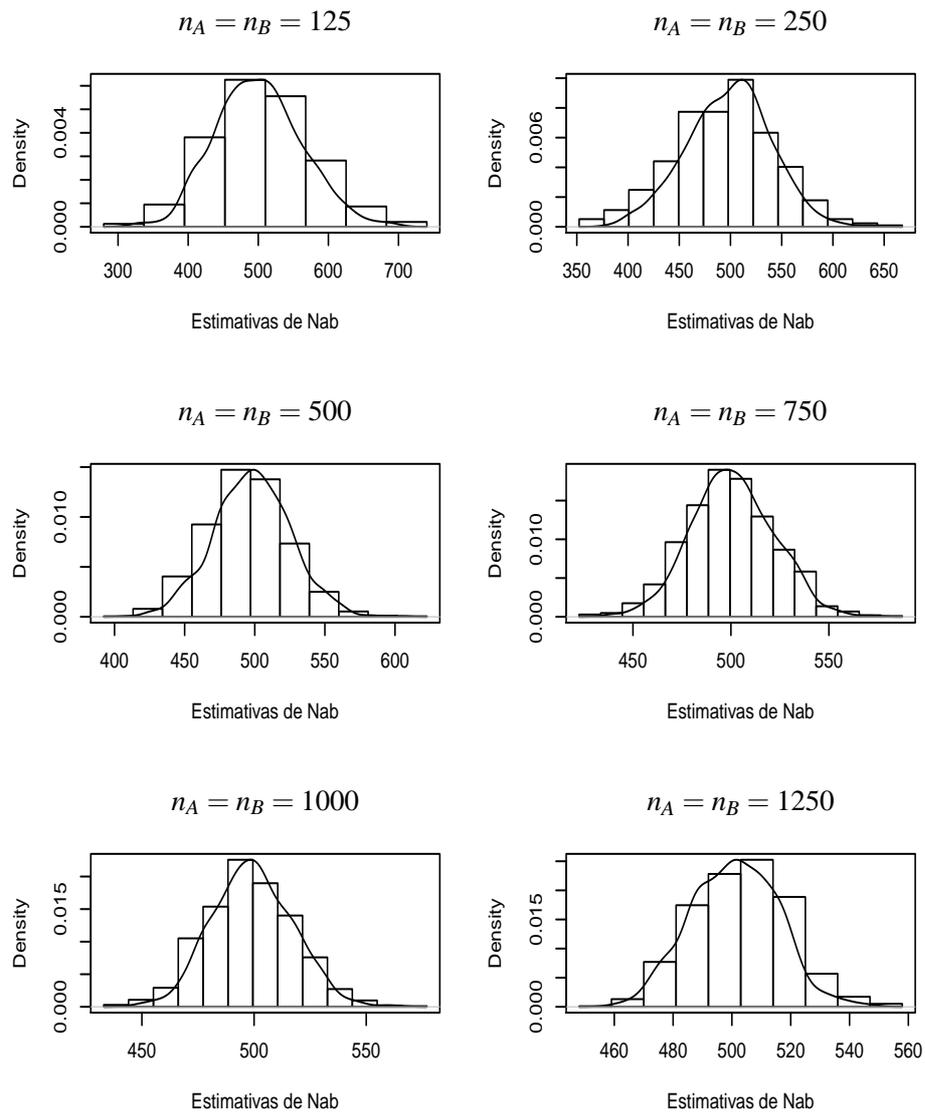


Figura 6.14: Densidade Estimada de $\hat{N}_{ab,PML}$, com valor de referência $N_{ab} = 1000$

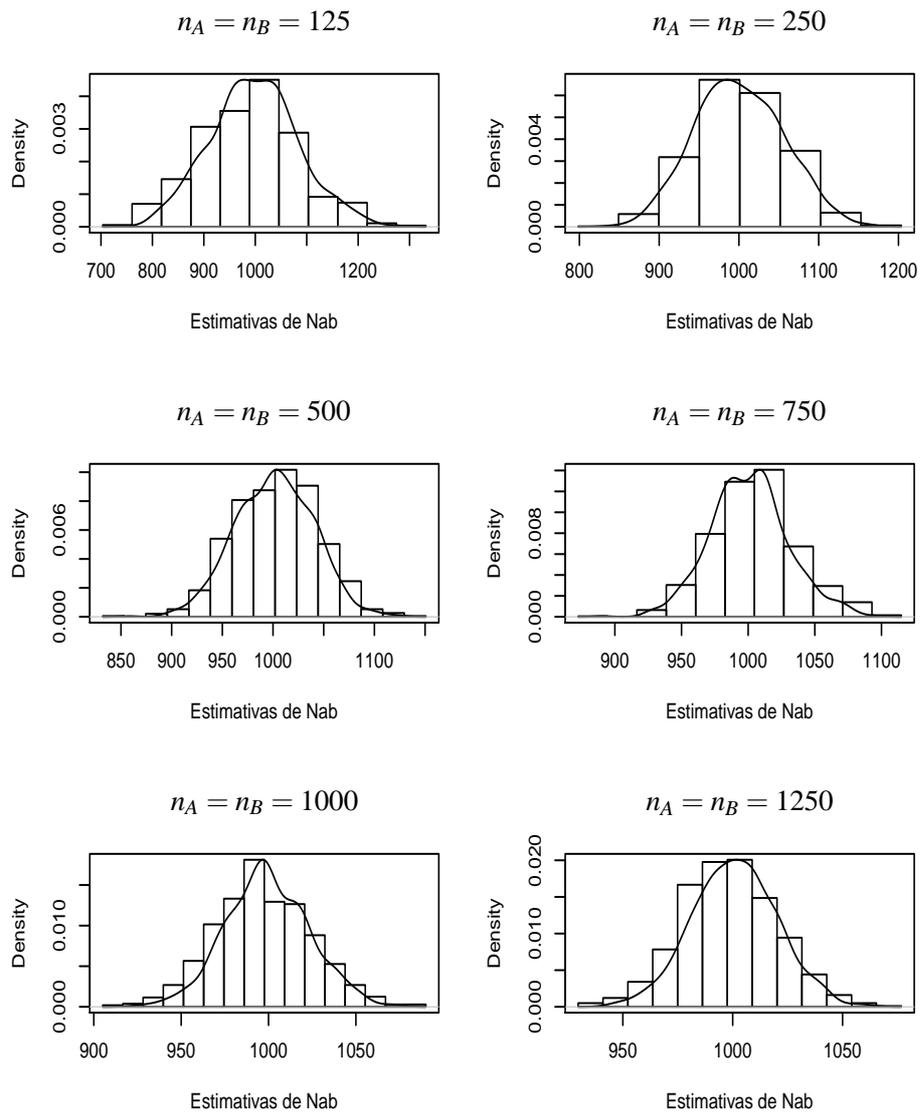
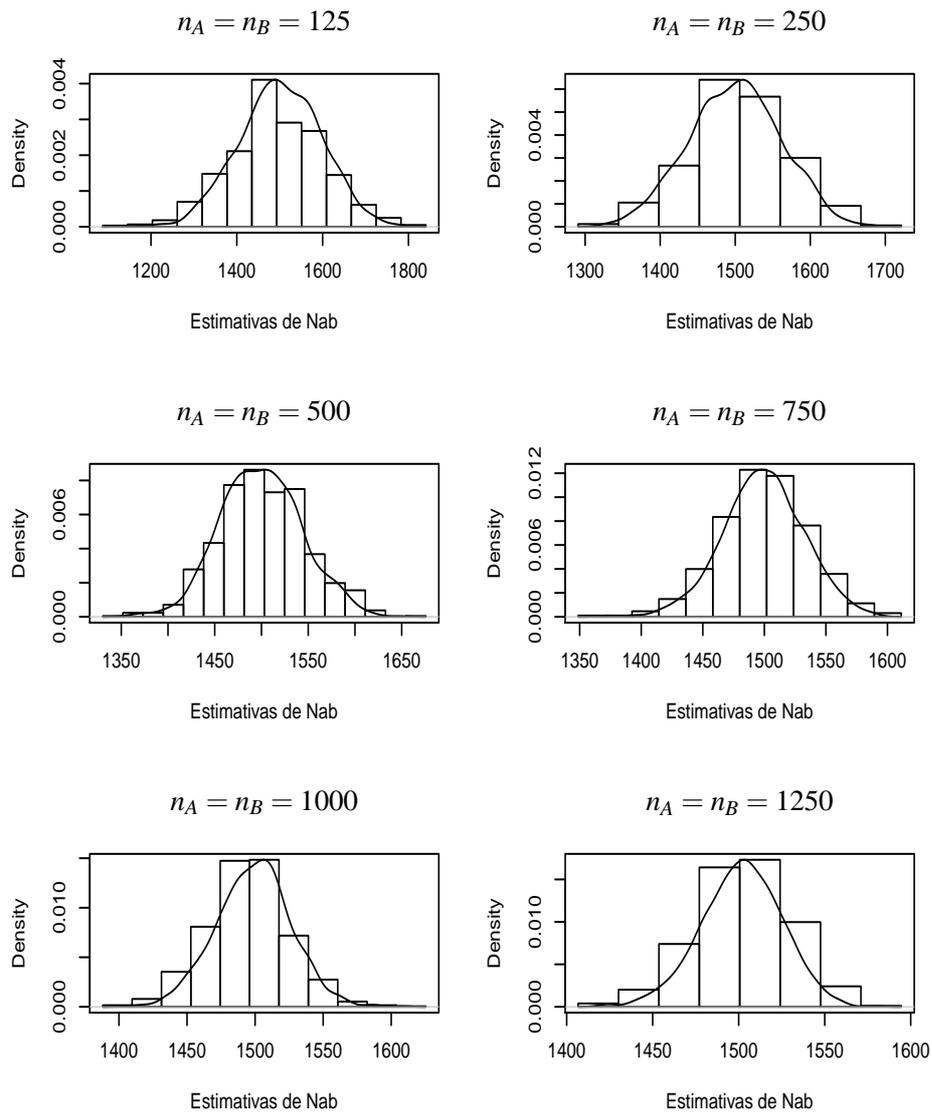


Figura 6.15: Densidade Estimada de $\hat{N}_{ab,PML}$, com valor de referência $N_{ab} = 1500$



CAPÍTULO 7

Considerações Finais

No desenvolvimento desta tese, foram apresentados quatro novos estimadores do tipo regressão generalizados para uso na abordagem de cadastro duplo. Estes estimadores foram propostos a partir das estratégias apresentadas por Hartley (1962), Fuller & Burmeister (1972), Bankier, Lepkowski & Groves (1986) e Skinner & Rao (1996). A forma da variância aproximada de todos os estimadores propostos foi derivada e apresentada em termos matriciais na forma $\mathbf{d}^T \Sigma \mathbf{d}$, onde \mathbf{d} e Σ são dadas de acordo com cada uma das estratégias apresentadas. Além disso, as propriedades de consistência, centralidade assintótica e expressões para o erro quadrático médio destes estimadores foram apresentadas de forma rigorosa com o intuito de motivar o uso de variáveis auxiliares na abordagem de cadastro duplo, uma vez que é esperado que informações auxiliares forneçam uma melhora no processo de estimação de um parâmetro populacional. Além disso, foi provado através do método delta que, assintoticamente, os estimadores propostos possuem distribuição normal, e que forneceu suporte para cálculo das variâncias aproximadas dos estimadores propostos.

Os resultados obtidos nesta tese generalizam todos os resultados para o estimador do tipo razão propostos por Coelho (2007), que considerou informação de apenas uma variável auxiliar em apenas um dos cadastros, além de fornecer novas formas de estimador do tipo razão sob as estratégias propostas por Bankier, Lepkowski & Groves (1986) e Skinner & Rao (1996).

O desempenho dos novos estimadores regressão para a abordagem de cadastro duplo foi avaliado através de um estudo de simulação sob os planos de amostragem aleatória simples, comparando-os com os estimadores já disponíveis na literatura. Os resultados das simulações mostraram que os estimadores regressão propostos apresentaram melhor desempenho em comparação aos estimadores propostos na literatura, apresentando considerável redução de

variância e erro quadrático médio, como era esperado.

Os estimadores propostos sob as versões de Hartley (1962) e Fuller & Burmeister (1972) apresentaram menor sensibilidade à heterogeneidade das observações dentro dos cadastros e domínios e ao aumento do tamanho populacional no domínio N_{ab} , quando comparados aos estimadores de Bankier, Lepkowski & Groves (1986) e Skinner & Rao (1996), que só apresentaram desempenho equivalente quando $N_{ab}/N \approx 0.10$, valor este que pode servir como uma recomendação exigida para uso destes estimadores quando se desejar facilidade de implementação do método.

Os estimadores regressão propostos sob as estratégias de Hartley (1962) e Fuller (2009) são mais eficientes do que os demais estimadores apresentados sob o cenário 2. Porém, para tamanhos de amostras elevados na situação descrita pelo cenário 2, é verificada vantagem de uso do estimador RBLG.

Os estimadores regressão sob as estratégias de Fuller & Burmeister e PML apresentaram melhor desempenho que os demais estimadores apresentados sob o cenário 3. O estimador regressão sob a estratégia de Fuller & Burmeister foi o que teve maior destaque por apresentar maior redução de sua variância quando o tamanho da amostra aumenta.

De um modo geral, ao desconsiderar os cenários os estimadores regressão propostos sob as estratégias de Fuller & Burmeister e PML foram os estimadores que apresentaram melhor desempenho, por apresentar maior redução da variância quando o tamanho da amostra aumenta em ambos os cadastros. Dessa forma, em um cenário de aplicação onde o tamanho populacional dos domínios é conhecido (cenário 2), recomenda-se o uso do estimador de Bankier, por ser mais eficiente que a estratégia de Hartley, além de ser relativamente simples de ser implementado. Se a situação envolver o não conhecimento do tamanho populacional dos domínios (cenário 3), recomenda-se o uso do estimador de Fuller & Burmeister.

Esta tese dá abertura para várias pesquisas sobre estratégias de inferência assistida por modelos sob a abordagem de cadastro duplo. Dentre as propostas de pesquisas futuras, por exemplo, podemos citar as seguintes:

- Apresentar os estimadores propostos sob o contexto da estratificação, considerando-se estimadores regressão nas versões separado e combinado, de modo a generalizar os resultados apresentados por Coelho & Ferraz (2007);
- Estender toda a teoria desenvolvida nesta tese para o estimador regressão logístico generalizado, utilizado para estimação em populações assistida por modelos variáveis dicotômicas, discutido por Poveda (2007);
- Apresentar os estimadores propostos para diversos planos amostrais e realizar estudos de alocação ótima de amostras aos cadastros sob estes planos;
- Analisar o desempenho dos estimadores propostos sob planos amostrais complexos;

- Analisar o desempenho dos estimadores sob técnicas de imputação de dados para tratamento da não-resposta.

Apêndice A

Resumo de conceitos importantes

A seguir, será apresentada uma revisão de alguns conceitos utilizados ao longo desta tese.

1 Ordens de Magnitude de seqüências de números reais (\mathbf{O} e \mathbf{o})

Sejam a_n e b_n seqüências de números reais. Dizemos que

1. $a_n = \mathbf{O}(b_n)$ se existirem $M > 0$ e $n_0 = n(M)$ tal que $\left| \frac{a_n}{b_n} \right| \leq M, \forall n \geq n_0$;
2. $a_n = \mathbf{o}(b_n)$ se para todo $\varepsilon > 0$ existir um número positivo $n_0 = n_0(\varepsilon)$ tal que $\left| \frac{a_n}{b_n} \right| < \varepsilon, \forall n \geq n_0$

1.1 Teorema 1

Sejam a_n, b_n, c_n e d_n seqüências de números reais. Então:

1. Se $a_n = \mathbf{o}(b_n)$, então $a_n = \mathbf{O}(b_n)$;
2. Se $a_n = \mathbf{O}(b_n)$ e $c_n = \mathbf{O}(d_n)$, então:
 - (a) $a_n c_n = \mathbf{O}(b_n d_n)$;
 - (b) $|a_n|^s = \mathbf{O}(|b_n|^s)$, para todo $s > 0$;
 - (c) $a_n + c_n = \mathbf{O}(\max\{|b_n|, |d_n|\})$
3. Se $a_n = \mathbf{o}(b_n)$ e $c_n = \mathbf{o}(d_n)$, então:
 - (a) $a_n c_n = \mathbf{o}(b_n d_n)$;
 - (b) $|a_n|^s = \mathbf{o}(|b_n|^s)$, para todo $s > 0$;
 - (c) $a_n + c_n = \mathbf{o}(\max\{|b_n|, |d_n|\})$

2 Desigualdade de Markov

Seja X uma variável aleatória qualquer. Então para todo $t > 0$, temos:

$$P(|X| \geq \varepsilon) \leq \frac{E(|X|^t)}{\varepsilon^t}$$

3 Convergência Estocástica

3.1 Ordens de Magnitude de seqüências de números reais ($\mathbf{O}_p(\cdot)$ e $\mathbf{o}_p(\cdot)$)

Sejam $\{X_n\}_{n \geq 1}$ seqüência de variáveis aleatórias e $\{b_n\}_{n \geq 1}$ uma seqüência de números reais (ou variáveis aleatórias). Dessa forma, temos que:

1. $X_n = \mathbf{O}_p(b_n)$ se para todo $\varepsilon > 0$ existir M tal que $M = M(\varepsilon)$ e um número inteiro positivo n_0 onde $n_0 = n_0(\varepsilon)$ tais que

$$P\left(\left|\frac{X_n}{b_n}\right| \geq M\right) \leq \varepsilon, \quad \forall n \geq n_0$$

ou seja, a seqüência $\left\{\frac{X_n}{b_n}\right\}_{n \geq 1}$ é limitada em probabilidade para todo n suficientemente grande.

2. $X_n = \mathbf{o}_p(b_n)$ se para todo $\varepsilon > 0$ e para todo $\lambda > 0$ existir n_0 tal que $n_0 = n_0(\varepsilon, \lambda)$ tais que

$$P\left(\left|\frac{X_n}{b_n}\right| \geq \varepsilon\right) < \lambda, \quad \forall n \geq n_0,$$

ou seja, $X_n = \mathbf{o}_p(b_n)$ se para todo $\varepsilon > 0$, $P\left(\left|\frac{X_n}{b_n}\right| \geq \varepsilon\right) \rightarrow 0$, quando $n \rightarrow \infty$.

3.2 Teorema 2

Sejam $\{X_n\}_{n \geq 1}$, $\{Y_n\}_{n \geq 1}$ seqüência de variáveis aleatórias e $\{a_n\}_{n \geq 1}$, $\{b_n\}_{n \geq 1}$ seqüências de números reais (ou variáveis aleatórias). Temos então que:

1. Se $X_n = \mathbf{o}_p(a_n)$, então $X_n = \mathbf{O}_p(a_n)$;
2. Se $X_n = \mathbf{O}_p(a_n)$ e $Y_n = \mathbf{O}_p(b_n)$, então:
 - (a) $X_n Y_n = \mathbf{O}_p(a_n b_n)$;
 - (b) $|X_n|^s = \mathbf{O}_p(|a_n|^s)$, para todo $s > 0$;
 - (c) $X_n + Y_n = \mathbf{O}_p(\max\{|b_n|, |d_n|\})$

3. Se $X_n = \mathbf{o}_p(a_n)$ e $Y_n = \mathbf{o}_p(b_n)$, então:

- (a) $X_n Y_n = \mathbf{o}_p(a_n b_n)$;
- (b) $|X_n|^s = \mathbf{o}_p(|a_n|^s)$, para todo $s > 0$;
- (c) $X_n + Y_n = \mathbf{o}_p(\max\{|b_n|, |d_n|\})$
- (d) Se $X_n = \mathbf{O}_p(a_n)$ e $Y_n = \mathbf{o}_p(b_n)$, então: $X_n Y_n = \mathbf{o}_p(a_n b_n)$

4 Tipos de convergência estocástica

1. Se $X_n \xrightarrow{q.c.} X$, então $X_n \xrightarrow{P} X$;

2. Sejam $\{X_n; n \geq 1\}$ e $\{Y_n; n \geq 1\}$ sequências de variáveis aleatórias. Se $|X_n - Y_n| \xrightarrow{P} 0$ e $Y_n \xrightarrow{d} X$, então $X_n \xrightarrow{d} X$;

3. Se $X_n \xrightarrow{P} X$, então $X_n \xrightarrow{d} X$;

4. Seja c uma constante. Se $X_n \xrightarrow{d} c$, então $X_n \xrightarrow{P} c$;

5. Se $X_n \xrightarrow{d} X$, então $X_n \xrightarrow{P} 0$, então, $X_n Y_n \xrightarrow{P} 0$;

6. Teorema de Slutsky:

- Se $X_n \xrightarrow{d} X$ e $Y_n \xrightarrow{P} 0$, então $X_n Y_n \xrightarrow{P} 0$;
- Seja c uma constante. Se $X_n \xrightarrow{d} X$ e $Y_n \xrightarrow{P} c$, então $X_n + Y_n \xrightarrow{d} X + c$, $X_n Y_n \xrightarrow{d} Xc$, $\frac{X_n}{Y_n} \xrightarrow{d} \frac{X}{c}$;

7. Seja $g: \mathbb{R} \rightarrow \mathbb{R}$ uma função contínua. Então,

- Se $X_n \xrightarrow{q.c.} X$, então $g(X_n) \xrightarrow{q.c.} g(X)$;
- Se $X_n \xrightarrow{P} X$, então $g(X_n) \xrightarrow{P} g(X)$;
- Se $X_n \xrightarrow{d} X$, então $g(X_n) \xrightarrow{d} g(X)$;

8. (**Método Delta**) Suponha que $\sqrt{n}(X_n - \mu) \xrightarrow{d} N(0, \sigma^2)$. Se g é uma função diferenciável em μ , então:

$$\sqrt{n}(g(X_n) - g(\mu)) \xrightarrow{d} N\left(0, \sigma^2 [g'(\mu)]^2\right).$$

Apêndice B

Demonstrações de resultados propostos

A seguir serão apresentadas as demonstrações dos resultados enunciados nesta tese de doutorado.

Consistência do estimador de Hartley

Resultado: O estimador \bar{y}_H^{AUB} é consistente, ou seja:

$$P(|\bar{y}_H^{AUB} - \bar{y}_v| > \varepsilon | \mathcal{F}_v) \longrightarrow 0, \text{ quando } v \rightarrow \infty$$

Prova: A prova deste resultado é construtiva, e depende das condições a seguir:

C1. Dado $N_v = N_{Av} + N_{Bv} - N_{abv}$, seja θ_v^{AUB} o parâmetro de interesse. Então, existe sequência de estimadores de θ_v^{AUB} , denotada por $\{\hat{\theta}_v; v \geq 1\}$ que satisfaz a seguinte propriedade: $\forall \varepsilon > 0$, existe $\delta_\varepsilon \in (0, \infty)$ tal que

$$P(|\hat{\theta}_v^{AUB} - \theta_v^{AUB}| > \delta_\varepsilon) < \varepsilon$$

C2. No cadastro A , $\forall N_{Av}$, $\forall k \in \mathcal{A}_v$, $\min_{k \in \mathcal{A}_v} \pi_{vk}^A \geq \pi_v^* > 0$, onde $N_{Av} \pi_v^* \rightarrow \infty$ e existe $s > 0$ tal que $(N_{Av})^{1/2+s} (\pi_{vk}^A) \rightarrow \infty$ e

$$\max_{k \in U_v^A} \sum_{l \in U_v^A: k \neq l} (\pi_{vkl} - \pi_{vk}^A \pi_{vl}^A)^2 = O((N_{Av})^{-2s}), N_{Av} \rightarrow \infty.$$

C3. As observações da variável de interesse no cadastro A , $\{y_k\}_{k \in U_v^A}$ satisfazem:

$$\limsup_{N_{Av} \rightarrow \infty} \sum_{k \in S_v^A} y_k^2 < \infty,$$

C4. As mesmas condições **C2** e **C3** são admitidas ao se considerar as correspondentes quantidades no cadastro B .

Breidt & Opsomer (2008) mostraram que a condição suficiente para garantir a consistência de acordo com o plano é a convergência em média quadrática, ou seja, basta provar que

$$\bar{y}^{A \cup B} \xrightarrow{(2)} \mu_v^{A \cup B}$$

Logo,

$$\begin{aligned} \bar{y}_{Hv} - \mu_v &= w_{Av} \bar{y}_{HTv}^A + w_{Bv} \bar{y}_{HTv}^B - w_{abv} \bar{y}_{HTv}^{A \cap B} \\ &- (w_{Av} \mu_{Av} + w_{Bv} \mu_{Bv} - w_{abv} \mu_{(A \cap B)v}) \\ &= w_{Av} (\bar{y}_{Av} - \mu_{Av}) + w_{Bv} (\bar{y}_{Bv} - \mu_{Bv}) - w_{abv} (\bar{y}^{A \cap B} - \mu_{Av}) \\ &= w_{Av} (\bar{y}_{Av} - \mu_{Av}) + w_{Bv} (\bar{y}_{Bv} - \mu_{Bv}) - w_{abv} (\bar{y}_{(A \cap B)v} - \mu_{(A \cap B)v}) \end{aligned}$$

$$\begin{aligned} (\bar{y}_{Hv} - \mu_v)^2 &= (w_{Av})^2 (\bar{y}_{Av} - \mu_{Av})^2 + (w_{Bv})^2 (\bar{y}_{Bv} - \mu_{Bv})^2 \\ &+ 2w_{Av}w_{Bv} (\bar{y}_{Av} - \mu_{Av}) (\bar{y}_{Bv} - \mu_{Bv}) \\ &- 2w_{Av}w_{abv} (\bar{y}_{(A \cap B)v} - \mu_{(A \cap B)v}) (\bar{y}_{Av} - \mu_{Av}) \\ &- 2w_{Bv}w_{abv} (\bar{y}_{(A \cap B)v} - \mu_{(A \cap B)v}) (\bar{y}_{Bv} - \mu_{Bv}) \end{aligned}$$

$$\begin{aligned} E(\bar{y}_{Hv} - \mu_v)^2 &= (w_{Av})^2 E(\bar{y}_{Av} - \mu_{Av})^2 + (w_{Bv})^2 E(\bar{y}_{Bv} - \mu_{Bv})^2 \\ &+ 2w_{Av}w_{Bv} E(\bar{y}_{Av} - \mu_{Av}) E(\bar{y}_{Bv} - \mu_{Bv}) \\ &- 2w_{Av}w_{abv} E[(\bar{y}_{(A \cap B)v} - \mu_{(A \cap B)v}) (\bar{y}_{Av} - \mu_{Av})] \\ &- 2w_{Bv}w_{abv} E[(\bar{y}_{(A \cap B)v} - \mu_{(A \cap B)v}) (\bar{y}_{Bv} - \mu_{Bv})] \end{aligned}$$

Para as quantidades referentes ao cadastro A, tem-se que

$$\begin{aligned} E((\bar{y}_{Av} - \mu_{Av})^2) &= E(\bar{y}_{Av} - \mu_{Av})^2 = \frac{1}{(N_{Av})^2} \sum_{k,l \in U_v^A} (\pi_{vkl} - \pi_{vk}\pi_{vl}) \frac{y_k y_l}{\pi_{vk}\pi_{vl}} \\ &\leq \frac{1}{N_{Av}\pi_v^*} \sum_{k \in U_v^A} \frac{y_k^2}{N_{Av}} + \frac{1}{(N_{Av})^2 (\pi_v^*)^2} \left(\sum_{k,l \in U_v^A: k \neq l} (\pi_{vkl} - \pi_{vk}^2 \pi_{vl}^2)^2 \right)^{1/2} \left(\sum_{k,l \in U_v^A: k \neq l} y_k^2 y_l^2 \right)^{1/2} \\ &\leq \frac{1}{N_{Av}\pi_v^*} \sum_{k \in U_v^A} \frac{y_k^2}{N_{Av}} \end{aligned}$$

$$+ \frac{\pi_{N_{Av}}^2}{(N_{Av})^{1/2+s}} \left(\frac{N_{Av} \max_{k \in \mathcal{A}_v} \sum_{l \in \mathcal{A}_v: k \neq l} (\pi_{vkl} - \pi_{vk}\pi_{vl})^2}{(N_{Av})^{1-2s}} \right)^{1/2} \left\{ \left(\sum_{k \in U_v^A} \frac{y_k^2}{N_{Av}} \right)^2 \right\}^{1/2},$$

que converge para zero quando $v \rightarrow \infty$. Analogamente, $E[(\bar{y}_{Bv} - \mu_{Bv})^2] \rightarrow 0, v \rightarrow \infty$.

Admitindo que

$$E \left[(\bar{y}_{Av} - \mu_{Av}) (\bar{y}_{(A \cap B)v} - \mu_{(A \cap B)v}) \right] < \infty$$

$$E \left[(\bar{y}_{Bv} - \mu_{Bv}) (\bar{y}_{(A \cap B)v} - \mu_{(A \cap B)v}) \right] < \infty$$

é possível notar que uma vez que N_{Av} e N_{Bv} convergem para infinito com magnitudes diferentes, então a quantidade $N_v = N_{Av} + N_{Bv} - N_{(A \cap B)v}$ (onde $N_{(A \cap B)v}$ é finita) converge para infinito a uma taxa maior que as de N_{Av} e N_{Bv} . Logo, $w_{Av} \rightarrow 0, w_{Bv} \rightarrow 0$ e $w_{abv} \rightarrow 0$, quando $v \rightarrow \infty$. Logo,

$$E(\bar{y}^{A \cup B} - \mu_v)^2 \rightarrow 0, \text{ quando } v \rightarrow \infty$$

Pela desigualdade de Markov, tem-se que

$$P(|\bar{y}_v^{A \cup B} - \bar{y}_N| > \varepsilon) \leq \frac{E(\bar{y}^{A \cup B} - \bar{y}_N)^r}{\varepsilon^r}$$

Quando $r = 2$ e $v \rightarrow \infty$, temos o resultado. ■

Resultados 6, 7, 8 e 9 - página 46

(i) Estimador \bar{y}_{RH} : Prova: Para quantidades nos cadastros e nos domínios, admite-se as condições apresentadas por Isaki & Fuller (1982). É possível notar que

$$\begin{aligned} \bar{y}_{RH(v)} - \mu_v &= w_{Av} \bar{y}_{RAv} + w_{Bv} \bar{y}_{RBv} - w_{abv} \left(p \bar{y}_{Rab(A)v} + (1-p) \bar{y}_{Rab(B)v} \right) \\ &- w_{Av} \mu_{Av} + w_{Bv} \mu_{Bv} - w_{abv} \left(p \mu_{ab(A)v} + (1-p) \mu_{ab(B)v} \right) \end{aligned}$$

$$\begin{aligned} \bar{y}_{RH(v)} - \mu_v &= w_{Av} (\bar{y}_{RAv} - \mu_{Av}) + w_{Bv} (\bar{y}_{RBv} - \mu_{Bv}) \\ &- w_{A \cap B} \left\{ p (\bar{y}_{Rab(A)v} - \mu_{ab(A)v}) + (1-p) (\bar{y}_{Rab(B)v} - \mu_{ab(B)v}) \right\} \end{aligned}$$

Tem-se que para A:

$$\bar{y}_R^A - \mu_{Av} = \frac{1}{N_{Av}} \sum_{k \in \mathcal{A}_v} y_k \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right) - \frac{1}{N_{Av}} \sum_{j=0}^q \sum_{k \in \mathcal{A}_v} \hat{\beta}_{vj}^A x_{kj} \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right)$$

Através do método de linearização de Taylor aplicado no segundo termo da expressão acima, tem-se que:

$$\bar{y}_{RA} - \mu_{Av} \approx \frac{1}{N_{Av}} \sum_{k \in \mathcal{A}_v} y_k \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right) - \sum_{j=0}^q \beta_j^A (\bar{x}_{jt}^A - \bar{x}_{N_{Av}})$$

A variância aproximada de $\bar{y}_{RAv} - \mu_{Av}$ é dada por:

$$\begin{aligned} A\text{Var}(\bar{y}_{RAv} - \mu_{Av}) &= \text{Var} \left\{ \sum_{k \in \mathcal{A}_v} y_k \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right) \right\} + \text{Var} \left\{ \sum_{j=0}^q \beta_j^A (\bar{x}_{jt}^A - \bar{x}_{N_{Av}}) \right\} \\ &\quad - \text{Cov} \left[\frac{1}{N_{Av}} \sum_{k \in \mathcal{A}_v} y_k \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right) ; \sum_{j=0}^q \beta_j^A (\bar{x}_{jt}^A - \bar{x}_{N_{Av}}) \right] \end{aligned}$$

Logo,

$$\begin{aligned} A\text{Var}(\bar{y}_{RAv} - \mu_{Av}) &\leq \text{Var} \left\{ \sum_{k \in \mathcal{A}_v} y_k \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right) \right\} + \text{Var} \left\{ \sum_{j=0}^q \beta_{vj}^A (\bar{x}_{jt}^A - \bar{x}_{N_{Av}}) \right\} \\ &= \text{Var} \left\{ \sum_{k \in \mathcal{A}_v} y_k \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right) \right\} + \sum_{j=0}^q (\beta_{vj}^A)^2 \text{Var}(\bar{x}_{jt}^A - \bar{x}_{N_{Av}}) \\ &\quad + \sum_{j \neq j'} \beta_j^A \beta_{j'}^A \text{Cov}[(\bar{x}_{vj}^A - \bar{x}_{N_{Av}}) ; (\bar{x}_{vj'}^A - \bar{x}_{N_{Av}})] \\ &= \alpha_1 + \alpha_2 + \alpha_3 \end{aligned}$$

Pelos resultados obtidos por Isaki & Fuller (1982), tem-se que

1. $\alpha_1 = \text{Var} \left\{ \sum_{k \in \mathcal{A}_v} y_k \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right) \right\} = O((n_v^A)^{-2\delta});$
2. $\alpha_2 = \sum_{j=0}^q (\beta_{vj}^A)^2 \text{Var}(\bar{x}_{jt}^A - \bar{x}_{N_{Av}}) = O((n_v^A)^{-2\delta})$

Além disso, tem-se que

$$\begin{aligned}
\alpha_3 &= \sum_{j \neq j'} \beta_j^A \beta_{j'}^A \text{Cov} [(\bar{x}_{Vj}^A - \bar{x}_{N_{Av}}) ; (\bar{x}_{Vj'}^A - \bar{x}_{N_{Av}})] \\
&\leq \sum_{j \neq j'} \beta_{Vj}^A \beta_{Vj'}^A \{ \text{Cov} [(\bar{x}_{Vj}^A - \bar{x}_{N_{Av}}) ; (\bar{x}_{Vj'}^A - \bar{x}_{N_{Av}})] \}^2 \\
&\leq \sum_{j \neq j'} \beta_{Vj}^A \beta_{Vj'}^A \text{Var} [(\bar{x}_{Vj}^A - \bar{x}_{N_{Av}})] \text{Var} [(\bar{x}_{Vj'}^A - \bar{x}_{N_{Av}})] = O \left((n_V^A)^{-4\delta} \right)
\end{aligned}$$

Portanto, $A\text{Var}(\bar{y}_{RAv} - \mu_{Av}) = \alpha_1 + \alpha_2 + \alpha_3 = O \left(\max \left\{ n_{t(A)}^{-2\delta}, (n_V^A)^{-4\delta} \right\} \right) = O \left((n_V^A)^{-2\delta} \right)$

Analogamente, é possível mostrar que

- $A\text{Var}(\bar{y}_{RBv} - \mu_{Bv}) = O \left((n_V^B)^{-2\delta} \right)$
- $A\text{Var} \left(\bar{y}_{Rab(A)v} - \mu_{ab(A)v} \right) = O \left[\left(n_V^{ab(A)} \right)^{-2\delta} \right]$
- $A\text{Var} \left(\bar{y}_{Rab(B)v} - \mu_{ab(B)v} \right) = O \left[\left(n_V^{ab(B)} \right)^{-2\delta} \right]$

É possível mostrar ainda que

$$\begin{aligned}
A\text{Var} \left(\bar{y}_{RH(v)} - \mu_v \right) &= (w_{Av})^2 [A\text{Var}(\bar{y}_{RAv} - \mu_{Av})] + (w_{Bv})^2 [A\text{Var}(\bar{y}_{RBv} - \mu_{Bv})] \\
&+ (w_{abv})^2 p^2 [A\text{Var}(\bar{y}_{Rab(A)v} - \mu_{ab(A)v})] + w_{abv}^2 (1-p)^2 [A\text{Var}(\bar{y}_{Rab(B)v} - \mu_{ab(B)v})] \\
&+ 2p w_{Av} w_{abv} \text{ACov} \left[(\bar{y}_{RAv} - \mu_{Av}) ; (\bar{y}_{Rab(A)v} - \mu_{ab(A)v}) \right] \\
&+ 2(1-p) w_{Bv} w_{abv} \text{ACov} \left[(\bar{y}_{RBv} - \mu_{Bv}) ; (\bar{y}_{Rab(B)v} - \mu_{ab(B)v}) \right] \\
&\leq (w_{Av})^2 [A\text{Var}(\bar{y}_{RAv} - \mu_{Av})] + (w_{Bv})^2 [A\text{Var}(\bar{y}_{RBv} - \mu_{Bv})] \\
&+ (w_{abv})^2 p^2 [A\text{Var}(\bar{y}_{Rab(A)v} - \mu_{ab(A)v})] + w_{abv}^2 (1-p)^2 [A\text{Var}(\bar{y}_{Rab(B)v} - \mu_{ab(B)v})] \\
&+ p w_{Av} w_{abv} \left\{ \text{ACov} \left[(\bar{y}_{RAv} - \mu_{Av}) ; (\bar{y}_{Rab(A)v} - \mu_{ab(A)v}) \right] \right\}^2 \\
&+ (1-p) w_{Bv} w_{abv} \left\{ \text{ACov} \left[(\bar{y}_{RBv} - \mu_{Bv}) ; (\bar{y}_{Rab(B)v} - \mu_{ab(B)v}) \right] \right\}^2 \\
&\leq (w_{Av})^2 [A\text{Var}(\bar{y}_{RAv} - \mu_{Av})] + w_B^2 [A\text{Var}(\bar{y}_{RBv} - \mu_{Bv})] \\
&+ (w_{abv})^2 p^2 [A\text{Var}(\bar{y}_{Rab(A)v} - \mu_{ab(A)v})] + w_{abv}^2 (1-p)^2 [A\text{Var}(\bar{y}_{Rab(B)v} - \mu_{ab(B)v})]
\end{aligned}$$

$$\begin{aligned}
& + pW_{Av}W_{abv}A\text{Var}[(\bar{y}_{RAv} - \mu_{Av})]A\text{Var}\left[\left(\bar{y}_{Rab(A)v} - \mu_{ab(A)v}\right)\right] \\
& + (1-p)W_{Bv}W_{abv}A\text{Var}[(\bar{y}_{RBv} - \mu_{Bv})]A\text{Var}\left[\left(\bar{y}_{Rab(B)v} - \mu_{ab(B)v}\right)\right] \\
& = \mathbf{O}\left\{\max\left((n_{Av})^{-2\delta}, (n_{Bv})^{-2\delta}\right)\right\}.
\end{aligned}$$

Com este resultado e conforme condições descritas para a situação de apenas um cadastro, apresentadas por Isaki & Fuller (1982), tem-se finalmente que

$$\bar{y}_{RH(v)} - \mu_v = \mathbf{O}_p\left(\max\left\{(n_{Av})^{-\delta}; (n_{Bv})^{-\delta}\right\}\right).$$

Robinson & Särndal (1983) mostraram que os estimadores regressão para cadastros e domínios são assintoticamente centrados, ou seja:

$$\lim_{v \rightarrow \infty} \{E(\bar{y}_R^C | \mathbf{Y}_v) - \mu_v\} = 0 \quad \text{e} \quad \lim_{v \rightarrow \infty} \{E(\bar{y}_R^d | \mathbf{Y}_v) - \mu_v\} = 0$$

Estes resultados são garantidos pela desigualdade de Markov, que fornece condições suficientes para mostrar que

$$\lim_{v \rightarrow \infty} E\left\{\bar{y}_{GREG(v)}^C - \bar{y}_v | \mathbf{Y}_v\right\} = 0,$$

da seguinte forma: seja

$$\bar{y}_R^C - \mu_v = \frac{1}{N_v^C} \sum_{k \in U_v^C} y_k \left(\frac{I_{vk}}{\pi_{vk}} - 1\right) - \sum_{j=1}^q \hat{\beta}_{jv} \frac{1}{N_v} \sum_{k \in U_v^C} x_{kj} \left(\frac{I_{vk}}{\pi_{vk}} - 1\right)$$

Logo,

$$\begin{aligned}
E\left\{\bar{y}_R^C - \mu_v | \mathbf{Y}_v\right\} & \leq \left\{E\left[\left(\frac{1}{N_v^C} \sum_{k \in U_v^C} y_k \left(\frac{I_{vk}}{\pi_{vk}} - 1\right)\right)^2 \middle| \mathbf{Y}_v\right]\right\}^{1/2} \\
& + \left\{E\left[\sum_{j=0}^q \hat{\beta}_{vj}^2 | \mathbf{Y}_v\right] \sum_{j=0}^q E\left[\left[\sum_{k \in U_v^C} x_{kj} \left(\frac{I_{vk}}{\pi_{vk}} - 1\right)\right]^2\right]\right\}^{1/2} \\
& = Q_1^{1/2} + Q_2^{1/2}.
\end{aligned}$$

Tem-se que

$$\begin{aligned}
 Q_1 &= E \left[\left(\frac{1}{N_v^C} \sum_{k \in U_v^C} y_k \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right) \right)^2 \middle| \mathbf{Y}_v \right] \\
 &= \frac{1}{N_v^2} E \left\{ \sum_{k \in U_v^C} y_k^2 \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right)^2 + 2 \sum_k \sum_{\substack{l \\ k \neq l}} y_k y_l \left(\frac{I_{vk} I_{vl}}{\pi_{vk} \pi_{vl}} - \frac{I_{vk}}{\pi_{vk}} - \frac{I_{vl}}{\pi_{vl}} + 1 \right) \right\} \\
 &= \frac{1}{N_v^2} \sum_{k \in U_v^C} y_k^2 \left(\frac{1}{\pi_{vk}} - 1 \right)^2 + \frac{2}{N_v^2} \sum_k \sum_{\substack{l \\ k \neq l}} y_k y_l \left(\frac{\pi_{vkl}}{\pi_{vk} \pi_{vl}} - 1 \right).
 \end{aligned}$$

Dessa forma,

$$\begin{aligned}
 \frac{1}{N_v^2} \sum_{k \in U_v^C} y_k^2 \left(\frac{1}{\pi_{vk}} - 1 \right)^2 &\leq \frac{N_v^{-1} \sum_{k \in U_v^C} y_k^2}{N_v \min_{1 \leq k \leq N_v} \pi_{vk}} \rightarrow 0, \text{ quando } v \rightarrow \infty \\
 \frac{2}{N_v^2} \sum_k \sum_{\substack{l \\ k \neq l}} y_k y_l \left(\frac{\pi_{vkl}}{\pi_{vk} \pi_{vl}} - 1 \right) &\leq \max_{1 \leq k \neq l \leq N_v} \left| \frac{\pi_{vkl}}{\pi_{vk} \pi_{vl}} - 1 \right| \left(\frac{\sum_{k \in U_v^C} |y_k|}{N_v} \right)^2 \\
 &\leq \max_{1 \leq k \neq l \leq N_v} \left| \frac{\pi_{vkl}}{\pi_{vk} \pi_{vl}} - 1 \right| \frac{\sum_{k \in U_v^C} y_k^2}{N_v} \rightarrow 0, \text{ quando } v \rightarrow \infty
 \end{aligned}$$

Analogamente, $Q_2 \rightarrow 0$, quando $v \rightarrow \infty$, pois

$$E \left\{ \left[\sum_{k \in U_v^C} x_{kj} \left(\frac{I_{vk}}{\pi_{vk}} - 1 \right) \right]^2 \right\} \rightarrow 0, \text{ quando } v \rightarrow \infty,$$

o que é garantido pelas condições **C6**, **C8**, **C9** e **C10**. Portanto, o estimador $\bar{y}_{RH(v)}$, dado em função de estimadores regressão para cadastros e domínios é assintoticamente centrado e consistente. ■

(ii) **Estimador** \bar{y}_{RFB} : É possível notar que através do método de linearização de Taylor, tem-se

que

$$\begin{aligned}
\bar{y}_{RFB(v)} &\approx (\hat{w}_{av}\mu_a + w_{av}\bar{y}_{RAv} - w_{av}\mu_a) + (\hat{w}_{bv}\mu_b + w_{bv}\bar{y}_{RBv} - w_{bv}\mu_b) \\
&+ p_1 \left(\hat{w}_{ab(A)v}\mu_{abv} + w_{abv}\bar{y}_{Rab(A)v} - w_{abv}\mu_{abv} \right) + (1 - p_1) (\hat{w}_{abv}\mu_{abv} + w_{abv}\mu_{abv}) \\
&+ N_v^{-1} p_2 (\hat{N}_{ab(A)v} - \hat{N}_{ab(B)v}) \\
&= w_{av}\bar{y}_{RAv} + w_{bv}\bar{y}_{RBv} + p_1 w^{abv} \bar{y}_{Rab(A)v} + (1 - p_1) w_{abv} \bar{y}_{Rab(B)v} \\
&+ \hat{w}_{av}\mu_{av} + \hat{w}_{bv}\mu_{bv} + p_1 \hat{w}_{abv}\mu_{abv} + (1 - p_1) \hat{w}_{abv}\mu_{abv} + N_v^{-1} p_2 (\hat{N}_{ab(A)} - \hat{N}_{ab(B)}) \\
&- w_{av}\mu_{av} - w_{bv}\mu_{bv} - w_{abv}\mu_{abv}
\end{aligned}$$

Admitindo que $\mu_v = w_{av}\mu_{av} + w_{bv}\mu_{bv} + p w_{abv}\mu_{abv} + (1 - p)w_{abv}\mu_{abv}$, tem-se que

$$\begin{aligned}
\bar{y}_{RFB(v)} - \mu_v &\approx w_a (\bar{y}_{RAv} - \mu_{av}) + w_b (\bar{y}_{RBv} - \mu_{bv}) + p_1 w_{abv} (\bar{y}_{Rab(A)v} - \mu_{abv}) + (1 - p_1) w_{abv} (\bar{y}_{Rab(B)v} - \mu_{abv}) \\
&+ N_v^{-1} (\hat{N}_a \mu_{av} + \hat{N}_b \mu_{bv} + p_1 \hat{N}_{ab}^A \mu_{abv} + (1 - p_1) \hat{N}_{ab}^B \mu_{abv} + p_2 \hat{N}_{ab}^A - p_2 \hat{N}_{ab}^B) \\
&= w_{av} (\bar{y}_{RAv} - \mu_{av}) + p_1 w_{abv} (\bar{y}_{Rab(A)v} - \mu_{abv}) \\
&+ w_{bv} (\bar{y}_{RBv} - \mu_{bv}) + (1 - p_1) w_{abv} (\bar{y}_{Rab(B)v} - \mu_{abv}) + \mathbf{O}(1).
\end{aligned}$$

Logo, a variância aproximada dessa diferença é dada por

$$\begin{aligned}
A\text{Var}(\bar{y}_{RFB(v)} - \mu_v) &= (w_{av})^2 A\text{Var}(\bar{y}_{RAv} - \mu_{av}) + (w_{bv})^2 A\text{Var}(\bar{y}_{RBv} - \mu_{bv}) \\
&+ (p_1 w_{abv})^2 A\text{Var}(\bar{y}_{Rab(A)v} - \mu_{abv}) + (1 - p_1)^2 (w_{abv})^2 A\text{Var}(\bar{y}_{Rab(B)v} - \mu_{abv}) \\
&+ 2p_1 w_{av} w_{abv} A\text{Cov} \left\{ (\bar{y}_{RAv} - \mu_{av}) (\bar{y}_{Rab(A)v} - \mu_{abv}) \right\}
\end{aligned}$$

$$\begin{aligned}
& + 2(1-p_1)w_{bv}w_{abv} \text{ACov} \left\{ (\bar{y}_{RBv} - \mu_{bv}) \left(\bar{y}_{Rab(B)v} - \mu_{abv} \right) \right\} \\
& + \mathbf{O}(N_v^{-1}) \\
& \leq (w_{av})^2 \text{AVar}(\bar{y}_{RAv} - \mu_{av}) + (w_{bv})^2 \text{AVar}(\bar{y}_{RBv} - \mu_{bv}) \\
& + (p_1 w_{abv})^2 \text{AVar}(\bar{y}_{Rab(A)v} - \mu_{abv}) + (1-p_1)^2 (w_{abv})^2 \text{AVar}(\bar{y}_{Rab(B)v} - \mu_{abv}) \\
& + 2p_1 w_{av} w_{abv} \text{AVar}(\bar{y}_{RAv} - \mu_{av}) \text{AVar}(\bar{y}_{Rab(A)v} - \mu_{abv}) \\
& + 2(1-p_1)w_{bv}w_{abv} \text{AVar}(\bar{y}_{RBv} - \mu_{bv}) \text{AVar}(\bar{y}_{Rab(B)v} - \mu_{abv}) + \mathbf{O}(N_v^{-1}) \\
& = \mathbf{O} \left\{ \max \left[(n_{av})^{-2\delta}, (n_{bv})^{-2\delta}, (n_v^{ab(A)})^{-2\delta}, (n_v^{ab(B)})^{-2\delta} \right] \right\} = \mathbf{O} \left\{ (n_v^*)^{-2\delta} \right\}
\end{aligned}$$

Além disso, para cadastros e domínios individualmente vale também os resultados apresentados Robinson & Särndal (1983), da mesma forma que a feita para o estimador $\bar{y}_{RH(v)}$. Portanto, sob a abordagem de cadastro duplo, o estimador $\bar{y}_{RFB(v)}$, função de estimadores regressão para cadastros e domínios também é assintoticamente centrado e consistente.

(iii) Estimador \bar{y}_{RBLG} :

Da mesma forma que as apresentadas para o teorema 3, são admitidas para os domínios as condições apresentadas por Isaki & Fuller (1982). Como a estratégia de Bankier é uma adaptação da estratégia de Horvitz-Thompson a partir da identificação dos elementos pertencentes aos domínios a , b e ab , tem-se que

$$\begin{aligned}
\bar{y}_{Rba(v)} &= (w_{av})\bar{y}_{RAv} + (w_{bv})\bar{y}_{Rbv} + w_{abv}\bar{y}_{Rabv} \\
\mu_v &= \frac{t_{yA} + t_{yB}}{N_v} = \frac{t_{ya} + t_{yab(A)} + t_{yb} + t_{yab(B)}}{N} = \frac{t_{ya} + t_{yb} + t_{yab}}{N}, \quad \text{onde } t_{yab} = t_{yab(A)} + t_{yab(B)} \\
\mu_v &= (w_{av})\mu_{av} + (w_{bv})\mu_{bv} + w_{abv}\mu_{abv}
\end{aligned}$$

Dessa forma, tem-se que

$$\begin{aligned}
\bar{y}_{RBa(v)} - \mu_v &= (w_{av})(\bar{y}_{Rav} - \mu_{av}) + (w_{bv})(\bar{y}_{Rbv} - \mu_{bv}) + w_{abv}(\bar{y}_{Rabv} - \mu_{abv}) \\
\text{AVar}(\bar{y}_{RBa(v)} - \mu_v) &= (w_{av})^2 \text{AVar}(\bar{y}_{Rav} - \mu_{av}) + (w_{bv})^2 \text{AVar}(\bar{y}_{Rbv} - \mu_{bv}) \\
&+ (w_{abv})^2 \text{AVar}(\bar{y}_{Rabv} - \mu_{abv}) \\
&+ 2(w_{av})w_{abv} \text{ACov}[(\bar{y}_{Rav} - \mu_{av}) ; (\bar{y}_{Rabv} - \mu_{abv})] \\
&+ 2(w_{bv})w_{abv} \text{ACov}[(\bar{y}_{Rbv} - \mu_{bv}) ; (\bar{y}_{Rabv} - \mu_{abv})] \\
&\leq (w_{av})^2 \text{AVar}(\bar{y}_{Rav} - \mu_{av}) + (w_{bv})^2 \text{AVar}(\bar{y}_{Rbv} - \mu_{bv}) \\
&+ (w_{abv})^2 \text{AVar}(\bar{y}_{Rabv} - \mu_{abv}) \\
&+ 2(w_{av})w_{abv} \{ \text{ACov}[(\bar{y}_{Rav} - \mu_{av}) ; (\bar{y}_{Rabv} - \mu_{abv})] \}^2 \\
&+ 2(w_{bv})w_{abv} \{ \text{ACov}[(\bar{y}_{Rbv} - \mu_{bv}) ; (\bar{y}_{Rabv} - \mu_{abv})] \}^2 \\
&\leq (w_{av})^2 \text{AVar}(\bar{y}_{Rav} - \mu_{av}) + (w_{bv})^2 \text{AVar}(\bar{y}_{Rbv} - \mu_{bv}) \\
&+ (w_{abv})^2 \text{AVar}(\bar{y}_{Rabv} - \mu_{abv}) \\
&+ 2(w_{av})w_{abv} \{ \text{AVar}(\bar{y}_{Rav} - \mu_{av}) \text{AVar}(\bar{y}_{Rabv} - \mu_{abv}) \} \\
&+ 2(w_{bv})w_{abv} \{ \text{AVar}(\bar{y}_{Rbv} - \mu_{bv}) \text{AVar}(\bar{y}_{Rabv} - \mu_{abv}) \} \\
&= \mathbf{O}(n_{a(v)}^{-2\delta}) + \mathbf{O}(n_{bv}^{-2\delta}) + \mathbf{O}(n_v^{ab})^{-2\delta}) + \mathbf{O}(n_{av})^{-2\delta}(n_v^{ab})^{-2\delta}) \\
&+ \mathbf{O}(n_{bv})^{-2\delta}(n_v^{ab})^{-2\delta}) = \mathbf{O}\left(\max\left\{(n_{av})^{-2\delta}, (n_{bv})^{-2\delta}, (n_v^{ab})^{-2\delta}\right\}\right),
\end{aligned}$$

Adaptando os resultados apresentados por Isaki & Fuller (1982) e Robinson & Särndal (1983) para as quantidades referentes à cadastros e domínios, da mesma forma que a feita para o estimador $\bar{y}_{RH(v)}$, tem-se que o estimador $\bar{y}_{RBa(v)}$ é assintoticamente centrado e consistente.

(iv) Estimador $\bar{y}_{RBLG|(\infty)}$:

Inicialmente, sejam $\omega_{ka}^{(\infty)}$, $\omega_{kb}^{(\infty)}$, $\omega_{kab}^{(\infty)}$ dados da mesma forma que a apresentada por Skinner (1991):

$$\begin{aligned}\omega_{ka}^{(\infty)} &= \frac{1}{\pi_{vk}^A} \left[\prod_{r \text{ par}} \left(\frac{N_{Av}}{\hat{N}^{A(r-1)}} \right) \right] = \frac{1}{\pi_{vk}^A} \alpha^A \\ \omega_{kb}^{(\infty)} &= \frac{1}{\pi_{vk}^B} \left[\prod_{r \text{ ímpar}} \left(\frac{N_{Bv}}{\hat{N}^{B(r-1)}} \right) \right] = \frac{1}{\pi_{vk}^B} \alpha^B \\ \omega_{kab}^{(\infty)} &= \frac{1}{\pi_{vk}^A + \pi_{vk}^B} \left[\prod_{r \text{ par}} \left(\frac{N_{Av}}{\hat{N}^{A(r-1)}} \right) \prod_{r \text{ ímpar}} \left(\frac{N_{Bv}}{\hat{N}^{B(r-1)}} \right) \right] = \left(\frac{1}{\pi_{vk}^A + \pi_{vk}^B} \right) \alpha^A \alpha^B\end{aligned}$$

α^A e α^B convergem para valores fixos quanto $r \rightarrow \infty$. Dessa forma, basta substituir as quantidades apresentadas em (3.10) e temos o resultado, da mesma forma que a feita para o estimador $\bar{y}_{RBa(v)}$.

(v) Estimador \bar{y}_{RPML} :

É possível notar que

$$\begin{aligned}
 \bar{y}_{RPML(v)} &\approx w_{av}\bar{y}_{Ra(v)} + w_{a(*)}\mu_{av} - w_{av}\mu_{av} + w_{bv}\bar{y}_{Rb(v)} + w_{abv(*)}\mu_{Bv} - w_v^b\mu_{Bv} \\
 &+ w_{abv} \left\{ \left(\frac{\pi_{vk}^A}{\pi_{vk}^A + \pi_{vk}^B} \right) \bar{y}_{Rab(A)v} + \left(\frac{\pi_{vk}^B}{\pi_{vk}^A + \pi_{vk}^B} \right) \bar{y}_{Rab(B)v} \right\} \\
 &+ w^{ab(*)}\mu_{abv} - w_{abv}\mu_{abv} \\
 &= w_{av}\bar{y}_{Ra(v)} + w_{bv}\bar{y}_{Rb(v)} \\
 &+ w_{abv} \left\{ \left(\frac{\pi_{vk}^A}{\pi_{vk}^A + \pi_{vk}^B} \right) \bar{y}_{Rab(A)v} + \left(\frac{\pi_{vk}^B}{\pi_{vk}^A + \pi_{vk}^B} \right) \bar{y}_{Rab(B)v} \right\} + O(1)
 \end{aligned}$$

Dessa forma, admitindo que $\mu_v = w_{av}\mu_{av} + w_{bv}\mu_b + w_{abv}\mu_{abv}$, e que

$$\begin{aligned}
 w_{abv} \left(\frac{\pi_{vk}^A}{\pi_{vk}^A + \pi_{vk}^B} \right) \bar{y}_{Rab(A)v} + w_{abv} \left(\frac{\pi_{vk}^B}{\pi_{vk}^A + \pi_{vk}^B} \right) \bar{y}_{Rab(B)v} &= \frac{w_{abv}}{\pi_{vk}^A + \pi_{vk}^B} \left(\pi_{vk}^A \bar{y}_{Rab(A)v} + \pi_{vk}^B \bar{y}_{Rab(B)v} \right) \\
 &= w_v^{ab(*)} \bar{y}^{ab(*)},
 \end{aligned}$$

tem-se que

$$\bar{y}_{RPML(v)} - \mu_v \approx w_{av} \left(\bar{y}_{Ra(v)} - \mu_{av} \right) + w_{bv} \left(\bar{y}_{Rb(v)} - \mu_{bv} \right) + w_{abv(*)} \left(\bar{y}_{Rab(v)} - \mu_{abv} \right) + O(1)$$

Dessa forma, a variância aproximada dessa diferença é dada por

$$\begin{aligned}
 AVar \left(\bar{y}_{RPML(v)} - \mu_v \right) &= (w_{av})^2 AVar \left(\bar{y}_{Ra(v)} - \mu_{av} \right) + (w_{bv})^2 AVar \left(\bar{y}_{Rb(v)} - \mu_{bv} \right) \\
 &+ (w_{abv(*)}) AVar \left(\bar{y}_{Rab(v)} - \mu_{abv} \right) + O(1) \\
 &= \mathbf{O} \left(\max \left\{ (n_{av})^{-2\delta}, (n_{bv})^{-2\delta}, (n_v^{ab(A)})^{-2\delta}, (n_v^{ab(B)})^{-2\delta} \right\} \right) = \mathbf{O} \left((n_v^*)^{-2\delta} \right),
 \end{aligned}$$

de acordo com os resultados obtidos Isaki & Fuller (1982) e Robinson & Särndal (1983) e apresentados anteriormente para os outros estimadores propostos.

Apêndice C

Programas utilizados

As avaliações computacionais realizadas ao longo desta tese foram feitas utilizando o ambiente de programação, análise de dados e gráficos R em sua versão 2.12.0, e que se encontra gratuitamente disponível em <http://www.r-project.org>. A opção por esta linguagem foi devido ao compromisso entre a flexibilidade oferecida por algumas linguagens compiladas e a conveniência dos tradicionais pacotes estatísticos, além da facilidade de implementação dos estimadores propostos. A seguir será apresentado o programa utilizado para a avaliação dos estimadores sob o plano de amostragem aleatória simples.

```
#####
#UNIVERSIDADE FEDERAL DE PERNAMBUCO
#CENTRO DE CIÊNCIAS EXATAS E DA NATUREZA
#PROGRAMA DE PÓS-GRADUAÇÃO EM ESTATÍSTICA
#DOUTORADO EM ESTATÍSTICA
#####
#INFERÊNCIA SOB PLANOS AMOSTRAIS DE CADASTROS DUPLOS
#####
#HEMÍLIO FERNANDES CAMPOS COELHO
#####

#####
#O PROGRAMA A SEGUIR FOI UTILIZADO PARA A APLICAÇÃO APRESENTADA NO CAPÍTULO 2
#####

#####
#Função para criar uma população artificial de tamanho N baseada em dados de uma
#amostra de tamanho n.
#####
# name -> Nome que a variável receberá
#####

artipop<-function(N,y,name)
{
y<-as.vector(unlist(y))
n<-length(y)
nv<-log(N/n)/(log(2))
nvr<-floor(nv+1)
```

```

for(k in 1:nvr){
y<-append(y,y, after=n)}
var<-data.frame(y[1:N])
names(var)<-c(name)
return(var)
}

#####
#Função para cálculo do p-otimo do estimador de Hartley caso normal.
#####
# N_A -> Tamanho da população do cadastro A
# NB -> Tamanho da população do cadastro B
# Nab -> Tamanho da população do domínio "ab"
# yabC -> Variável resposta na amostra para o domínio ab do cadastro C
#####

pH=function(N_A,NB,Nab,yabA,yabB)
{
var1=var(yabA)
var2=var(yabB)
fAinv=N_A/nA
fBinv=NB/nB
wabA=Nab/N_A
wabB=Nab/NB
a=(2*NB*fBinv*wabB*var2)/(2*N_A*fAinv*wabA*var1+2*NB*fBinv*wabB*var2)
return(a)
}

#####
#Função que calcula o estimador de Hartley
#####
# p -> Função "pH" aplicada anteriormente
# Na -> Tamanho da população do domínio "a"
# Nb -> Tamanho da população do domínio "b"
# yd -> Variável de interesse para o domínio d
# ydC -> Variável de interesse para o domínio d, cadastro C
#####

esthart=function(p,Na,Nb,ya,yb,yabA,yabB)
{
est=Na*mean(ya)+Nb*mean(yb)+Nab*p*mean(yabA)
+Nab*(1-p)*mean(yabB)
return(est)
}

##### CONSTRUÇÃO DO BANCO DE DADOS #####
library(TeachingSampling)
#Cadastro de Hortaliças (Cadastro B)
cadhor<- read.table("marco hortalizas.tab", header=TRUE, sep=" ", na.strings="NA", dec=".", strip.white=TRUE)

#Dados de área cultivada, área colhida e produção (variáveis de interesse)
acult<- read.table("asemh.tab", header=TRUE, sep=" ", na.strings="NA", dec=".", strip.white=TRUE)
acolhi<- read.table("acosh.tab", header=TRUE, sep=" ", na.strings="NA", dec=".", strip.white=TRUE)

```

```
producao<- read.table("prodh.tab", header=TRUE, sep=" ", na.strings="NA", dec=".", strip.white=TRUE)

#Cadastro do DANE (Cadastro A)
cadENA <- read.table("universo1.tab", header=TRUE, sep=" ", na.strings="NA", dec=".", strip.white=TRUE)

#Criação de um cadastro geral
N=nrow(cadENA)
Nab= 10000
Na=trunc(nrow(cadENA)/2)-Nab
Nb=trunc(nrow(cadENA)/2+1)-Nab
N_A=Na+Nab
N_B=Nb+Nab
domA=append(matrix(1,Na),matrix(0,Nab))
domB=append(matrix(0,Nab),matrix(1,Nb))
cult= artipop(N,acult,"area cultivada")
colhi=artipop(N,acolhi,"area colhida")
prod=artipop(N,producao,"produção")
cadgeral=data.frame(cadENA,cult,colhi,prod)

####CADASTROS
cadA=data.frame(cadgeral[1:N_A,],domA)
cadB=data.frame(cadgeral[(N_A+1-Nab):(N_A+N_B-Nab),],domB)
domab=subset(cadA,domA==0)

####PARÂMETROS DE INTERESSE
media_tot=(sum(cadA$produção)+sum(cadB$produção)-sum(domab$produção))/N

####COLETA DAS INFORMAÇÕES DA AMOSTRA
nA=1000
nB=800
selA=S.SI(N_A,nA)
amoA=cadA[selA,]
nabA=nrow(domabA)
selB=S.SI(N_B,nB)
amoB=cadB[selB,]

####DOMÍNIOS
doma=subset(amoA,domA==1)
domabA=subset(amoA,domA==0)
domb=subset(amoB,domB==1)
domabB=subset(amoB,domB==0)

####CÁLCULO DO ESTIMADOR DE HARTLEY
phart=pH(N_A,N_B,Nab,domabA$produção,domabB$produção)
hartley=esthart(phart,Na,Nb,doma$produção,domb$produção,domabA$produção,domabB$produção)/N
mean(doma$produção)
mean(domb$produção)
mean(domabA$produção)
mean(domabB$produção)
```

```
#####
# SIMULAÇÃO:
#O PROGRAMA A SEGUIR AVALIA OS ESTIMADORES DE HARTLEY, FULLER, BANKIER,
#PML, E OS ESTIMADORES PROPOSTOS ATRAVÉS DO USO DE ESTIMADORES GREG
#####

#####
#Função para cálculo do p-otimo do estimador de Hartley caso normal.
#####
# nA -> Tamanho da amostra obtida do cadastro A
# N_A -> Tamanho da população do cadastro A
# nB -> Tamanho da amostra obtida do cadastro B
# NB -> Tamanho da população do cadastro B
# sA -> Amostra obtida do cadastro A
# sB -> Amostra obtida do cadastro B
# Nab -> Tamanho da população do domínio "ab"
# dab -> Indicador numérico do domínio "ab"
#####

pH=function(nA,N_A,nB,NB,sA,sB,Nab,dab)
{
  domabA=subset(sA,sA[,2]==dab)
  doma=subset(sA,sA[,2]==(1-dab))
  domabB=subset(sB,sB[,2]==dab)
  domb=subset(sB,sB[,2]==(1-dab))
  var1=var(domabA[,1])
  var2=var(domabB[,1])
  fAinv=N_A/nA
  fBinv=NB/nB
  wabA=Nab/N_A
  wabB=Nab/NB
  a=(2*NB*fBinv*wabB*var2)/(2*N_A*fAinv*wabA*var1+2*NB*fBinv*wabB*var2)
  return(a)
}

#####
#Função que calcula o estimador de Hartley
#####
# p -> Função "pH" aplicada anteriormente
# nA -> Tamanho da amostra obtida do cadastro A
# N_A -> Tamanho da população do cadastro A
# nB -> Tamanho da amostra obtida do cadastro B
# NB -> Tamanho da população do cadastro B
# sA -> Amostra obtida do cadastro A
# sB -> Amostra obtida do cadastro B
# Nab -> Tamanho da população do domínio "ab"
# dab -> Indicador numérico do domínio "ab"
#####
```

```

esthart=function(p,n_A,nB,N_A,NB,sA,sB,Nab,dab)
{
domabA=subset(sA,sA[,2]==dab)
doma=subset(sA,sA[,2]==(1-dab))
domabB=subset(sB,sB[,2]==dab)
domb=subset(sB,sB[,2]==(1-dab))
Na=N_A-Nab
Nb=NB-Nab
est=Na*mean(doma[,4])+Nb*mean(domb[,4])+Nab*p*mean(domabA[,4])
+Nab*(1-p)*mean(domabB[,4])
return(est)
}

#####
#Função que estima Nab para o estimador de Fuller-Burmeister
#####
# N_A -> Tamanho da população do cadastro A
# NB -> Tamanho da população do cadastro B
# sA -> Amostra obtida do cadastro A
# sB -> Amostra obtida do cadastro B
# dA -> Variável indicadora de domínios (1 se pertence a "a"; 0 se pertence a "ab")
# dB -> Variável indicadora de domínios (1 se pertence a "b"; 0 se pertence a "ab")
# dab -> Indicador numérico do domínio "ab"
#####

Nabful<- function(N_A,NB,sA,sB,dA,dB,dab)
{
domabA<-subset(sA,dA==dab)
domabB<-subset(sB,dB==dab)
nabA<-length(domabA[,1])
nabB<-length(domabB[,1])
nA<-length(sA[,1])
nB<-length(sB[,1])
gA<-(N_A - nabA)/(N_A - 1)
gB<-(NB - nabB)/(NB - 1)
z3 <-(nA*gB)+(nB*gA)
z2 <-nA*nB*gB+nB*N_A*gA+nabA*N_A*gB+nabB*NB*gA
z1 <-(nabA*gB+nabB*gA)*N_A*NB
if(missing(z1) || is.null(z1) || is.na(z1) || !is.numeric(z1)){z1<-0}
if(missing(z2) || is.null(z2) || is.na(z2) || !is.numeric(z2)){z2<-0}
if(missing(z3) || is.null(z3) || is.na(z3) || !is.numeric(z3)){z3<-0}
roots<-polyroot(c(z1,-z2,z3))
NabFB <- min(Re(roots))
return(NabFB)
}

#####
#Função que calcula o estimador de Fuller & Burmeister
#####
#Precisa da função "Nabful"
#####
# nA -> Tamanho da amostra obtida do cadastro A
# N_A -> Tamanho da população do cadastro A

```

```

# nB -> Tamanho da amostra obtida do cadastro B
# NB -> Tamanho da população do cadastro B
# sA -> Amostra obtida do cadastro A
# sB -> Amostra obtida do cadastro B
# Nabf -> Tamanho da população do domínio "ab", estimada pela função "Nabful"
# dab -> Indicador numérico do domínio "ab"
#####

estful=function(n_A,nB,N_A,NB,sA,sB,Nabf,dab)
{
  domabA=subset(sA,sA[,2]==dab)
  doma=subset(sA,sA[,2]==(1-dab))
  domabB=subset(sB,sB[,2]==dab)
  domb=subset(sB,sB[,2]==(1-dab))
  nabA=length(domabA[,1])
  nabB=length(domabB[,1])
  fA=n_A/N_A
  fB=nB/NB
  w=nabA*(1-fB)/(nabA*(1-fB)+nabB*(1-fA))
  yabest=w*mean(domabA[,4])+(1-w)*mean(domabB[,4])
  Naest=(N_A-Nabf)
  Nbest=(NB-Nabf)
  estful=Naest*mean(doma[,4])+Nbest*mean(domb[,4])+Nabf*yabest
  return(estful)
}

#####
#Função que calcula o estimador de Bankier
#####
# nA -> Tamanho da amostra obtida do cadastro A
# N_A -> Tamanho da população do cadastro A
# nB -> Tamanho da amostra obtida do cadastro B
# NB -> Tamanho da população do cadastro B
# sA -> Amostra obtida do cadastro A
# sB -> Amostra obtida do cadastro B
# dab -> Indicador numérico do domínio "ab"
#####

estbank=function(n_A,nB,N_A,NB,sA,sB,dab)
{
  domabA=subset(sA,sA[,2]==dab)
  doma=subset(sA,sA[,2]==(1-dab))
  domabB=subset(sB,sB[,2]==dab)
  domb=subset(sB,sB[,2]==(1-dab))
  fA=n_A/N_A
  fB=nB/NB
  somapondA=(1/fA)*sum(doma[,4])
  somapondabA=(1/(fA+fB))*sum(domabA[,4])
  soma1=somapondA+somapondabA
  somapondB=(1/fB)*sum(domb[,4])
  somapondabB=(1/(fA+fB))*sum(domabB[,4])
  soma2=somapondB+somapondabB
  est=soma1+soma2
  return(est)
}

```

```
}

```

```
#####
#Função que calcula o estimador de Horvitz-Thompson para o caso do estimador Bankier Combinado
#####
#Esta função é uma alteração da função "E.SI" do pacote "TeachingSampling"
#####
E.SImod=function(N_A,NB, n_A,nB,nd,y)
{
y <- as.data.frame(y)
val1=n_A/N_A
val2=nB/NB
pik <- rep(1/(val1+val2),nd)
dk <- 1/pik
ty <- sum(y * dk)
return(ty)
}

```

```
#####
#Função que calcula o estimador de Nab para a estratégia PML
#####
# nA -> Tamanho da amostra obtida do cadastro A
# N_A -> Tamanho da população do cadastro A
# nB -> Tamanho da amostra obtida do cadastro B
# NB -> Tamanho da população do cadastro B
# sA -> Amostra obtida do cadastro A
# sB -> Amostra obtida do cadastro B
# dab -> Indicador numérico do domínio "ab"
#####

```

```
NabPML = function(n_A,nB,N_A,NB,sA,sB,dab)
{
domabA=subset(sA,sA[,2]==dab)
doma=subset(sA,sA[,2]==(1-dab))
domabB=subset(sB,sB[,2]==dab)
domb=subset(sB,sB[,2]==(1-dab))
fA=n_A/N_A
fB=nB/NB
vetor1=rep(1,length(domabA[,4]))
vetor2=rep(1,length(domabB[,4]))
nabA=sum(vetor1)
nabB=sum(vetor2)
NabA=(1/fA)*nabA
NabB=(1/fB)*nabB
z1=z2=z3=0
z1 = n_A*NabA*NB+nB*NabB*N_A
z2 = n_A*NB+nB*N_A+n_A*NabA+nB*NabB
z3 = n_A+nB
if(missing(z1) || is.null(z1) || is.na(z1) || !is.numeric(z1)){z1<-0}
if(missing(z2) || is.null(z2) || is.na(z2) || !is.numeric(z2)){z2<-0}
if(missing(z3) || is.null(z3) || is.na(z3) || !is.numeric(z3)){z3<-0}
roots<-polyroot(c(z1,-z2,z3))

```

```

NabPML = min(Re(roots))
return(NabPML)
}

#####
#Função que calcula o estimador para o total via PML
#####
# nA -> Tamanho da amostra obtida do cadastro A
# N_A -> Tamanho da população do cadastro A
# nB -> Tamanho da amostra obtida do cadastro B
# NB -> Tamanho da população do cadastro B
# sA -> Amostra obtida do cadastro A
# sB -> Amostra obtida do cadastro B
# dab -> Indicador numérico do domínio "ab"
#####

estPML = function(n_A,nB,N_A,NB,Nab,sA,sB,dab)
{
fA=n_A/N_A
fB=nB/NB
domabA=subset(sA,sA[,2]==dab)
doma=subset(sA,sA[,2]==(1-dab))
domabB=subset(sB,sB[,2]==dab)
domb=subset(sB,sB[,2]==(1-dab))
vetor1=rep(1,length(domabA[,1]))
vetor2=rep(1,length(domabB[,1]))
NabA=(1/fA)*sum(vetor1)
NabB=(1/fB)*sum(vetor2)
valor1=fA*NabA*mean(domabA[,4])+fB*NabB*mean(domabB[,4])
valor2=fA*NabA+fB*NabB
muabest=valor1/valor2
est=(N_A-Nab)*mean(doma[,4])+(NB-Nab)*mean(domb[,4])+Nab*muabest
return(est)
}

#####
#Função que disponibiliza o vetor de resíduos e Função que calcula a covariância estimada entre
#estimadores GREG
#####
# n -> Tamanho da amostra no cadastro de interesse
# N -> Tamanho da população no cadastro de interesse
# y -> Vetor com as informações da variável de interesse
# x -> Vetor ou matriz com informações das variáveis auxiliares
# tx -> Vetor com os totais para as variáveis auxiliares
# ck -> Por default, igual a 1. Vetor de pesos referente à variância do modelo
#####

residuo=function (N, n, y, x, tx, b, b0 = FALSE)
{
y <- as.data.frame(y)
x <- as.matrix(x)
pik <- rep(n/N, n)
dk <- 1/pik
if (b0 == TRUE) {

```

```

x <- as.matrix(cbind(1, x))
}
for (k in 1:dim(y)[2]) {
e <- y[, k] - (x %>% as.matrix(b[, k]))
}
return(e)
}

covarg=function(n,N,e1,e2)
{
a=n/(N^2)
b=n/N
c=(n-1)/(N-1)
V1=1-a/(b*c)
covar=V1*sum(e1%>%t(e2))
return(covar)
}

#####
#Função que calcula o p-ótimo para o estimador regressão com base na estratégia de Hartley
#####
#-> Precisa da função "covarg"
#-> Precisa da função "residuo"
#####
# N -> Total populacional referente a toda população
# N_A -> Total populacional referente a A
# NB -> Total populacional referente a B
# nA -> Tamanho da amostra referente ao cadastro A
# nB -> Tamanho da amostra referente ao cadastro B
# Nab -> Total populacional referente ao domínio ab
# nabA -> Tamanho da amostra em ab, coletada de A
# nabB -> Tamanho da amostra em ab, coletada de B
# eA -> Vetor de resíduos para o estimador regressão referente a A
# eB -> Vetor de resíduos para o estimador regressão referente a B
# eabA -> Vetor de resíduos para o estimador regressão referente a ab(cadastro A)
# eabB -> Vetor de resíduos para o estimador regressão referente a ab(cadastro B)
# vabA -> Variância estimada do estimador regressão para ab (cadastro A)
# vabB -> Variância estimada do estimador regressão para ab (cadastro B)
# covabA-> Covariância estimada entre o estimador para A e o estimador para ab(cadastro A)
# covabB-> Covariância estimada entre o estimador para B e o estimador para ab(cadastro B)
#####

preghart=function(N,nA,nB, N_A,NB,nabA,nabB,Nab,eA,eB,eabA,eabB)
{
wab=Nab/N
wA=N_A/N
wB=NB/N
q1=(Nab^2)*(1-(nabA/Nab))*var(eabA)/(nabA)
q2=(Nab^2)*(1-(nabB/Nab))*var(eabB)/(nabB)
q3=covarg(nA,N_A,eA,eabA)
q4=covarg(nB,NB,eB,eabB)
num=wA*q3-wB*q4
den=wab*(q1-q2)
p=num/den
}

```

```

return(p)
}

#####
#Função que calcula o estimador greg para um domínio
#####
#obs: Desempenha o mesmo papel que a função GREG.SI, apenas foi modificada para cálculo correto
#dentro dos domínios
#####

GREG.SIESPd=function (N, n, nd, y, x, tx, b, b0 = FALSE)
{
y <- as.data.frame(y)
x <- as.matrix(x)
pik <- rep(n/N, nd)
dk <- 1/pik
if (b0 == TRUE) {
x <- as.matrix(cbind(1, x))
}
Total <- matrix(NA, nrow = 3, ncol = dim(y)[2])
rownames(Total) = c("Estimation", "Variance", "CVE")
colnames(Total) <- names(y)
for (k in 1:dim(y)[2]) {
xHT <- t(x) %*% dk
yHT <- sum(y[, k] * dk)
ty <- yHT + (tx - t(xHT)) %*% as.matrix(b[, k])
e <- y[, k] - (x %*% as.matrix(b[, k]))
Vty <- (N^2) * (1 - (n/N)) * var(e)/(n)
CVE <- 100 * sqrt(Vty)/ty
Total[, k] <- c(ty, Vty, CVE)
}
return(Total)
}

#####
#Função que calcula o estimador regressão de acordo com a estratégia de Hartley
#####
#Precisa da função "preghart"
#####
# p -> Valor de "p" obtido pela função "preghart"
# amo1 -> Amostra de A
# amo2 -> Amostra de B
# pop1 -> População do cadastro A
# pop2 -> População do cadastro B
# N_A -> Tamanho da população do cadastro A
# nA -> Tamanho da amostra obtida do cadastro A
# NB -> Tamanho da população do cadastro B
# nB -> Tamanho da amostra obtida do cadastro B
# Nab -> Tamanho da população do domínio "ab"
# dab -> Identificador do domínio ab
# name1 -> Identificador para a existência de intercepto no modelo para a pop. de A
# ("TRUE" ou "FALSE")

```

```

# name2 -> Identificador para a existência de intercepto no modelo para a pop. de B
# ("TRUE" ou "FALSE")
#####

GREG.SIH=function(p,pop1,pop2, amo1,amo2,N_A, nA, NB,nB,Nab,dab,name1,name2)
{
  amo1 <- as.data.frame(amo1)
  amo2 <- as.data.frame(amo2)
  popxA=data.frame(pop1[,1],pop1[,3])
  txA=c(length(popxA[,1]),sum(pop1[,1]),sum(pop1[,3]))
  popa=subset(pop1,pop1[,2]==(1-dab))
  popxa=data.frame(popa[,1],popa[,3])
  txA=c(length(popa[,1]),sum(popa[,1]),sum(popa[,3]))
  popxB=data.frame(pop2[,1],pop2[,3])
  txB=c(length(popxB[,1]),sum(pop2[,1]),sum(pop2[,3]))
  popb=subset(pop2,pop2[,2]==(1-dab))
  popxb=data.frame(popb[,1],popb[,3])
  txb=c(length(popb[,1]),sum(popb[,1]),sum(popb[,3]))
  xA=data.frame(amo1[,1],amo1[,3])
  xB=data.frame(amo2[,1],amo2[,3])
  Na=length(popa[,1])
  Nb=length(popb[,1])
  fA=nA/N_A
  fB=nB/NB
  pikA=rep(fA,nA)
  pikB=rep(fB,nB)
  popabA=subset(pop1,pop1[,2]==dab)
  popabB=subset(pop2,pop2[,2]==dab)
  popxabA=data.frame(popabA[,1],popabA[,3])
  popxabB=data.frame(popabB[,1],popabB[,3])
  domabA=subset(amo1,amo1[,2]==dab)
  xabA=data.frame(domabA[,1],domabA[,3])
  txabA=c(Nab,sum(popabA[,1]),sum(popabA[,3]))
  doma=subset(amo1,amo1[,2]==(1-dab))
  xa=data.frame(doma[,1],doma[,3])
  domabB=subset(amo2,amo2[,2]==dab)
  xabB=data.frame(domabB[,1],domabB[,3])
  txabB=c(Nab,sum(popabB[,1]),sum(popabB[,3]))
  domb=subset(amo2,amo2[,2]==(1-dab))
  xb=data.frame(domb[,1],domb[,3])
  pika=rep(fA,length(doma[,1]))
  pikb=rep(fB,length(domb[,1]))
  pikabA=rep(fA,length(domabA[,1]))
  pikabB=rep(fB,length(domabB[,1]))
  na=length(doma[,1])
  nb=length(domb[,1])
  ba=E.Beta(doma[,4],xa,pika,1,b0=name1)
  bb=E.Beta(domb[,4],xb,pikb,1,b0=name2)
  babA=E.Beta(domabA[,4],xabA,pikabA,1,b0=name1)
  babB=E.Beta(domabB[,4],xabB,pikabB,1,b0=name2)
  somapondA=GREG.SIESPd(N_A,nA,length(doma[,4]),doma[,4],xa,txA,ba,b0=name1)
  somapondAB=GREG.SIESPd(N_A,nA,length(domabA[,1]),domabA[,4],xabA,txabA,babA,b0=name2)
  somapondB=GREG.SIESPd(NB,nB,length(domb[,4]),domb[,4],xb,txb,bb,b0=name2)
  somapondABB=GREG.SIESPd(NB,nB,length(domabB[,1]),domabB[,4],xabB,txabB,babB,b0=name2)
}

```

```

esthart=somapondA[1]+p*somapondabA[1]+somapondB[1] + (1-p)*somapondabB[1]
return(esthart[1])
}

```

```

#####
#Função que calcula o estimador regressão de acordo com a estratégia de Bankier
#####
# Desempenha mesmo papel que a função GREG.SI. Aqui, a função foi modificada para poder
# receber os pesos conforme determina a estratégia de Bankier. Precisa da função "GREG.SIESPab"
#####
# amo1 -> Amostra de A
# amo2 -> Amostra de B
# pop1 -> População do cadastro A
# pop2 -> População do cadastro B
# N_A -> Tamanho da população do cadastro A
# nA -> Tamanho da amostra obtida do cadastro A
# NB -> Tamanho da população do cadastro B
# nB -> Tamanho da amostra obtida do cadastro B
# Nab -> Tamanho da população do domínio "ab"
# dab -> Identificador do domínio ab
# name1 -> Identificador para a existência de intercepto no modelo para a pop. de A ("TRUE" ou "FALSE")
# name2 -> Identificador para a existência de intercepto no modelo para a pop. de B ("TRUE" ou "FALSE")
#####
#-> Utiliza a primeira função a ser apresentada, a "GREG.SIESPab" (Função que calcula o
# estimador greg para o domínio ab sob a estratégia de bankier)
#-> Precisa da função "GREG.SIESPd" para os domínios
#####

GREG.SIESPab=function (N, n, N1,n1, nab,y, x, tx, b, b0 = FALSE)
{
y <- as.data.frame(y)
x <- as.matrix(x)
val1=n1/N1
val2=n/N
pik <- rep(1/(val1+val2), nab)
dk <- 1/pik
if (b0 == TRUE) {
x <- as.matrix(cbind(1, x))
}
Total <- matrix(NA, nrow = 3, ncol = dim(y)[2])
rownames(Total) = c("Estimation", "Variance", "CVE")
colnames(Total) <- names(y)
for (k in 1:dim(y)[2]) {
xHT <- t(x) %*% dk
yHT <- sum(y[, k] * dk)
ty <- yHT + (tx - t(xHT)) %*% as.matrix(b[, k])
e <- y[, k] - (x %*% as.matrix(b[, k]))
Vty <- (N^2) * (1 - (n/N)) * var(e)/(n)
CVE <- 100 * sqrt(Vty)/ty
Total[, k] <- c(ty, Vty, CVE)
}
return(Total)
}

```

```

GREG.SIBA=function(pop1,pop2,amo1,amo2,N_A, nA, NB, nB,Nab,dab,name1,name2)
{
  amo1 <- as.data.frame(amo1)
  amo2 <- as.data.frame(amo2)
  popxA=data.frame(pop1[,1],pop1[,3])
  txA=c(length(popxA[,1]),sum(pop1[,1]),sum(pop1[,3]))
  popa=subset(pop1,pop1[,2]==(1-dab))
  popxa=data.frame(popa[,1],popa[,3])
  txA=c(length(popa[,1]),sum(popa[,1]),sum(popa[,3]))
  popxB=data.frame(pop2[,1],pop2[,3])
  txB=c(length(popxB[,1]),sum(pop2[,1]),sum(pop2[,3]))
  popb=subset(pop2,pop2[,2]==(1-dab))
  popxb=data.frame(popb[,1],popb[,3])
  txB=c(length(popb[,1]),sum(popb[,1]),sum(popb[,3]))
  xA=data.frame(amo1[,1],amo1[,3])
  xB=data.frame(amo2[,1],amo2[,3])
  Na=length(popa[,1])
  Nb=length(popb[,1])
  fA=Na/N_A
  fB=Nb/NB
  pikA=rep(fA,nA)
  pikB=rep(fB,nB)
  popabA=subset(pop1,pop1[,2]==dab)
  popabB=subset(pop2,pop2[,2]==dab)
  popxabA=data.frame(popabA[,1],popabA[,3])
  popxabB=data.frame(popabB[,1],popabB[,3])
  domabA=subset(amo1,amo1[,2]==dab)
  xabA=data.frame(domabA[,1],domabA[,3])
  txabA=c(Nab,sum(popabA[,1]),sum(popabA[,3]))
  doma=subset(amo1,amo1[,2]==(1-dab))
  xa=data.frame(doma[,1],doma[,3])
  domabB=subset(amo2,amo2[,2]==dab)
  xabB=data.frame(domabB[,1],domabB[,3])
  txabB=c(Nab,sum(popabB[,1]),sum(popabB[,3]))
  domb=subset(amo2,amo2[,2]==(1-dab))
  xb=data.frame(domb[,1],domb[,3])
  pika=rep(fA,length(doma[,1]))
  pikb=rep(fB,length(domb[,1]))
  pikabA=rep(1/(fA+fB),length(domabA[,1]))
  pikabB=rep(1/(fA+fB),length(domabB[,1]))
  na=length(doma[,1])
  nb=length(domb[,1])
  ba=E.Beta(doma[,4],xa,pika,1,b0=name1)
  bb=E.Beta(domb[,4],xb,pikb,1,b0=name2)
  babA=E.Beta(domabA[,4],xabA,pikabA,1,b0=name1)
  babB=E.Beta(domabB[,4],xabB,pikabB,1,b0=name2)
  somaponda=GREG.SIESPd(N_A,nA,length(doma[,4]),doma[,4],xa,txa,ba,b0=name1)
  somapondab=GREG.SIESPab(N_A,nA,NB,nB,length(domabA[,1]),domabA[,4],xabA,txabA,babA,b0=name2)
  soma1=somaponda[1]+somapondabA[1]
  somapondB=GREG.SIESPd(NB,nB,length(domb[,4]),domb[,4],xb,txb,bb,b0=name2)
  somapondabB=GREG.SIESPab(N_A,nA,NB,nB,length(domabB[,1]),domabB[,4],xabB,txabB,babB,b0=name2)
  soma2=somapondB[1]+somapondabB[1]
  est=soma1+soma2
}

```

```

return(est[1])
}

#####
#Função que calcula o estimador regressão de acordo com a estratégia de Fuller
#####
#-> Precisa da função "GREG.SIESPd" para os domínios
#####
# amo1 -> Amostra de A
# amo2 -> Amostra de B
# pop1 -> População do cadastro A
# pop2 -> População do cadastro B
# N_A -> Tamanho da população do cadastro A
# nA -> Tamanho da amostra obtida do cadastro A
# NB -> Tamanho da população do cadastro B
# nB -> Tamanho da amostra obtida do cadastro B
# Nabf -> Tamanho da população do domínio "ab"
# dab -> Identificador do domínio ab
# name1 -> Identificador para a existência de intercepto no modelo para a pop. de A
# ("TRUE" ou "FALSE")
# name2 -> Identificador para a existência de intercepto no modelo para a pop. de B
# ("TRUE" ou "FALSE")
#####

GREG.SIFUL=function(pop1,pop2,amo1,amo2,N_A, nA, NB,nB,Nabf,dab,name1,name2)
{
amo1 <- as.data.frame(amo1)
amo2 <- as.data.frame(amo2)
popxA=data.frame(pop1[,1],pop1[,3])
txA=c(length(popxA[,1]),sum(pop1[,1]),sum(pop1[,3]))
popa=subset(pop1,pop1[,2]==(1-dab))
popxa=data.frame(popa[,1],popa[,3])
txa=c(length(popa[,1]),sum(popa[,1]),sum(popa[,3]))
popxB=data.frame(pop2[,1],pop2[,3])
txB=c(length(popxB[,1]),sum(pop2[,1]),sum(pop2[,3]))
popb=subset(pop2,pop2[,2]==(1-dab))
popxb=data.frame(popb[,1],popb[,3])
txb=c(length(popb[,1]),sum(popb[,1]),sum(popb[,3]))
xA=data.frame(amo1[,1],amo1[,3])
xB=data.frame(amo2[,1],amo2[,3])
Na=length(popa[,1])
Nb=length(popb[,1])
fA=nA/N_A
fB=nB/NB
pikA=rep(fA,nA)
pikB=rep(fB,nB)
popabA=subset(pop1,pop1[,2]==dab)
popabB=subset(pop2,pop2[,2]==dab)
popxabA=data.frame(popabA[,1],popabA[,3])
popxabB=data.frame(popabB[,1],popabB[,3])
domabA=subset(amo1,amo1[,2]==dab)
xabA=data.frame(domabA[,1],domabA[,3])
txabA=c(length(popabA[,1]),sum(popabA[,1]),sum(popabA[,3]))

```

```

doma=subset(amo1,amo1[,2]==(1-dab))
xa=data.frame(doma[,1],doma[,3])
domabB=subset(amo2,amo2[,2]==dab)
xabB=data.frame(domabB[,1],domabB[,3])
txabB=c(length(popabB[,1]),sum(popabB[,1]),sum(popabB[,3]))
domb=subset(amo2,amo2[,2]==(1-dab))
xb=data.frame(domb[,1],domb[,3])
pika=rep(fA,length(doma[,1]))
pikb=rep(fB,length(domb[,1]))
pikabA=rep(fA,length(domabA[,1]))
pikabB=rep(fB,length(domabB[,1]))
na=length(doma[,1])
nb=length(domb[,1])
nabA=length(domabA[,1])
nabB=length(domabB[,1])
ba=E.Beta(doma[,4],xa,pika,1,b0=name1)
bb=E.Beta(domb[,4],xb,pikb,1,b0=name2)
babA=E.Beta(domabA[,4],xabA,pikabA,1,b0=name1)
babB=E.Beta(domabB[,4],xabB,pikabB,1,b0=name2)
somaponda=GREG.SIESPd(N_A,na,length(doma[,4]),doma[,4],xa,txa,ba,b0=name1)
somapondabA=GREG.SIESPd(N_A,na,length(domabA[,1]),domabA[,4],xabA,txabA,babA,b0=name2)
somapondB=GREG.SIESPd(NB,nb,length(domb[,4]),domb[,4],xb,txb,bb,b0=name2)
somapondabB=GREG.SIESPd(NB,nb,length(domabB[,1]),domabB[,4],xabB,txabB,babB,b0=name2)
w=(nabA*(1-fB)/(nabA*(1-fB)+nabB*(1-fA)))
estful=somaponda[1]+somapondB[1]+w*somapondabA[1]+(1-w)*somapondabB[1]
return(estful)
}

```

```

#####
#Função que calcula o estimador regressão de acordo com a estratégia PML
#####
#-> Precisa da função "GREG.SIESPd" para os domínios
#-> Precisa da função "NabPML" para estimar Nab
#####
# amo1 -> Amostra de A
# amo2 -> Amostra de B
# pop1 -> População do cadastro A
# pop2 -> População do cadastro B
# N_A -> Tamanho da população do cadastro A
# nA -> Tamanho da amostra obtida do cadastro A
# NB -> Tamanho da população do cadastro B
# nB -> Tamanho da amostra obtida do cadastro B
# Nabpml -> Estimador PML para Nab
# dab -> Identificador do domínio ab
# name1 -> Identificador para a existência de intercepto no modelo para a pop. de A
# ("TRUE" ou "FALSE")
# name2 -> Identificador para a existência de intercepto no modelo para a pop. de B
# ("TRUE" ou "FALSE")

```

```
#####
```

```
GREG.SIPML=function(pop1,pop2,amo1,amo2,N_A, nA, NB, nB,Nabpml,dab,name1,name2)
{
  amo1 <- as.data.frame(amo1)
  amo2 <- as.data.frame(amo2)
  popxA=data.frame(pop1[,1],pop1[,3])
  txA=c(length(popxA[,1]),sum(pop1[,1]),sum(pop1[,3]))
  popa=subset(pop1,pop1[,2]==(1-dab))
  popxA=data.frame(popa[,1],popa[,3])
  txA=c((N_A-Nabpml),sum(popa[,1]),sum(popa[,3]))
  popxB=data.frame(pop2[,1],pop2[,3])
  txB=c(length(popxB[,1]),sum(pop2[,1]),sum(pop2[,3]))
  popb=subset(pop2,pop2[,2]==(1-dab))
  popxB=data.frame(popb[,1],popb[,3])
  txB=c((NB-Nabpml),sum(popb[,1]),sum(popb[,3]))
  xA=data.frame(amo1[,1],amo1[,3])
  xB=data.frame(amo2[,1],amo2[,3])
  Na=length(popa[,1])
  Nb=length(popb[,1])
  fA=Na/N_A
  fB=Nb/NB
  pikA=rep(fA,Na)
  pikAu=pikA[1]
  pikB=rep(fB,Nb)
  pikBu=pikB[1]
  popabA=subset(pop1,pop1[,2]==dab)
  popabB=subset(pop2,pop2[,2]==dab)
  popxabA=data.frame(popabA[,1],popabA[,3])
  popxabB=data.frame(popabB[,1],popabB[,3])
  domabA=subset(amo1,amo1[,2]==dab)
  xabA=data.frame(domabA[,1],domabA[,3])
  doma=subset(amo1,amo1[,2]==(1-dab))
  xa=data.frame(doma[,1],doma[,3])
  domabB=subset(amo2,amo2[,2]==dab)
  xabB=data.frame(domabB[,1],domabB[,3])
  domb=subset(amo2,amo2[,2]==(1-dab))
  xb=data.frame(domb[,1],domb[,3])
  pika=rep(fA,length(doma[,1]))
  pikb=rep(fB,length(domb[,1]))
  pikabA=rep(1/pikAu,length(domabA[,1]))
  pikabB=rep(1/pikBu,length(domabB[,1]))
  NabestabA=sum(pikabA)
  NabestabB=sum(pikabB)
  txabA=c(NabestabA,sum(popabA[,1]),sum(popabA[,3]))
  txabB=c(NabestabB,sum(popabB[,1]),sum(popabB[,3]))
  na=length(doma[,1])
  nb=length(domb[,1])
  ba=E.Beta(doma[,4],xa,pika,1,b0=name1)
  bb=E.Beta(domb[,4],xb,pikb,1,b0=name2)
  babA=E.Beta(domabA[,4],xabA,pikabA,1,b0=name1)
  babB=E.Beta(domabB[,4],xabB,pikabB,1,b0=name2)
  somapondA=GREG.SIESPd(N_A,nA,length(doma[,4]),doma[,4],xa,txa,ba,b0=name1)
}
```

```

somapondab=GREG.SIESPd(N_A,nA,length(domabA[,1]),domabA[,4],xabA,txabA,babA,b0=name2)
somapondB=GREG.SIESPd(NB,nB,length(domb[,4]),domb[,4],xb,txb,bb,b0=name2)
somapondabB=GREG.SIESPd(NB,nB,length(domabB[,1]),domabB[,4],xabB,txabB,babB,b0=name2)
tgregabpml=(pikAu*somapondabA[1]+pikBu*somapondabB[1])/(pikAu*NabestabA+pikBu*NabestabB)
pmlgreg=somapondA[1]+tgregabpml+somapondB[1]
return(pmlgreg)
}

```

```

#####INÍCIO DO PROGRAMA#####
library(TeachingSampling)
RNGkind("Marsaglia-Multicarry")
seed=.Random.seed
set.seed(seed)
#####
##ARMAZENANDO VALORES POPULACIONAIS E
##TAMANHOS DE AMOSTRA
#####
N=c(2000,1750,1500,1250,1000,750)
Nab=c(500,1000,1500,2000,2500,3000)
nA=c(125,250,500,750,1000,1250)
nB=c(125,250,500,750,1000,1250)
#####
#####

#####INÍCIO DO LOOP PARA OS TAMANHOS DE "N"
for(i in 1:length(N))
{
Total=matrix(NA,nrow=2,ncol=3)
colnames(Total)=c("NC","Nab","Soma")
Total[2,]=c(N[i],Nab[i],N[i]+Nab[i])
print(Total[2,])
print("#####")

####Carregando o banco de dados
xA11=4+rnorm(N[i])
xA21=8+rnorm(N[i])
xA31=2+rnorm(Nab[i])
xA32=1+rnorm(Nab[i])
xA1=append(xA11,xA31)
xA2=append(xA21,xA32)
yAae=2+xA1[1:N[i]]+xA2[1:N[i]]+rnorm(N[i])
yabA=2+xA31+xA32+rnorm(Nab[i])
yA=append(yAae,yabA)
domA=c(append(matrix(1,N[i]),matrix(0,Nab[i])))
N_A=length(domA)
dadosA=data.frame(XA1=xA1,DA=domA,XA2=xA2,YA=yA)
xB11=4+rnorm(N[i])
xB21=8+rnorm(N[i])
xB31=xA31
xB32=xA32
xB1=append(xB11,xB31)
xB2=append(xB21,xB32)

```



```

dataA=data.frame(amoA[,1],amoA[,3])
pikA=rep(nA[j]/N_A,nA[j])
bA=E.Beta(amoA[,4],dataA,pikA,1,b0=TRUE)
resA=residuo(N_A, nA[j], amoA[,4], dataA, totalA, bA, b0 = TRUE)
dataB=data.frame(amoB[,1],amoB[,3])
totalB=c(sum(amoB[,1]),sum(amoB[,3]))
pikB=rep(nB[j]/NB,nB[j])
bB=E.Beta(amoB[,4],dataB,pikB,1,b0=TRUE)
resB=residuo(NB, nB[j], amoB[,4], dataB, totalB, bB, b0 = TRUE)
covg=covarg(nA[j]+nB[j],N_A+NB,resA,resB)
domabA=subset(amoA,amoA[,2]==0)
domabB=subset(amoB,amoB[,2]==0)
nabA=length(domabA[,1])
nabB=length(domabB[,1])
dataabA=data.frame(domabA[,1],domabA[,3])
dataabB=data.frame(domabB[,1],domabB[,3])
babA=E.Beta(domabA[,4],dataabA,rep(nA[j]/N_A,nabA),1,b0=TRUE)
babB=E.Beta(domabB[,4],dataabB,rep(nB[j]/NB,nabB),1,b0=TRUE)
dadosabA=subset(dadosA,dadosA[,2]==0)
#Nab=length(dadosabA[,1])
dadosabB=subset(dadosB,dadosB[,2]==0)
totalabA=c(sum(dadosabA[,1]),sum(dadosabA[,3]))
totalabB=c(sum(dadosabB[,1]),sum(dadosabB[,3]))
resabA=residuo(Nab[i],nabA,domabA[,4],dataabA,totalabA,babA,b0=TRUE)
resabB=residuo(Nab[i],nabB,domabB[,4],dataabB,totalabB,babB,b0=TRUE)
pgh=preghart(N_A+NB-Nab[i],nA[j],nB[j], N_A,NB,nabA,nabB,Nab[i],resA,resB,resabA,resabB)
greghart[r]=GREG.SIH(pgh,dadosA,dadosB,amoA,amoB,N_A, nA[j], NB,nB[j],Nab[i],0
,"TRUE","TRUE")/total
gregbankier[r]=GREG.SIBA(dadosA,dadosB,amoA,amoB,N_A, nA[j], NB,nB[j],Nab[i],0
,"TRUE","TRUE")/total
gregfuller[r]=GREG.SIFUL(dadosA,dadosB,amoA,amoB,N_A, nA[j], NB,nB[j],Nab[i],0
,"TRUE","TRUE")/total
gregestpml[r]=GREG.SIPML(dadosA,dadosB,amoA,amoB,N_A, nA[j], NB,nB[j],Nab_PML[r],
0,"TRUE","TRUE")/total

}

print(c("ESTIMADORES NORMAIS","valor do parâmetro:",par))
Total=matrix(NA,nrow=4,ncol=4)
colnames(Total)=c("Média","Desvio Padrão", "V.R.", "EQMR")
Total[1,]=c(mean(hartley),sd(hartley),(mean(hartley)-par)/par,var(hartley)
+(mean(hartley)-par)^2)
Total[2,]=c(mean(fuller),sd(fuller),(mean(fuller)-par)/par,var(fuller)
+(mean(fuller)-par)^2)
Total[3,]=c(mean(bankier),sd(bankier),(mean(bankier)-par)/par,var(bankier)
+(mean(bankier)-par)^2)
Total[4,]=c(mean(estpml),sd(estpml),(mean(estpml)-par)/par,var(estpml)
+(mean(estpml)-par)^2)
print(Nab_fuller)

print(Total)
print(c("ESTIMADORES GREG","valor do parâmetro:",par))
Total=matrix(NA,nrow=4,ncol=4)
colnames(Total)=c("Média","Desvio Padrão", "V.R.", "EQMR")

```

Apêndice D

Cenários de seqüências de populações finitas sob a abordagem de cadastro duplo

É possível apresentar exemplos das estruturas das populações finitas que podem ser consideradas sob a abordagem de cadastro duplo. Estes exemplos estão apresentados nas figuras a seguir para $v = 2$.

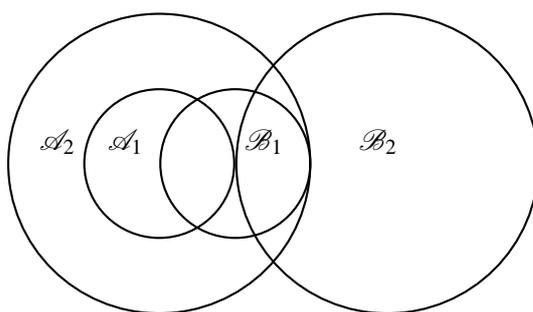


Figura 1: Cenário (b) $A_1 \subset A_2$ e $B_1 \subset A_2$, e $A_1 \cup B_1 \subset A_2$

O cenário (b) ilustra a união entre os primeiros elementos de cada seqüência encaixados no segundo elemento da seqüência referente ao cadastro A.

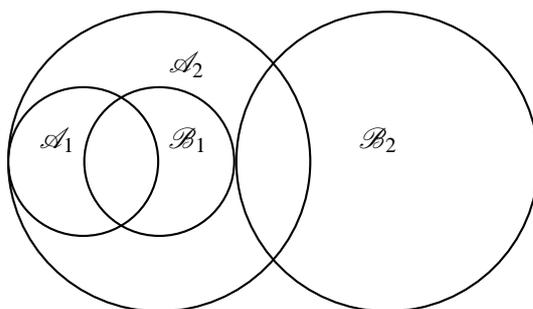


Figura 2: Cenário (c) $A_1 \subset A_2$ e $A_1 \cup B_1 \subset A_2$, com $(A_1 \cup B_1) \cap B_2 = \emptyset$

O cenário (c) ilustra situação equivalente ao cenário (b), mas com o detalhe de que a união dos cadastros não contemplar elementos da interseção entre os elementos das sequências.

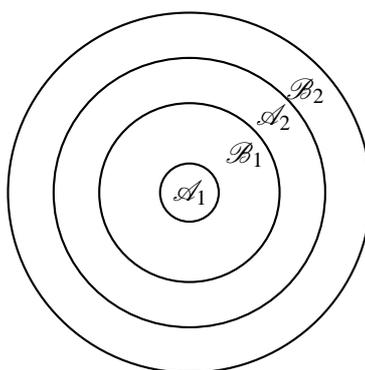


Figura 3: Cenário (d) $A_1 \subset A_2 \subset B_1 \subset B_2$

No cenário (d) é possível observar a situação em que o cadastro A está contido no cadastro B , e vice-versa.

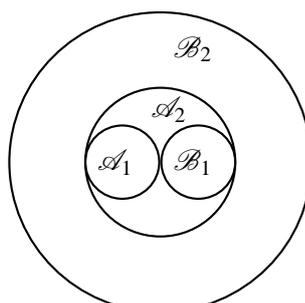


Figura 4: Cenário (e) $(A_1 \cup A_2) \subset A_2$ e $A_2 \subset B_2$

O cenário (e) apresenta como característica a não existência de interseção entre os primeiros elementos das sequências referentes a A e B . Porém, a união destes elementos é embutida nos demais elementos das sequências.

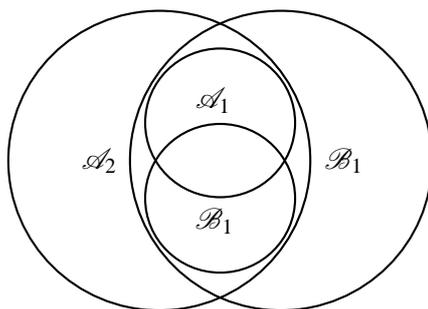


Figura 5: Cenário (f) $(A_1 \cup B_1) \subset A_2 \cap B_2$

O cenários (f) e (g) apresentam situação similar ao apresentado anteriormente, apenas com o fato de que a interseção entre os primeiros elementos das sequências existe.

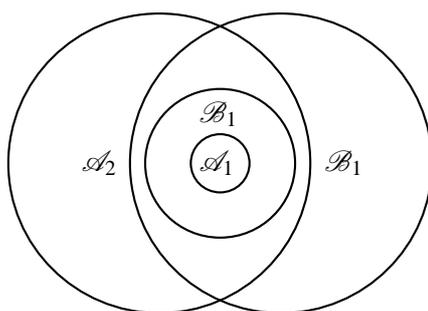


Figura 6: Cenário (g) $A_1 \subset B_1$ e $(A_1 \cup B_1) \subset A_2 \cap B_2$

REFERÊNCIAS

ALPIZAR-JARA, R.; POLLOCK, K. H.; HAINES, D. E. **Mark-recapture estimator for dual frame population size of prominent nesting structures: the effect of uncertain detection probability.** *Environmental and Ecological Statistics*, 2005.

ARMSTRONG, B. **Test of multiple frame sampling techniques for agricultural surveys: New brunswick, 1978.** *Proceedings of the Section on Survey Research Methods of the American Statistical Association*, pp. 295-300, 1979.

BANKIER; LEPKOWSKI, J. M.; GROVES, R. M. **Estimators based on several stratified samples with applications to multiple frame surveys.** *Journal of the American Statistical Association*, 1986.

BOSECKER, R. R.; FORD, B. L. **Multiple frame estimation with stratified overlap domain.** *Proceedings of the Social Statistics Section, American Statistical Association*, 1976.

BRACKSTONE, G. J.; RAO, J. N. K. **An investigation of raking ratio estimation.** *Sankhya: The Indian Journal of Statistics*, 16, Supp., 47-55, 1979.

BREIDT, F. J.; OPSOMER, J. D. **Endogenous post-stratification in surveys: Classifying with a sample-fitted model.** *Annals of Statistics*, 2008.

CARFAGNA, E. **List frames, area frames and administrative data, are complementary or in competition?** *invited paper to the Meeting MEXSAI Conference organized by Eurostat, FAO, OCSE, UN/ECE, NASS/USDA, JRC, ISI, Istat, SAGARPA*, 2004.

- CASADY, R. J.; LEPKOWSKI, J. M. **Optimal allocation for stratified telephone survey designs**. *Proceedings of the Section on Survey Research Methods of the American Statistical Association*, pp. 111-116, 1991.
- CASADY, R. J.; LEPKOWSKI, J. M. **Stratified telephone survey designs**. *Survey Methodology*, Vol. 19, pp. 103-113, 1993.
- CHOUDHRY, G. H.; PARK, I.; LI, W. T. **Dual frame sample design for a national sample of veterans**. *Joint Statistical Meetings, Section of Survey Research Methods*, 2002.
- COELHO, H. F. C. **A Abordagem de Cadastro Duplo: Estimação Assistida por Modelos com Aplicações em Pesquisas Agropecuárias**. Dissertação (Mestrado) — Universidade Federal de Pernambuco, Departamento de Estatística, 2007.
- COELHO, H. F. C.; FERRAZ, C. **Ratio type estimators for stratified dual frame surveys**. *Proceedings of the 56th Session of the ISI*, 2007.
- FULLER, W. A. **Introduction to statistical time series**. Wiley-Interscience, 1996.
- FULLER, W. A. **Sampling statistics**. *Wiley series in Survey Methodology*, 2009.
- FULLER, W. A.; BURMEISTER, L. F. **Estimators for samples selected from two overlapping frames**. *Proceedings of the Social Statistics Section, American Statistical Association*, 245-249, 1972.
- HAINES, D. E.; POLLOCK, K. H. **Estimating the number of active and successful bald eagle nests: an application of the dual frame method**. *Environmental and Ecological Statistics*, 1998.
- HARTLEY, H. O. **Multiple frame surveys**. *Proceedings of the Social Statistics Association-ASA*, 1962.
- ISAKI, C. T.; FULLER, W. A. **Survey design under the regression superpopulation model**. *Journal of American Statistical Association*, 1982.

LOHR, S. L. **Inference from dual frame surveys**. *Journal of American Statistical Association*, 2007.

LÓPEZ, O. F. M. **Comparación de la eficiencia de estimadores para marcos duales**. Dissertação (Mestrado) — Universidad Nacional de Colombia, Facultad de Ciencias-Departamento de Estadística, 2010.

LU, Y. **Longitudinal Estimation in Dual Frame Surveys**. Tese (Doutorado) — Arizona State University, 2007.

LUND, R. E. **Estimators in multiple frame surveys**. *Proceedings of the Social Statistics Section, American Statistical Association*, 1968.

PARK, M. **Regression Estimation of the Mean in Survey Sampling**. Tese (Doutorado) — Iowa State University, 2002.

POTTER, F. J. *et al.* **List-assisted rdd surveys**. *Proceedings of the Section on Survey Research Methods of the American Statistical Association*, pp. 117-122, 1991.

POVEDA, L. M. R. **Estimação em populações assistidas por modelos para variáveis dicotômicas**. Dissertação (Mestrado) — Universidade Federal de Pernambuco, Departamento de Estatística, 2007.

ROBINSON, P. M.; SÄRNDAL, C. E. **Asymptotic properties of the generalized regression estimator in probability sampling**. *Sankhya: The Indian Journal of Statistics*, 45, Series B, 240-248, 1983.

SEN, P. K.; SINGER, J. M.; LIMA, A. C. Pedroso de. **From finite sample to asymptotic methods in statistics**. Cambridge Series in Statistical and Probabilistic Mathematics (No. 29), 2010.

SHIN, H. C.; MOLINARI, N. A.; WOLTER, K. **A dual frame design for the national immunization survey**. *Section on Survey Research Methods - AAPOR*, 2008.

SINGH, A. C.; WU, S. **An extension of generalized estimator to dual frame surveys**. *Join*

Statistical Meetings - Section on Survey Research Methods, 2003.

SKINNER, C. J. **On the efficiency of raking ratio estimation for multiple frame surveys.** *Journal of the American Statistical Association*, 1991.

SKINNER, C. J.; RAO, J. N. K. **Estimation in dual frame surveys with complex designs.** *Journal of the American Statistical Association*, 1996.

SÄRNDAL, C. E. **Model assisted survey sampling.** *New York: Springer*, 1992.

SÄRNDAL, C. E.; SWENSSON, B.; WRETMAN, J. **Model assisted survey sampling.** *Springer Series in Statistics*, 2003.

WU, C.; RAO, J. N. K. **Pseudo-empirical likelihood ratio confidence intervals for complex surveys.** *The Canadian Journal of Statistics*, 34, 359-375, 2006.