



Nádia Maria Coelho Rodrigues

**Concatenação Estatística de Dados e Afiliação
Estocástica**

Tese de Doutorado

Tese apresentada ao Programa de Pós-Graduação em Engenharia Elétrica da PUC_Rio como requisito parcial para obtenção do título de Doutor em Engenharia Elétrica.

Orientador: Reinaldo C. Souza
Co-orientador: Kaizô I. Beltrão

Rio de Janeiro, 13 de outubro de 2009.



Nádia Maria Coelho Rodrigues

**Concatenação Estatística de Dados e Afiliação
Estocástica**

Tese apresentada como requisito parcial para obtenção do grau de Doutor pelo Programa de Pós-Graduação em Engenharia Elétrica do Departamento de Engenharia Elétrica do Centro Técnico Científico da PUC-Rio. Aprovada pela Comissão Examinadora abaixo assinada.

Prof. Reinaldo C. Souza
Orientador

Departamento de Engenharia Elétrica – PUC-Rio

Prof. Kaizô I. Beltrão
Co-Orientador
ENCE/IBGE

Profa. Monica Barros
Departamento de Engenharia Elétrica – PUC-Rio

Prof. Victor Hugo de Carvalho Gouvêa
UFF

Prof. Sérgio da Costa Cortes
Departamento de Informática – PUC-Rio

Profa. Ana Carolina Letichevsky
Fundação Cesgranrio

Prof. José Eugênio Leal
Coordenador Setorial do Centro
Técnico Científico

Rio de Janeiro, 13 de outubro de 2009.

Todos os direitos reservados. É proibida a reprodução total ou parcial do trabalho sem autorização da universidade, da autora e do orientador.

Rodrigues, Nádia Maria Coelho

Na Escola Nacional de Ciências Estatísticas graduou-se em Estatística em 1976 e no curso de mestrado em Estudos Populacionais e Pesquisas Sociais em 2003.

Ficha Catalográfica

Rodrigues, Nádia Maria Coelho

Concatenação estatística de dados e afiliação estocástica / Nádia Maria Coelho Rodrigues ; orientador: Reinaldo C. Souza ; co-orientador: Kaizô I. Beltrão. – 2009.

159 f. ; 30 cm

Tese (Doutorado em Engenharia Elétrica)–Pontifícia Universidade Católica do Rio de Janeiro, Rio de Janeiro, 2009.

Inclui bibliografia

1. Engenharia elétrica – Teses. 2. Emparelhamento estatístico. 3. Suposição de independência condicional (CIA). 4. Informação auxiliar. 5. Procedimento hot deck. I. Souza, Reinaldo C. II. Beltrão, Kaizô I. III. Pontifícia Universidade Católica do Rio de Janeiro. Departamento de Engenharia Elétrica. IV. Título.

CDD: 621.3

A meus pais e aos meus filhos por seu amor, carinho e seu contínuo apoio e estímulo.

Agradecimentos

Aos meus orientadores Professores Reinaldo C. Souza e Kaizô I. Beltrão, pelo apoio, confiança e pela sua dedicada orientação.

Ao professores examinadores Ana Carolina Letichevsky, Monica Barros, Sérgio da Costa Cortes e Victor Hugo de Carvalho Gouvêa.

Ao IBGE pelo apoio financeiro. Aos funcionários e colegas do IBGE, principalmente da Diretoria de Informática.

À Pontificia Universidade Católica de Rio de Janeiro, a todos os professores e funcionários do Departamento.

A meus amigos Mário Gama e Emilia Matos por contínuo seu apoio e colaboração.

A todos os amigos e familiares que de uma forma ou de outra me estimularam ou ajudaram.

Resumo

Rodrigues, Nádia Maria Coelho; Souza, Reinaldo C.. **Concatenação Estatística de Dados e Afiliação Estocástica**. Rio de Janeiro, 2009. 159p. Tese de Doutorado - Departamento de Engenharia Elétrica, Pontifícia Universidade Católica do Rio de Janeiro.

O emparelhamento estatístico é a técnica de combinar informações de duas ou mais fontes de dados, possivelmente de pesquisas independentes, que possuam um subconjunto comum de variáveis, para produzir uma informação mais abrangente e coerente, em um arquivo de dados síntese, onde as variáveis observadas nas diferentes amostras são gravadas conjuntamente. Métodos computacionalmente intensivos viabilizam novas formas de emparelhamento estatístico. Em um estudo de simulação, um único arquivo de dados da PNAD é dividido em duas bases de dados para emular um caso de amostragem matricial. Esses dois arquivos são emparelhados estatisticamente utilizando quatro metodologias, e os resultados das mesmas são comparados com as do arquivo único original. A CIA (*conditional independence assumption*) não parece ser válida. Para evitar a suposição de independência condicional (CIA), os três métodos de emparelhamento estatístico desenvolvidos são baseados na hipótese do relacionamento teórico de afiliação estocástica entre as variáveis contínuas renda e aluguel, e dois deles usam também a informação de resíduos. Os métodos são comparados entre si e com o método clássico, que se baseia na suposição de independência condicional – CIA. Em uma abordagem não-paramétrica, com um objetivo micro, os métodos de emparelhamento estatístico propostos são irrestritos, e realizam-se em classes, definidas pela variável de número de cômodos. Usam um procedimento *distance hot deck*, além de adicionar os resíduos supostos conhecidos. Esse estudo investiga os resultados de viés e do EQM, dos quatro métodos, investigando a preservação da correlação original entre renda e aluguel.

Palavras-chave

Emparelhamento estatístico, suposição de independência condicional (CIA), informação auxiliar, procedimento *hot deck*

Abstract

Rodrigues, Nádía Maria Coelho; Souza, Reinaldo C.. **Statistical Matching using Stochastic Affiliation**. Rio de Janeiro, 2009. 159p. DSc Thesis - Departamento de Engenharia Elétrica, Pontifícia Universidade Católica do Rio de Janeiro.

Statistical matching is the art of combining information from two or more datasets, collected in independent samples, but a subset of the variables is common to both samples, to produce both coherent and comprehensive information in a synthetic data file, where variables observed in different samples are jointly recorded. Modern computing can make possible, under techniques described here, an advance in the application of Statistical Matching. It is reported on a simulation that splits a single file into two pieces, statistically matched the pieces using four methodologies and compares the results to the original single file. The Conditional Independence Assumption (CIA) did not seem a valid assumption. To avoid CIA, it is suggested two methods of statistical matching where kind of auxiliary information, based on Stochastic Affiliation relationship between income and rent, and residuals are used, in conjunction with nonparametric regression and hot deck distance. Both are compared not only with the classic method based on the CIA and regression without residuals but also with each other. In a nonparametric approach, with a micro objective, an unconstrained statistical matching is applied using hot-deck nearest neighbor within classes, using the household logarithm income and logarithm rent percentile groups. We are satisfied with the performance of creating a correlation coefficient of Y and Z , as measured using MSE.

Keywords

Statistical matching, conditional independence assumption (*CIA*), auxiliary information, hot deck imputation procedure.

Sumário

1 INTRODUÇÃO	14
1.1. CONSIDERAÇÕES GERAIS	14
1.2. ESTRUTURA DO EMPARELHAMENTO ESTATÍSTICO	16
1.3. OBJETIVO, ABORDAGEM E REPRESENTATIVIDADE DO EMPARELHAMENTO ESTATÍSTICO	20
1.3.1. <i>Objetivo e abordagem do emparelhamento estatístico</i>	20
1.3.2. <i>Representatividade do arquivo síntese predito</i>	24
1.3.3. <i>O caso particular da CIA</i>	25
1.4. VARIANTE IMPORTANTE DO EMPARELHAMENTO ESTATÍSTICO	29
1.4.1. <i>Aplicações de Amostragem Matricial</i>	30
1.5. O “PROBLEMA DE IDENTIFICAÇÃO” DO MÉTODO	39
1.6. SOLUÇÕES PROPOSTAS	42
1.7. CONTRIBUIÇÃO DESSA TESE	46
1.8. ESTRUTURA DA TESE	48
2 REVISÃO DA LITERATURA SOBRE EMPARELHAMENTO ESTATÍSTICO 50	
2.1. INTRODUÇÃO	50
2.2. EMPARELHAMENTO ESTATÍSTICO RESTRITO E IRRESTRITO	56
2.2.1. <i>Emparelhamento irrestrito</i>	59
2.2.2. <i>Emparelhamento restrito</i>	60
2.2.3. <i>Comparação do emparelhamento irrestrito e restrito</i>	64
2.3. HARMONIZAÇÃO DAS PESQUISAS ANTES DO EMPARELHAMENTO	65
2.4. AFILIAÇÃO E DEPENDÊNCIA POSITIVA	69
2.5. ADIÇÃO DE RESÍDUOS ALEATÓRIOS	73
2.5.1. <i>Método de Kadane</i>	73
2.6. TRANSFORMAÇÃO PERCENTIL MONOTÔNICA	78
3 PRESSUPOSTO DE INDEPENDÊNCIA CONDICIONAL - CIA	79
3.1. PROCEDIMENTO MICRO NÃO-PARAMÉTRICO SOB A CIA	81
3.2. TÉCNICAS <i>HOT DECK</i>	83
3.2.1. <i>Random hot deck</i>	83
3.2.2. <i>Rank hot deck</i>	84
3.2.3. <i>Distance hot deck</i>	85
3.2.4. <i>Matching noise</i>	88
4 INFORMAÇÃO AUXILIAR	89
4.1. DIFERENTES TIPOS DE INFORMAÇÃO AUXILIAR	89
4.2. O CASO DA NORMAL MULTIVARIADA	91

COVARIÂNCIAS MARGINAIS CONHECIDAS	92
4.3. PROCEDIMENTO MICRO NÃO-PARAMÉTRICO USANDO INFORMAÇÃO AUXILIAR	94
4.4. MÉTODOS MISTOS	96
5 MÉTODO PROPOSTO E SIMULAÇÃO REALIZADA.....	97
5.1. INTRODUÇÃO AO MÉTODO PROPOSTO	99
5.2. METODOLOGIA DA SIMULAÇÃO.....	102
5.2.1. <i>Base reais de dados</i>	102
5.2.2. <i>Propostas de procedimentos não-paramétricos de emparelhamento estatístico</i>	107
6 ESTUDO DA SIMULAÇÃO.....	114
6.1. INTRODUÇÃO	114
6.1.1. <i>Qualidade da informação</i>	118
6.1.2. <i>Resultados</i>	124
7 CONSIDERAÇÕES FINAIS	135
8 REFERÊNCIAS BIBLIOGRÁFICAS	138
9 APÊNDICES.....	149
9.1. APÊNDICE A.....	149
9.2. APÊNDICE B	151
9.3. APÊNDICE C.....	151
9.4. APÊNDICE D.....	159

Lista de Figuras

<i>Figura 1 - Problema geral: duas amostras independentes com n_A e n_B unidades, com suas variáveis específicas Y e Z, com variáveis comuns X</i>	<i>17</i>
<i>Figura 2.....</i>	<i>33</i>
<i>Figura 3 - Box-plot Lrenda, Laluguel, p_Y e p_Z por número de cômodos.....</i>	<i>106</i>
<i>Figura 4 - Distribuição uniforme bidimensional dos resíduos, limitados pelas retas que garantem os percentis preditos entre 0 e 1.....</i>	<i>109</i>
<i>Figura 5 - Histograma e Box-plot da distribuição (Y,Z) e seus percentis.....</i>	<i>120</i>
<i>Figura 6 - Distribuição uniforme bidimensional dos percentis p_Y e p_Z, sem “indicação de ajuste” linear.</i>	<i>120</i>
<i>Figura 7 - Distribuição das variáveis da Lrenda e Laluguel por número de cômodos.....</i>	<i>121</i>
<i>Figura 8 - Distribuição das variáveis da renda e aluguel e seus percentis.....</i>	<i>122</i>
<i>Figura 9 - Resíduos por número de cômodos.....</i>	<i>122</i>
<i>Figura 10 – Distribuição de renda e aluguel para domicílios de 5 cômodos.</i>	<i>124</i>
<i>Figura 11 – Distribuição do aluguel antes e depois do emparelhamento</i>	<i>133</i>
<i>Figura 12 - Distribuição do aluguel antes e depois do emparelhamento</i>	<i>134</i>

Lista de Tabelas

<i>Tabela 1: Registros do arquivo A</i>	57
<i>Tabela 2: Estatísticas descritivas de A</i>	57
<i>Tabela 3: Registros do arquivo B</i>	58
<i>Tabela 4: Estatísticas descritivas de B</i>	58
<i>Tabela 5: Resultado do emparelhamento irrestrito.</i>	59
<i>Tabela 6: Estatísticas descritivas do emparelhamento irrestrito B doador</i>	60
<i>Tabela 7: Resultado do emparelhamento restrito.</i>	62
Tabela 8 - Número de registros nas amostras, dos arquivos, por número de cômodos	123
<i>Tabela 9 - Freqüência de Lrenda e Laluguel para domicílios de 5 cômodos.</i>	123
<i>Tabela 10 - Valores originais das covariâncias e correlações.</i>	125
<i>Tabela 11 - Estimativas das covariâncias do emparelhamento estatístico segundo os diferentes métodos.</i>	126
<i>Tabela 12 - Estimativas das correlações do emparelhamento estatístico segundo os diferentes métodos.</i>	127
<i>Tabela 13 - Estimativas das correlações do emparelhamento estatístico B doador.</i>	159
<i>Tabela 14 - Estimativas das correlações do emparelhamento estatístico A doador.</i>	159

Lista de Quadros

Quadro 1 - Abordagens e objetivos do emparelhamento estatístico.22

Quadro 2 - Valores do intervalo que contém \hat{p}_{Y_i} 110

Lista de Siglas

EQM	Erro quadrado médio
IBGE	Instituto Brasileiro de Geografia e Estatística
i.i.d	Independente e identicamente distribuídos
MEC	Ministério de Educação e Cultura
MRJ	Método de replicação <i>Jackknife</i>
NAEP	National Center for Education Statistics
NHIS	National Health Interview Survey
PNAD	Pesquisa Nacional por Amostra de Domicílios
POF	Pesquisa de Orçamentos Familiares
SAEB	Sistema Nacional de Avaliação de Educação Básica
INEP	Instituto Nacional de Estudos e Pesquisas Educacionais Anísio Teixeira
PPV	Pesquisa sobre Padrões de Vida
CIA	<i>conditional independence assumption</i>
SAM	<i>Social Accounting Matrix</i>
TNR_{lr}	Transformação não-paramétrica relacional intervalar
TNR_{io}	Transformação não-paramétrica relacional com interpolação
EMV	Estimador de máxima verossimilhança
SQS	<i>Split Questionnaire Surveys</i>
PLRD	<i>positive likelihood ratio dependence</i>