**UNIVERSIDADE FEDERAL DO CEARÁ**
**FACULDADE DE ECONOMIA, ADMINISTRAÇÃO, ATÚARIA E CONTABILIDADE**
**PROGRAMA DE PÓS-GRADUAÇÃO EM ECONOMIA - CAEN**


**DIEGO DE MARIA ANDRÉ**


**THREE ESSAYS ON APPLIED MICROECONOMETRICS WITH SPATIAL EFFECTS**


FORTALEZA

2016

DIEGO DE MARIA ANDRÉ

THREE ESSAYS ON APPLIED MICROECONOMETRICS WITH SPATIAL EFFECTS

Tese de Doutorado submetida à Coordenação do Programa de Pós-Graduação em Economia – CAEN, da Faculdade de Economia, Administração, Atúaria e Contabilidade da Universidade Federal do Ceará, como requisito parcial para a obtenção do título de Doutor em Ciências Econômicas. Área de concentração: Econometria aplicada.

Prof. Dr. José Raimundo de Araújo Carvalho Júnior

FORTALEZA

2016

DIEGO DE MARIA ANDRÉ

THREE ESSAYS ON APPLIED MICROECONOMETRICS WITH SPATIAL EFFECTS

> Tese de Doutorado submetida à Coordenação do Programa de Pós-Graduação em Economia – CAEN, da Faculdade de Economia, Administração, Atúaria e Contabilidade da Universidade Federal do Ceará, como requisito parcial para a obtenção do título de Doutor em Ciências Econômicas. Área de concentração: Econometria aplicada.

Aprovada em: 18 de Maio de 2016.

BANCA EXAMINADORA

_____

Prof. Dr. José Raimundo de Araújo Carvalho
Júnior (Orientador)
Universidade Federal do Ceará (UFC)

_____

Prof. Dr. João Mário Santos de França
Universidade Federal do Ceará (UFC)

_____

Prof. Dr. Emerson Luís Lemos Marinho
Universidade Federal do Ceará (UFC)

_____

Prof. Dr. Victor Hugo de Oliveira Silva
Instituto de Pesquisa e Estratégia Econômica do
Ceará (IPECE)

_____

Prof. Dr. Cleyber Nascimento de Medeiros
Instituto de Pesquisa e Estratégia Econômica do
Ceará (IPECE)

Aos meus pais e à minha esposa.

# AGRADECIMENTOS

Nesse momento tão especial da minha vida, não poderia deixar de agradecer as pessoas que de alguma forma contribuiram para que eu alcançasse o meu objetivo. Em primeiro lugar, agradeço a Deus, por nos dar a oportunidade de vivenciar esse grande mistério que é a vida.

Agradeço aos meus pais, Haroldo e Célia, que através de muito esforço e dedicação, conseguiram me dar educação e me ensinaram a valorizar as pequenas coisas da vida, mostrando-me que só tem valor aquilo que é conseguido com o suor do nosso trabalho. Se hoje sou Doutor em economia, devo isso eles. Obrigado!

A minha esposa, Talita, que durante os nossos quase 9 anos juntos, sempre tem me dado apoio, amor e carinho nos momentos díficeis, sempre me incentivando a continuar a busca por soluções quando estou quase desistindo. Sem esses incentivos, com certeza, não teria conseguido. Te amo e Obrigado!

Ao professor José Raimundo, com quem tenho aprendido bastante nesses 6 anos em que temos trabalhado juntos (Mestrado e Doutorado), e cujos ensinamentos levarei para a vida toda.

Aos professores João Mário, Emerson Marinho, Victor Hugo e Cleyber Nascimento por disponibilizarem seu tempo para participar da banca de avaliação e por suas contribuições ao trabalho.

Aos demais professores do CAEN por suas contribuições à minha formação acadêmica.

A todos os funcionários do CAEN, em especial ao Cléber e o S. Adelino, que com suas conversas e brincadeiras ajudam a tornar o CAEN um lugar especial.

Agradeço ainda à todos os colegas do Laboratório de Econometria e Otimização (LECO), não só os atuais (Sylvia, Abel, Luan, Marcelino, Sara) mas a todos os que em algum momento passaram pelo LECO durante os 6 anos que trabalho lá (Luis Carlos, Yuri, Isadora), pelo convívio e trocas de conhecimento e experiências. Agradeço também a todos os colegas de Mestrado e Doutorado do CAEN.

Por fim, agradeço ao CAEN por disponibilizar sua estrutura para a realização deste trabalho e a CAPES pela bolsa concedida.

A todos, os meus sinceros agradecimentos.

# RESUMO

A presente tese é composta por três capítulos, independentes entre si, em microeconometria aplicada. O primeiro capítulo aplica o instrumental teórico e empírico da econometria espacial para analisar os determinantes da demanda residencial de água para a cidade de Fortaleza (Brasil). Estimamos três modelos econométricos, que tem como variáveis explicativas o preço médio/marginal, a diferença, renda, número de homens e mulheres residentes, número de banheiros, sob diferentes especificações espaciais: O modelo de erro espacial (SEM), o modelo espacial autorregressivo (SAR) e o modelo espacial autorregressivo de médias móveis (SARMA), sendo o modelo SARMA o que melhor se ajusta aos dados. Os resultados indicaram que não controlar pelos efeitos espaciais é uma fonte de erro de especificação, subestimando o efeito de quase todas as variáveis. Algumas vezes, essas diferenças podem chegar a 24.66% e 13.32% para a elasticidade-preço no modelo de preço médio e no modelo de McFadden, respectivamente. No segundo capítulo estima-se a disposição a pagar (WTP) pela redução estocástica de primeira ordem no risco de ser roubado, para a cidade de Fortaleza (Brasil). Inspirado por Cameron e DeShazo (2013), desenvolveu-se um modelo simples de escolha que aninha o processo de avaliação contingente (CV) entre loterias e estimou-se por máxima verossimilhança paramétrica e pelo modelo de regressão geograficamente ponderada (GWR). Para o modelo global, isto é, sem efeitos espaciais, estimou-se uma disposição a pagar média de R$ 23.35 por mês/por residência, e um valor implícito de um roubo estatístico de R$ 11,969 por crime evitado. Para o modelo local (GWR), implementou-se o protocolo da krigagem para calcular uma superfície de disposição a pagar. Os resultados sugerem que embora na periferia a disposição a pagar seja menor, à medida que vamos para o centro da cidade existe muita heterogeneidade na distribuição espacial da disposição a pagar para a redução do risco de roubo. No terceiro capítulo analisou-se como o rendimento acadêmico de alunos universitários é afetado pelos seus colegas de sala, através de um desenho descontínuo. Utilizando dados da Universidade Federal do Ceará (UFC), empregamos o modelo de regressão descontínua (RDD) para estimar a diferença entre entrar na turma do primeiro ou do segundo semestre. Devido à quantidade de cursos disponível na nossa base de dados, classificamos os cursos em quatro categorias, de acordo com as notas de entrada no vestibular. Então, procedemos com a estimação de um modelo multi-tratamento. Os resultados mostram que os alunos que foram classificados um pouco acima do limite de vagas (turma do primeiro semestre) têm rendimento acadêmico 2% menor (-0.19) do que alunos que tiveram classificação um pouco abaixo desse limite (turma do segundo semestre). Ademais, encontramos não linearidades nesses efeitos, assim como Sacerdote (2001) e Zimmerman (2003), com intervalos entre 2.5 e -0.18.

**Palavras-chave:** Demanda por água. Efeitos Espaciais. Crime. Avaliação Contingente. *Peer Effects*. Regressão Descontínua.

# ABSTRACT

This Thesis consists of three independent essays on applied microeconometrics. The first chapter applies theoretical and empirical tools of spatial econometrics to analyze the determinants of residential water demand function for the city of Fortaleza (Brazil). We estimated three econometric models, which have as explanatory variables the average/marginal price, the difference, income, number of male and female residents and the number of bathrooms, under different spatial specifications: the Spatial Error Model (SEM), the Spatial Autoregressive model (SAR), and finally, the Spatial Autoregressive Moving Average model (SARMA), which is the model that best fitted the data. Results suggest that not controlling for spatial effects is a key specification error, underestimating the effect of almost all variables in the model. Sometimes, these differences can be as high as 24.66 % and 13.32 % for price elasticity in the Average Price and the McFadden models, respectively. In the second chapter we estimated willingness to pay (WTP) for a first order stochastic reduction on the risk of robbery, for the city of Fortaleza (Brazil). Inspired by Cameron and DeShazo (2013), we develop a simple choice model that nests a process of contingent valuation (CV) among lotteries and estimate it by both parametric maximum likelihood and geographically weighted regression (GWR). For the global model (i.e., without spatial effects), we estimated an average WTP of R\$ 23.35 per month/household, and an implicit value of a statistical robbery approximately equal to R\$ 11,969 per crime avoided. For the local model (GWR), we implement a protocol to calculate a surface of WTP using Kriging techniques. The results suggests that although peripheries present lower willingness to pay, as long as we go inwards there is plenty of heterogeneity on its spatial distribution for risk reductions. In the third chapter we analyzed how undergraduate students' academic performance is affected by theirs classmates, by means of a "discontinuity design". With data from Ceará Federal University (UFC), we employed regression discontinuity design (RDD) to estimate the difference between entering in the first semester class or second semester class. Due to the great courses availability, we assign each course into one of four categories depending on its admitted students' results at the entrance exam. Then, we proceed the estimation exercise using a multi-treatment effect model. Results show that students who were ranked just above the cutoff (first semester class) had an academic performance 2% smaller (-0.19) than students who were ranked just below the cutoff (second semester class). Moreover, we found non-linearities in this effect, as well as Sacerdote (2001) and Zimmerman (2003), with intervals between 0.5 to -0.18.

**Keywords:** Water demand. Spatial Effects. Crime. Contingent Valuation. Peer Effects. Regression Discontinuity Design.

# LIST OF FIGURES

# LIST OF TABLES

# CONTENTS

# 1 GENERAL INTRODUCTION

This Thesis is a collection of three independent essays on microeconometrics using data from the city of Fortaleza (Brazil). In the first chapter, I examine the determinants of urban residential water demand. The second chapter studies willingness to pay for a reduction on the risk of robbery. The third chapter studies how "peer effects" determines academic performance in high education.

Fortaleza is the state capital of Ceará. Located in northeastern Brazil, a dry climate region, the city has a population of about to 2.5 million and it is the fifth largest city in Brazil with an area of 313 square kilometers, boasting one of the highest demographic densities in the country (8,001 per $km^2$). Fortaleza's economy is mainly based on trade, service, and tourism and its gross domestic product is the largest in northeastern Brazil. For being a large urban center, Fortaleza has many problems and criminality is one of these. The city is ranked as number one in intentional lethal violent crimes against life in Brazil and one of the most violent city in the world, which creates a great sense of insecurity among the population. In educational sector, Fortaleza hosted one of the best university in Brazil, besides many others privates universities, which attracts students from several cities of the region. In this sense, this thesis analyze three important aspects of the city of Fortaleza: The Water scarcity, the criminality and the higher education.

Figure 1.1 – Geographical location of the city of Fortaleza

The first chapter is "Spatial determinants of urban residential water demand in Fortaleza, Brazil". This essay is a enhanced version of my dissertation, and was published with Professor José Raimundo Carvalho in Water Resources Management[1]. In this essay we estimated a residential water demand function for the city of Fortaleza, Brazil, considering the potential impact of including spatial effects in the model. The empirical evidence is a unique micro-data set obtained through a household water consumption survey carried out in 2007. We estimated three econometric models, which have as explanatory variables the average/marginal price, the difference, income, number of male and female residents and the number of bathrooms, under different spatial specifications: the Spatial Error Model (SEM), the Spatial Autoregressive model (SAR), and finally, the Spatial Autoregressive Moving Average model (SARMA). Results suggest that the SARMA model is the "best" as shown by a series of tests. Such results contradict conclusions drawn by Chang et al. (2010), House-Peters et al. (2010), and Ramachandran and Johnston (2011). This means, among other things, that not controlling for spatial effects is a key specification error, underestimating the effect of almost all variables in the model. Sometimes, these differences can be as high as 24.66 % and 13.32 % for price elasticity in the Average Price and the McFadden models, respectively.

The second chapter is "Spatial willingness to pay for a first order stochastic reduction on the risk of robbery"[2]. In this essay we estimated willingness to pay (WTP) for a first order stochastic reduction on the risk of crimes for residents of a large and dense urban center. Inspired by Cameron and DeShazo (2013), we develop a simple structural choice model that nests a process of contingent valuation (CV) among lotteries and estimate it by both parametric maximum likelihood and geographically weighted regression (GWR). Our empirical support is a unique and rich micro data set about victimization in Fortaleza, CE (Brazil). For the global model (i.e., without spatial effects), we estimated an average WTP of R$ 23.35 per month/household, and an implicit value of a statistical robbery approximately equal to R$ 11,969 per crime avoided. By means of geographically weighted regression (GWR), we find that variables Sex, Age and Education present a reasonable amount of spatial heterogeneity and, as expected, follow the very inertial city's socioeconomic spatial distribution profile. We implement as well a protocol to calculate a surface of WTP using Kriging techniques. Income, age, and crime spatial distributions have important effects on the surface of WTP. Although peripheries present lower willingness to pay, as long as we go inwards there is plenty of heterogeneity on its spatial distribution for risk reductions. Our results supports a theory of crime with an active role for victim (costly) precautions.

[1] Spatial Determinants of Urban Residential Water Demand in Fortaleza, Brazil. Water Resources Management , v. 28, p. 2401-2414, 2014.
[2] This essays was presented at the $42^{nd}$ economics national meeting - ANPEC 2014, at the VII CAEN-EPGE Public Policy and Economic Growth meeting (2015) and accepted for presentation at the Spatial Econometrics Association Annual meeting - IX world conference SEA 2015 (Miami, USA).

The third chapter is "peer effects and academic performance in higher education - a regression discontinuity design approach". We estimated peer effects in undergraduate students' academic performance at a Brazilian university. Our empirical evidence comes from a micro data set containing information of 1550 undergraduate students enrolled in 27 courses at the Federal University of Ceará. In light of this great courses availability, we assign each course into one of four categories depending on its admited students' results at the entrance exam. Then, we proceed the estimation exercise using a multi-treatment effect model. In this fashion, using $IRA$ as a measure of academic performance, we obtain a negative effect (-0.19) for being in a first semester class, which means a 2% smaller $IRA$ for firt semester students, *vis-a-vis* members of second semester classes. Moreover, we found non-linearities in this effect, since, for example, it ranges between 0.5 to -0.18. This results are in accordance with Sacerdote (2001) and Zimmerman (2003), also finding non-linearities in "peer effects".

# 2 SPATIAL DETERMINANTS OF URBAN RESIDENTIAL WATER DEMAND IN FORTALEZA, BRAZIL

## 2.1 INTRODUCTION

The literature on residential water demand estimation has grown considerably, indicating the main variables that affect water consumption and the estimation techniques to be employed. However, that line of research has not been yet capable of fully exploring how spatial effects might influence water demand. Franczyk and Chang (2009) point out that " water consumption standards cannot be explained by economic and population growth only, but also through biophysical and socioeconomic factors that usually have spatial dependence". Following the same line, House-Peters, Pratt and Chang (2010) suggest that "residential water consumption is not affected by climate, socioeconomic and physical variables only. It is also affected by geographical location and its interaction with nearby regions."

Therefore, incorporating spatial effects into the analysis of residential water demand could provide a wider and more accurate explanation on its consumption variations. Papers like Chang, Parandvash and Shandas (2010), Wentz and Gober (2007), Franczyk and Chang (2009), House-Peters, Pratt and Chang (2010), Ramachandran and Johnston (2011) have recently included spatial effects in their studies, increasing the significance of their models when compared to other models that do not consider such effects. Based on that series of papers, we believe that our endeavor has its own merits. Firstly, because estimates on water micro-demand models with spatial effects are quite new in international literature and absent in the national research. Secondly, because by aggregating new methodological procedures we can better understand the factors that affect residential water demand.

Therefore, this paper aims at analyzing water demand using spatial econometric techniques in an exploratory way. For this, we have at our disposal information from a study field in the city of Fortaleza, Brazil (The state capital of Ceará, located in Northeastern Brazil - WGS84 coordinates $3^0 43' 6'' South$ and $38^0 32' 34'' West$). The city has a population of about to 2.5 million and it is the fifth largest city in Brazil with an area of 313 square kilometers, boasting one of the highest demographic densities in the country (8,001 per $km^2$).

From a series of test procedures and econometric exercises, we can confirm the importance of considering spatial effects, since the exclusion of those effects underestimate the impact that income and the number of bathrooms per residence can have over water

demand. More importantly, it underestimate considerably the impact of average and marginal prices on water demand. Although our results are of an exploratory nature, we believe they will enable us to better understand how space matters for water-micro demand estimation.

Besides this introduction and a section discussing final considerations, this paper has four more sections. Section 2 offers a brief literature review on residential water micro-demand estimation. Section 3 introduces the database used and the results of an exploratory spatial data analysis. Section 4 sets up the demand function to be estimated with a non-linear tax structure. We also introduce the econometric models used in this paper together with the tests used for (mis)specification analysis. Finally, Section 5 shows the results and Section 6 draws some final considerations.

## 2.2 WATER MICRO DEMAND ESTIMATION WITH SPATIAL EFFECTS

To understand the way in which charging water use affects its consumption, it is necessary to know the factors that determine water demand. Since Gottlieb (1963) and Howe and Linaweaver (1967), several researchers in many countries have carried out studies to estimate a residential water demand function for their particular regions in order to provide technical work as a support to implement policies aimed at controlling and promoting its rational use and preservation.

Agthe, Billings and Dobra (1986) for Tucson, Arizona (USA), Rietveld, Rouwendal and Zwart (2000) in Indonesia, Polycarpou and Zachariadis (2013) in Cyprus, and Miyawaki, Omori and Hibiki (2013) for the cities of Tokyo and Chiba, are just some representative examples. In Brazil, literature on residential water demand estimation is still new. One of the first papers that approached this issue was written by Andrade et al. (1995). For the city of Piracicaba, São Paulo, Mattos (1998) and Melo and Neto (2007) estimated the function of residential water demand for Northeastern Brazil.

Although very heterogenous in terms of methodologies and scopes, all studies briefly cited above share an important shortcoming: they do not consider spatial effects in water demand estimation. Aware of this problem, some authors recently started to include spatial effect in their analysis, seeking to explain the spatial association pattern for water consumption. Wentz and Gober (2007) used GWR model (Geographic Weighted Regression) in a study for Phoenix, USA, in order to verify if there was any additional spatial effect contribution to the results obtained through the OLS model (Ordinary Least Square). The authors verified through the GWR model that the importance of spatial effects reduces to two the variables that determine water demand. The variables are residence size and the existence or not of a swimming pool in the property.

In the state of Oregon (USA), Franczyk and Chang (2009) realized that water demand was not just related to population and economic growth, but also to other biophysical and socioeconomic factors that in general present spatial dependence. The authors used the spatial error model (SEM), besides the OLS model in order to include spatial autocorrelation effects in the study. They applied Moran-I statistics and showed that there is spatial dependence on errors.

In the city of Portland (Oregon, USA) Chang, Parandvash and Shandas (2010) identified a spatial association pattern for water demand. They verified that the areas where water consumption was higher coincided with the areas in which average home sizes were larger and both the building density and property age averages were low. House-Peters, Pratt and Chang (2010) carried out a study for the city of Hillsboro (Oregon, USA) analyzing climate effects on water demand. Using spatial analysis techniques, the authors found that although water demand in that area was not sensitive to dry conditions at all, some specific areas presented higher water consumption levels under such conditions.

Ramachandran and Johnston (2011) studied how the spatial effect influenced residential water demand for external use in the city of Ipswich (Massachusetts, USA) while a restricted use of water policy was being implemented. They argued that decisions on house landscapes, and therefore, the use of water in order to maintain these landscapes would depend on economic factors such as if the landscape affects the house selling price or social factors such as imitation reasons, as people tend to copy the landscaping and vegetation used in gardens of nearby residences.

## 2.3 DATA SET

### 2.3.1 The Sample

The database contains information from a scientific project carried out by a group of researchers from UECE and UFC (Ceará State University and Ceará Federal University) requested by CAGECE (Ceará Water and Sewage Company). It collected more than 3,000 questionnaires containing information on socioeconomic and physical characteristics from different households in Fortaleza. After deletion of missing observations, we end up with 2,891 usable observations, as shown in Table 2.1.

The data introduced shows that residential water consumption in February 2007 for the city of Fortaleza was 16.41 $m^3$, on average. The median of 14 $m^3$ indicates that half of residences in Fortaleza are either in CAGECE's first or second consumption block ([0 $m^3$, 10 $m^3$] or (10 $m^3$, 15 $m^3$]). This might not be good for CAGECE if the tax structure is poorly designed. As for the socioeconomic characteristics, on average, each household has 2.09 and 1.76 male and female residents respectively. Families have an average monthly income of 2.43, which indicates that they spread out between two income classes: class 2

Table 2.1 – Descriptive Statistics

|                            | Mean  | Std. Dev. | Min   | Median | Max    |
|----------------------------|-------|-----------|-------|--------|--------|
| Water Consumption ($m^3$)  | 16.41 | 9.71      | 2.00  | 14.00  | 60.00  |
| Effective Price (R$)       | 24.49 | 21.44     | 9.80  | 16.04  | 159.35 |
| Average Price (R$)         | 1.45  | 0.56      | 0.98  | 1.25   | 4.90   |
| Marginal Price (R$)        | 1.61  | 1.02      | 0.10  | 1.56   | 4.95   |
| Difference (R$)            | 9.46  | 18.68     | -8.71 | 5.80   | 137.65 |
| Family Income (class)      | 2.43  | 1.04      | 1.00  | 2.00   | 5.00   |
| Type of Property (class)   | 2.69  | 0.55      | 1.00  | 3.00   | 4.00   |
| Male Residents (number)    | 2.09  | 1.16      | 0.00  | 2.00   | 8.00   |
| Female Residents (number)  | 1.76  | 1.12      | 0.00  | 2.00   | 8.00   |
| Bathrooms (number)         | 1.55  | 0.86      | 0.00  | 1.00   | 8.00   |
| Gardens (dummy)            | 0.24  | 0.43      | 0.00  | 0.00   | 1.00   |

Source: Elaborated by the Authors

(whose people earn from one minimum wage (R$ 350.00 or US$ 219.45 - Purchase Parity Power) up to 2 minimum wages), and class 3(those earning from 2 to 5 minimum wages).

With regards to the physical characteristics of households, the residences have, on average, 1.55 bathrooms and 24% of them have a garden. The average type of property is 2.69, which allows us to classify residences between the medium and regular categories according to CAGECE's standards. In terms of pricing, CAGECE applies an increasing tariff system through consumption blocks: $[0, 10]$, $(10, 15]$, $(15, 20]$, $(20, 50]$ and $(50, \infty)$. The rates, in 2007, for each block were: R\$ 9.80 (Fixed fee), 1.56 $R\$/m^3$, 1.65 $R\$/m^3$, 2.80 $R\$/m^3$ and 4.95 $R\$/m^3$, respectively. In next section we shall introduce a spatial exploratory analysis applied to our data set.

## 2.3.2 Spatial Exploratory Data Analysis

In order to check the hypothesis that spatial effect plays an important role to explain residential water demand, we will verify if water consumption presents any spatial association pattern at all. Therefore, we stick to the literature and use the Moran-I statistic to test for global spatial association and the local Moran-I statistic to test for local spatial association, besides the significances and clusters maps (see, Anselin (1995)).

Moran-I statistics for five well know weighting matrices (distance, 5, 10, 15 and 20 nearest neighborhood) were calculated. All figures for the Moran I statistic belong to the interval $(0.0, 0.15)$. These values exceed their statistical averages but they are close to zero, which apparently indicates no spatial autocorrelation in water consumption. However, although these values are close to zero, they are statistically different from zero, once the pseudo-p-value is extremely low (in fact it is undistinguishable from zero). That is an indication that we cannot reject the hypothesis of lack of positive spatial autocorrelation, even if in a reduced magnitude and for any common weighting matrix. These first results prompt us to carry on.

Not always the global pattern of spatial association reflects the local pattern of spatial association, though. In this sense, LISA indexes (Local Indicator Spatial Association)

are used to overcome this obstacle and capture local patterns of linear association. The most well known LISA statistic, the local Moran-I, is derived from a global indicator of autocorrelation that decomposes the local contribution of each observation into four categories.

The results of dispersion diagram show that there is a tendency to a positive autocorrelation with the observations distributed in the first and third quadrants for all weighting matrices, however with a lower value (0.0487) for the distance weighting matrix[1].

As for the significance and dispersion maps, they all give similar results for all weighting matrices: there is a high concentration of water consumption at the top center in the city of Fortaleza that covers the downtown area and the richest neighborhoods in the city, as well as low consumption clusters in suburban areas[2]. Hence, such exploratory spatial analysis confirmed our idea that not only there are global and local spatial autocorrelation patterns in our sample but also that such spatial autocorrelation might be important. Next section deals with the demand function models specifications and estimations.

## 2.4 ECONOMETRIC MODEL

### 2.4.1 Non-Spatial Specification

It is well known that in a setup with non linear prices, the consumer's budget restriction will be non linear as well. In such cases, the solution to the optimization problem faced by a consumer, i.e. maximization of utility given (non-linear) budget constraint, will give us the water demand as a function of prices and income, according to Moffitt (1986). However, the residential water demand is not a function of water price and consumers income only.

We need to add other variables, which are important to explain residential water demand. Although there is no consensus over the "best" econometric specification for modeling household water demand, we stick to two of the most prominent ones (see, Arbués, García-Valiñas and Espiñeira (2003), Olmstead, Hanemann and Stavins (2007) and Worthington and Hoffman (2008)):

$\ln(\text{QC}_i) = \beta_1 + \beta_2 ln(Pavg_i) + \beta_3 Inc_i + \beta_4 Male_i + \beta_5 Female_i + \beta_6 Bath_i + \beta_7 Garden_i + \varepsilon_i$

$\ln(\text{QC}_i) = \beta_1 + \beta_2 ln(Pmg_i) + \beta_3 Diff_i + \beta_4 Inc_i + \beta_5 Male_i + \beta_6 Female_i + \beta_7 Bath_i + \beta_8 Garden_i + \varepsilon_i$

where,

- $QC$ = Amount of consumed water in February 2007 in $m^3$

- $Pavg$ = Average price in February 2007

---

[1]    Moran-I statistics and all dispersion diagrams can be obtained from the authors upon request.
[2]    Both maps can be obtained from the authors upon request. A possible explanation for the configuration of such clusters is that the income distribution in the city of Fortaleza is very unequal.

- $Pmg$ = Marginal price in February 2007

- $Diff$ = Difference variable[3]

- $Inc$ = Family Income

- $Male$ = Number of male residents in the household

- $Female$ = Number of female residents in the household

- $Bath$ = Number of bathrooms in the household

- $Garden$ = Dummy for the presence of a garden in the household

- $\varepsilon$ = Error term

We decided to start our modeling[4] exercise with both the average price and marginal price coupled with the difference variable specifications based on the following premises. Firstly, the average price *versus* marginal price (with difference) continues to be an open issue, yet to be settled. Hence, from a methodological point of view it is good practice to rely on statistical methodology and not on any *ad hoc* personal choice of specification *ex ante* the modeling exercise.

Secondly, we agree with Saleth and Dinar (2001) when these authors claim that the average price *versus* marginal price issue has not been casted in a correct way when stressing the question of the lack of perfect information on the water tariff structure or the inexpressive value of water bill compared to household total income. Rather, Saleth and Dinar (2001) argue quite convincingly that *the price perception debate is not as much of a controversy on the price specification itself as it is with regards to the relative relevance of the positive [Pavg] versus normative [Pmg and Diff] approach to consumer behavior under block rate pricing.*

Although we may end up choosing the "best" specification, the ... *comparison of demand functions under these prices can be used to at least show the effects of the change in price levels due to a shift in the price perception.* Therefore, the issue on average price *versus* marginal price has important behavioral implications beyond the simply traditional econometric specification debate, resulting in an almost necessary topic to deal with by estimating both specifications.

---

[3]   The difference between the bill that would result if each $m^3$ of water consumed was priced by the marginal price and the actual bill. For increasing block tariffs, the difference is negative for households located on the first block, meaning that their water consumption receives subsidies. See, among others, Nordin (1976)

[4]   Note that specifications like equations 2.4.1 and 2.4.1 are by no means the only ones. For example, there is growing interest in modeling water demand as composed of two parts; a fixed and a residual component, seeking to capture consumption niches that are non-responsive to pricing (see, Dharmaratna and Harris (2012)).

The choice of socioeconomic variables and the physical characteristics of residences agree with the main studies carried out on water demand estimation, very well summarized in Arbués, García-Valiñas and Espiñeira (2003). We expect that *family income*, *number of male residents* and *number of female residents* in the household, as well as the *number of bathrooms* shall exert a positive effect on water demand, since an increase in these variables will increase water demand. As for the independent variable *presence of garden*, we also expect a positive value; however this variable will play an important role when discussing spatial effects. Finally, with regards to average, marginal and difference price, it is expected that water reacts as a normal good.

It is known that the estimation of a demand function in a non linear tax context creates, a priori, a problem of endogeneity. We are aware of the potential deleterious impacts of endogeneity in our estimates (especially on the average price specification) as well as of authors that find no significant difference between simple least square estimations and instrumental variable approaches such the one carried out by Jones and Morris (1984). However, we do control for endogeneity issues, although in a rather traditional way, by estimating the marginal price model through a methodology developed by McFadden, Puig and Kirschner (1977).

Before we proceed to discuss spatial specification issues, it is important to stress the fact that even though Fortaleza is a city located in a developing country, its urban water market is very similar to those of cities located in developed countries. Besides being a large and dense urban center, it has been served by the same water company since 1970, and its customers are used to their tariff system. Also, the market presents an index of hydrometration above 98%. However, when we move away from the coast, towards inner cities in the state of Ceara, the water demand and supply conditions can change quickly and drastically. On such settings, our model could be a considerable specification error[5].

## 2.4.2 Spatial Specification

To verify if the inclusion of spatial effects affect residential water demand, we used three models (see, Anselin (1988)): SEM (Spatial Error Model), which is used when we believe that spatial dependence is caused by autocorrelation in error terms; mixed SAR (Spatial Autoregressive), that aggregates explicative variables and it is used when the spatial dependence is contained in the dependent variable and finally, the SARMA model (Spatial Autoregressive and Moving Average), that is used when we believe that spatial dependence is contained both in error terms and in the dependent variable. The SARMA model is represented by:

---

[5]    We would like to thank a referee for pointing out that to us.

$$Y = \rho W_1 Y + X\beta + \varepsilon \qquad (2.1)$$

$$\varepsilon = \lambda W_2 \varepsilon + u \qquad (2.2)$$

$Y$ is an $n \times 1$ vector that contains observations on water demand in logarithms. $X$ is an $n \times m$ vector of explicative variables, the same used in previous models, and $\beta$ is an $m \times 1$ parameter vector to be estimated, where $m$ is the number of independent variables and $n$, the sample size. $W_1$ and $W_2$ are the spatial weighting matrixes, $u$ is the random error term in standard normal distribution with mean equal zero and a constant variance, and $\lambda$ is the autoregressive parameter associated to error term. Finally, $\rho$ is the autoregressive parameter associated to the lagged dependent variable.

In order to help us decide which of the three specifications capture in a more accurate way the spatial effect over residential water demand, we applied Lagrange multipliers tests, both for lag ($LM_\rho$) and for spatial error ($LM_\lambda$), as well as their robust Lagrange multipliers (RLM) versions. To detect the correct functional form, Florax, Folmer and Rey (2003) suggest the use of the "hybrid identification" strategy, using both the classical and robust tests for spatial autocorrelation.

### 2.4.3   Why Spatial Effects in Water Demand?

As logical as SARMA model might appear, it subsumes a host of possible theoretical or, according to econometric parlance, "structural", explanations for its channels of causation. However, establishing clear cut causal linkages for spatial models is not an easy task. In fact, according to Corrado and Fingleton (2012), literature on spatial statistics, as well as spatial econometrics, appear to be dominated by data-analytic considerations only during the model specification phase, to the detriment of causal modeling. However, data-driven protocols are indispensable approaches to perform, especially during the exploratory analysis of statistical and econometric models. A unique reliance on data-analytic considerations trades off against a better understanding of the important behavioral and policy implications of the model. Consequently, a more equilibrated modeling strategy has to be chosen and this requires a justification on how spatial effects might be important for water demand estimation.

The justification for the use of the SARMA model comes from the belief that the spatial effect might work through both the error terms and/or lags of the dependent variable. Such factors would be the climate-related, biophysical, socioeconomic and geographical, as well as the infrastructure of the water distribution system. Two theoretical justifications for spatial effects that have gained wider acceptance are: i) imitation of consumption in neighboring residences, and ii) water supply network dependencies.

Some authors such as Ramachandran and Johnston (2011) believe that there is imitation of water consumption in neighboring residences, especially in gardening activities, due

to the attempt to imitate the shape and type of plants used by neighbors. Although a seminal idea, their papers fall short after computing descriptive spatial dependence indexes. Others, like Wentz and Gober (2007) and Janmaat (2013) found similar effects. Janmaat (2013) calls that an *emulation effect* in water use behavior after modeling water demand for the city of Okanagan, Canada, through a geographically weighted regression. Observe, however, that he is very careful to imply a more elaborated causal link beyond asserting that *I do not have an explicit theory on how neighbors influence each other ... beyond neighbors noticing each other's water use.* Anyway, we conjecture that imitation or emulation is a possible effect that incorporate (positive) water consumption spatial dependence.

Another possible justification for spatial effects comes from the infrastructure of water distribution network systems. A network may create a negative consumption autocorrelation, once the pressure over the distribution system causes a given residential consumption that affects the consumption of nearby residences. Such channel of spatial effect has a much longer history (see Jones and Morris (1984) for a justification along these lines).

## 2.5 RESULTS

### 2.5.1 Non-Spatial Specification

Table 2.2 presents first the results related to the econometric model for residential water demand function with no spatial effects. We estimated three specifications: Average Price (AV model), Marginal Price *cum* Difference (MP model), and Marginal Price *cum* Difference with McFadden (McFadden model) method. According to results, the estimated coefficients for all variables (excluding log(Pavg), log(Pmg), Diff and Garden) showed expected positive signals and are statistically significant. However, there are important intra and inter-models differences.

The AP model presents quite intuitive estimated parameters and an overall fit ($R^2 = 0.17$) compatible with estimations of models based on micro-data sets. The elasticity of the average price is negative (-0.3503) and conforms to past empirical exercises. All figures are in accordance with theoretical predictions. Water is a (slightly) normal good, as reflected by the estimated parameter of Income (0.0631). Male and Female exerts a different impact on water demand with Females (0.0850) consuming less than males (0.1140). The number of bathrooms, as expected, have a positive impact (0.1351) on water demand *ceteris paribus* as well as the presence of garden (0.0530). the Garden variable will play an import role when discussing channels of spatial effects. In the meantime, let us comment on the MP model.

Overall, the MP model presents estimated parameters for common variables with

Table 2.2 – Estimates Average Price, Marginal Price and Mc Fadden, No Spatial Effects

| | AV | | MP | | McFadden | |
|---|---|---|---|---|---|---|
| | Estimate | S. E. | Estimate | S. E. | Estimate | S. E. |
| Intercept | 1.9713*** | 0.0365 | 2.1329*** | 0.0256 | 1.7578*** | 0.0367 |
| log(Pavg) | −0.3503*** | 0.0364 | - | - | - | - |
| log(Pmg) | - | - | 0.056*** | 0.0096 | −0.4098*** | 0.0326 |
| Diff | - | - | 0.0216*** | 0.0006 | 0.0392*** | 0.0013 |
| Income | 0.0631*** | 0.0114 | 0.0239** | 0.0080 | 0.0447*** | 0.0079 |
| Male | 0.1140*** | 0.0092 | 0.0573*** | 0.0065 | 0.1119*** | 0.0076 |
| Female | 0.0850*** | 0.0094 | 0.0318*** | 0.0067 | 0.0657*** | 0.0069 |
| Bath | 0.1351*** | 0.0139 | 0.0261** | 0.0098 | 0.0456*** | 0.0097 |
| Garden | 0.0530* | 0.0257 | 0.0147 | 0.0180 | 0.0249 | 0.0176 |
| adjusted $R^2$ | 0.1760 | | 0.5980 | | 0.6150 | |
| $F$ statistic | 104 | | 616 | | 660 | |

Source: Elaborated by the Authors
Note: Signif. codes: 0 *** 0.001 ** 0.01 * 0.05 . 0.1

regards to the AP model (say, Income, Male, Female, Bath and Garden) that have the right signal and are statistically significant (except for Garden), although with much less size. Interestingly, we are able to find the elusive intramarginal effect (see, Nordin (1976)), since we cannot reject equality between the estimated parameters of Diff (0.0216) and Income (0.0239). Despite these initial achievements, the MP model shows a key weakness: the coefficient of log(Pavg) is positive and significant! This is so despite the much "better" $R^2$ and $F - statistic$ compared to the AP model. Therefore, we go straight to endogeneity issues and this lead us to the McFadden model.

The estimated parameters of the McFadden model, not surprisingly, are quite different from the ones of the MP model. Overall, they all inflate the values. The intra-marginal effect is preserved, as again, we cannot reject equality between the estimated parameters of Diff (0.0392) and Income (0.0447). The good news is the sensible estimated effect of the elasticity of marginal prices (-0.4098). A Hausman test (-14.9518) rejects thoroughly the null of exogeneity under any sensible level of significance. From this point on, we feel confident to eliminate the MP specification and proceed comparing only the AP and McFadden models[6].

The most striking, although not necessarily surprising result is the small difference between the elasticity of average price (-0.3503) and the elasticity of marginal price (-0.4098). Also, all common estimated parameters present similar results and are statistically significant, except the variable Garden that is both lower and not significant in the McFadden specification.

---

[6]  We would like to stress that we conducted the famous Oppaluch testing approach (see, Opaluch (1984)) and could not discard either specification, say, AP and MP. However, we back our choice on pragmatic grounds reflected on the results from the AP and McFadden model.

## 2.5.2   Spatial Specification

Since there is room for spatial dependence on water consumption, we run the Moran-I test for the residues estimated for both the AP and McFadden models. The Moran-I statistics is significant for both models and weighting matrices types[7]. This means that the probability for the spatial association pattern being random is close to zero, supporting the hypothesis that the residues are spatially dependent. Moreover, the positive value indicates that the autocorrelation is positive, as expected due to an "imitation" channel of spatial causation.

After confirming the presence of spatial autocorrelation in the residues, we run Lagrange multipliers tests in their classic and robust versions to define which model is more appropriate. Following the methods proposed by Florax, Folmer and Rey (2003), we compare the $LM_\lambda$ and $LM_\rho$ values first[8]. The values are significant for both the AP model and the McFadden model. This indicates that there is spatial dependence associated both to lag in the dependent variable as well as to non-modeled effects, the latter represented by error term. This means that the SAR specification should be estimated. However, after analyzing the SARMA tests results, we can see that the SAR in not the best model. In fact, they point out that the SARMA model is the correct way to model spatial effect on residential water demand in the city of Fortaleza.

The results[9] shown in Table 2.3 demonstrate that the coefficients for both models are rather comparable. Except for the variable Bathrooms, they agree on sign and size and are all statistically significant. The elasticity of water consumption with regards to prices is very similar, with a slightly higher value (module) for the MaFadden model (0.4090) compared to the AV model (0.33654).

The specific spatial parameters, say $\rho$ (0.3541) and $\lambda$ (-0.2517) are both significant for the AP model, which backs the explanation based on imitation effects for the spatial dependence. Also consider the fact that the variable Garden became non-significant once we decided to include a spatial lag of the dependent variable. However, the fact that $\lambda$ is negative cannot be underestimated. This might be the result of network effects not controlled by observed covariates, operating through the error term. For the McFadden model $\rho$ (0.0524) (borderline significant ($p - value = 0.11$)) as well as $\lambda$ (0.1431). Again a positive $\rho$ makes more convincing explanations of spatial effects based on imitation of behavior. The positive and significant effect represented by $\lambda$ has the opposite sign compared to the AP model. We are not able to rationalize these differences and we believe it is more interesting to move forward and compare our spatial results with the non-spatial

---

[7]   However, the matrix **5 Nearest Neighbors** presented the highest value of that statistic. Hence, from now on, all econometric manipulations will consider only that weighting matrix choice.

[8]   Both the Moran-I statistics and Lagrange Multiplier tests can be obtained from the authors upon request.

[9]   It is worth mentioning that in the first round of estimations, the variable Garden was not significant, so we remove it and estimated the model again.

Table 2.3 – Estimates Average Price and McFadden, with Spatial Effects

| | Average Price | | Mc Fadden | |
|---|---|---|---|---|
| | Estimate | S. E. | Estimate | S. E. |
| Intercept | 1.1166*** | 0.1122 | 1.6403*** | 0.0911 |
| log(Pavg) | −0.3365*** | 0.0348 | - | - |
| log(Pmg) | - | - | −0.4090*** | 0.0318 |
| Diff | - | - | 0.0389*** | 0.0012 |
| Income | 0.0535*** | 0.0104 | 0.0442*** | 0.0081 |
| Male | 0.1096*** | 0.0087 | 0.1102*** | 0.0074 |
| Female | 0.0810*** | 0.0089 | 0.0650*** | 0.0069 |
| Bath | 0.1194*** | 0.0131 | 0.0424*** | 0.0098 |
| Rho | 0.3541*** | 0.0463 | 0.0524 | 0.0332 |
| lambda | −0.2517*** | 0.0794 | 0.1431*** | 0.0422 |
| LR test | 59.705*** | | 45.84*** | |
| Log likelihood | -2432.501 | | -1351.896 | |

Source: Elaborated by the Authors

Note: Signif. codes: 0 *** 0.001 ** 0.01 * 0.05 · 0.1

estimates next.

In spatial models, changes in an independent variable of a specific location ($\Delta x_i$) impact on a given dependent variable ($\Delta y_i$) (direct effect). However, since the dependent variables of other locations depend on $y_i$, by means of the weighting matrix, there will be a change on other dependent variables ($\Delta y_j, j \neq i$) which for the same reason will affect $y_i$ (indirect effect). To address this issue, Table 2.4 shows the total effects (direct + indirect effect) for the AP and McFadden models respectively, through a procedure implemented by means of the **spdep** package developed by Bivand and contributions (2011).

Table 2.4 – Total Impact on AP and McFadden Models

| | Average Price | | | Mc Fadden | | |
|---|---|---|---|---|---|---|
| | Direct | Indirect | Total | Direct | Indirect | Total |
| log(Pavg) | -0.3560 | -0.0807 | -0.4367 | - | - | - |
| log(Pmg) | - | - | - | -0.3993 | -0.0651 | -0.4644 |
| Diff | - | - | - | 0.0385 | 0.0063 | 0.0448 |
| Income | 0.0599 | 0.0136 | 0.0735 | 0.0408 | 0.0067 | 0.0475 |
| Male | 0.1130 | 0.0256 | 0.1386 | 0.1101 | 0.0180 | 0.1280 |
| Female | 0.0860 | 0.0195 | 0.1055 | 0.0659 | 0.0108 | 0.0767 |
| Bath | 0.1286 | 0.0292 | 0.1578 | 0.0389 | 0.0063 | 0.0452 |

Source: Elaborated by the Authors

Observe that the package procedure to estimate the total effect in **spdep** is only correct asymptotically. Therefore, values appearing on the column labeled "Direct" in Table 2.4 are slightly different from the estimated parameters entered in Table 2.3. The total effects are clearly different from their corresponding effects (the estimated parameters) appearing in Table 2.3. Although all estimated (total) effects in Table 2.4 remain with the same signal, figures are considerable different from those that do not take into consideration this subtle but very important issue on interpretation of estimated parameters from spatial models.

We still have to see how the inclusion of the spatial effect changes the estimated effects, *vis a vis* the model without spatial effects. Table 2.5 does exactly that by comparing these effects with those from the non-spatial estimation (see, Table 2.2).

Table 2.5 – Percentual Variation (%) - $100 \times \frac{\widehat{\beta}_{TotalSpatial} - \widehat{\beta}_{NonSpatial}}{\widehat{\beta}_{NonSpatial}}$

|  | Average Price | Mc Fadden |
|---|---|---|
| log(Pavg) | 24.66 | - |
| log(Pmg) | - | 13.32 |
| Diff | - | 14.28 |
| Income | 16.48 | 6.26 |
| Male | 21.57 | 14.38 |
| Female | 24.11 | 16.74 |
| Bath | 16.80 | -0.87 |

Source: Elaborated by the Authors

Now, after calculating the correct total effects, all estimated parameters from the non-spatial models (except that from variable Bath, on the McFadden model) are in fact underestimated! Indeed, the underestimation is by no means negligible[10]. For instance, for the price-elasticity effect, these figures are 24.66% and 13.32% for the AP and McFadden models respectively. The same happens for the Income variable. Clearly, the absence of spatial effects appears to be an important shortcoming in water demand estimation, at least if one is using a micro data set.

## 2.6   FINAL CONSIDERATIONS

This paper sought to apply a new methodological approach to estimate residential water demand models in a large urban center in Brazil that includes spatial effects in the analysis. We showed that the determinant factors explaining residential water consumption in the city of Fortaleza are average price, marginal price, difference, income, number of male and female residents and total number of bathrooms per residence, as long as we add spatial effects. Most importantly, our results point out that not considering spatial effects might be a key specification error in water micro-demand analysis.

Our empirical methodology built a sort of detailed approach showing the main steps on how to start from a non-spatial model and achieve a "good" econometric model with spatial effects. Through this approach, we are not able to discard neither the average price model nor the marginal price model *a la* McFadden. We see that as an advantage in the sense that rather than focusing on having a necessary unique choice of specification, keeping a "dichotomy" between these two models might be a sensible way to approach the problem.

---

[10]   We thank a lot a referee for prompting us to get deeper on the real difference, in terms of estimated parameter magnitudes, between our spatial model and the traditional non-spatial models.

As expected, for both spatial and non-spatial specifications, the average and marginal prices variables had a negative impact on water consumption. Also, water behaved as a normal good. Income, total of male and female residents, and total of bathrooms resulted in a positive effect. As to the long debate on endogeneity issues, we found no considerable differences between the AP and McFadden models. Interestingly, we were able to find the intramarginal effect (see Nordin (1976)).

Lagrange multipliers and SARMA tests showed both in classic and robust versions that the "best specification" to estimate residential water demand is the SARMA model, instead of the SEM. We address that by estimating a SARMA model for both the average price and the McFadden procedure. Now, after correcting the direct and indirect effects of the estimated parameters, the advantage of using a spatial approach appears to be more evident. Not including spatial features underestimates almost all variables in absolute terms when compared to their non-spatial counterparts. For instance, including spatial effects increases the price-elasticity in the AP price in 24.66% and the price-elasticity for the McFadden model in 13.32%!

As suggestions for future studies, we believe that both the incorporation of spatial heterogeneity and the inclusion of water quality variables are worth pursuing. Also, a detailed study of spatial effects on markets that are not well served by water companies and that rely on alternative non-market sources of water seems to be a mandatory task. Another interesting line of research would be applying spatial models to longitudinal data. Finally, replicating our empirical exercise on different data sets coming from different institutional backgrounds might be something worth pursuing in order to validate our approach.

# REFERENCES

AGTHE, D. E.; BILLINGS, R. B.; DOBRA, J. L. A simultaneous equation demand model for block rates. *Water Resources Research*, v. 1, p. 1 – 4, 1986.

ANDRADE, T. et al. Saneamento urbano: a demanda residencial por água. *Pesquisa e Planejamento Econômico*, v. 25, n. 3, p. 427–448, 1995.

ANSELIN, L. *Spatial econometrics: methods and models.* [S.l.: s.n.], 1988. 304 p. ISBN 978-90-247-3735-2.

ANSELIN, L. Local indicators of spatial association-lisa. *Geographical analysis*, v. 27, n. 2, p. 93–115, 1995.

ARBUÉS, F.; GARCÍA-VALIÑAS, M. A.; ESPIÑEIRA, R. M. Estimation of residential water demand: a state-of-the-art review. *Journal of Socio-Economics*, v. 32, n. 1, p. 81–102, 2003.

BIVAND, R.; CONTRIBUTIONS with. *spdep: spatial dependence: weighting schemes, statistics and models.* [S.l.], 2011. R package version 0.5-40. Disponível em: <http://CRAN.R-project.org/package=spdep>.

CHANG, H.; PARANDVASH, G. H.; SHANDAS, V. Spatial variations of single-family residential water consumption in portland, oregon. *Urban Geography*, v. 31, n. 7, p. 953–972, 2010.

CORRADO, L.; FINGLETON, B. Where is the economics is spatial econometrics? *Journal of Regional Science*, Wiley Blackwell (Blackwell Publishing), v. 52, n. 2, p. 210–239, May 2012.

DHARMARATNA, D.; HARRIS, E. Estimating residential water demand using the stone-geary functional form: the case of sri lanka. *Water Resources Management*, Springer Netherlands, v. 26, n. 8, p. 2283–2299, 2012. ISSN 0920-4741.

FLORAX, R.; FOLMER, H.; REY, S. Specification searches in spatial econometrics: the relevance of hendry's methodology. *Regional Science and Urban Economics*, Elsevier, v. 33, n. 5, p. 557–579, 2003.

FRANCZYK, J.; CHANG, H. Spatial analysis of water use in oregon, usa, 1985-2005. *Water Resources Management*, v. 23, n. 4, p. 755–774, 2009.

GOTTLIEB, M. Urban domestic demand of water in the united states. *Land Economics*, v. 39, n. 2, p. 204–210, 1963.

HOUSE-PETERS, L.; PRATT, B.; CHANG, H. Effects of urban spatial structure, sociodemographics, and climate on residential water consumption in hillsboro, oregon. *JAWRA Journal of the American Water Resources Association*, v. 46, n. 3, p. 461–472, 2010.

HOWE, C. W.; LINAWEAVER, F. P. The impact of price on residential water demand and its relation to system design and price structure. *Water Resources Research*, v. 3, n. 1, p. 13, 1967.

JANMAAT, J. Spatial patterns and policy implications for residential water use: An example using kelowna, british columbia. *Water Resources and Economics*, Elsevier, v. 1, p. 3–19, Jan 2013.

JONES, C. V.; MORRIS, J. R. Instrumental price estimates and residential water demand. *Water Resources Research*, Wiley Blackwell (John Wiley &amp; Sons), v. 20, n. 2, p. 197–202, Feb 1984.

MATTOS, Z. d. B. Uma análise da demanda residencial por água usando diferentes métodos de estimação. *Pesquisa and Planejamento Econômico*, v. 28, n. 1, p. 207–224, 1998.

MCFADDEN, D.; PUIG, C.; KIRSCHNER, D. Determinants of the long-run demand for electricity. In: *Proceedings of the American Statistical Association*. [S.l.: s.n.], 1977. v. 1, n. 1, p. 109–19.

MELO, J.; NETO, P. Estimação de funções de demanda residencial de agua em contexto de preços não-lineares. *Pesquisa and Planejamento Econômico*, v. 37, n. 1, p. 149–173, 2007.

MIYAWAKI, K.; OMORI, Y.; HIBIKI, A. Exact estimation of demand functions under block rate princing. *Econometric Reviews*, 2013.

MOFFITT, R. The econometrics of piecewise-linear budget constraints: a survey and exposition of the maximum likelihood method. *Journal of Business & Economic Statistics*, JSTOR, v. 4, n. 3, p. 317–328, 1986.

NORDIN, J. A proposed modification of taylor's demand analysis: comment. *The Bell Journal of Economics*, v. 7, n. 2, p. 719–721, 1976.

OLMSTEAD, S. M.; HANEMANN, W. M.; STAVINS, R. N. Water demand under alternative price structures. *Journal of Environmental Economics and Management*, Elsevier, v. 54, n. 2, p. 181–198, Sep 2007.

OPALUCH, J. J. A test of consumer demand response to water prices: reply. *Land Economics*, v. 60, n. 4, p. 417–421, 1984.

POLYCARPOU, A.; ZACHARIADIS, T. An econometric analysis of residential water demand in cyprus. *Water Resources Management*, Springer Netherlands, v. 27, n. 1, p. 309–317, 2013. ISSN 0920-4741.

RAMACHANDRAN, M.; JOHNSTON, R. J. Quantitative restrictions and residential water demand : a spatial analysis of neighborhood effects. 2011.

RIETVELD, P.; ROUWENDAL, J.; ZWART, B. Block rate pricing of water in indonesia: an analysis of welfare effects. *Bulletin of Indonesian Economic Studies*, Informa UK (Taylor &amp; Francis), v. 36, n. 3, p. 73–92, Dec 2000.

SALETH, R. M.; DINAR, A. Preconditions for market solution to urban water scarcity: empirical results from hyderabad city, India. *Water Resources Research*, Wiley Blackwell (John Wiley &amp; Sons), v. 37, n. 1, p. 119–131, Jan 2001.

WENTZ, E. a.; GOBER, P. Determinants of small-area water consumption for the city of phoenix, arizona. *Water Resources Management*, v. 21, n. 11, p. 1849–1863, 2007.

WORTHINGTON, A. C.; HOFFMAN, M. An empirical survey of residential water demand modelling. *Journal of Economic Surveys*, v. 22, n. 5, p. 842–871, Dec 2008.

# 3 SPATIAL WILLINGNESS TO PAY FOR A FIRST-ORDER STOCHASTIC REDUCTION ON THE RISK OF ROBBERY

## 3.1 INTRODUCTION

Contingent Valuation (CV) is a method widely used in recent decades. Its foremost objective is to infer, by means of public opinion surveys, the value of certain goods which are not readily tradable on traditional markets, such as public goods and natural resources. This method consists in constructing a hypothetical market for a certain good, as realistic and structured as possible, such that, by performing a survey, researchers can extract the maximum willingness to pay (WTP) of individuals for that good[1]. Bowen (1943) and Ciriacy-Wantrup (1947) were the pioneers to propose the use of public opinion surveys specially developed for the valuation of *social goods* or *collective goods* (CARSON; HANEMANN, 2005). These authors believed that voting would be the closest substitute to consumer choice, so they considered that the public opinion surveys would be a valid instrument for valuation of these goods (HOYOS; MARIEL, 2010; CARSON; HANEMANN, 2005).

Although the main goal of CV is to measure the monetary value of a certain good for an individual (CARSON; HANEMANN, 2005), there is a much more powerful insight on top of it: welfare analysis. According to Hoyos and Mariel (2010), by means of CV surveys, it is possible to directly obtain a monetary measure (Hicksian) of welfare associated with a discrete change in the provision of an environmental good, either by the substitution of one good for another or by the marginal substitution of different attributes of an existing good.

To understand the measurement of this value for the agent, we follow Whitehead and Blomquist (2006) and Carson and Hanemann (2005). Define a utility function that, for simplicity, only depends on a good $x$ and contingent good $q$, given by $u(x, q)$. Thus, assuming that good $q$ is desirable, and that $q^0$ is the state in which the consumer does not have the good and $q^1$ is the state in which the consumer has access to the good, the consumer will pay to consume the good if, and only if, the utility obtained with the consumption of the good is greater than the utility obtained without the consumption of the good, i.e., $u^1(x, q^1) > u^0(x, q^0)$.

---

[1] There is also the concept of minimum willingness to accept, where the individual reports the minimum amount he/she would be willing to accept to give up consuming a good that he/she would have been entitled. However, we will not cover this side.

So, the consumer will maximize their utility function $u(x, q)$, subject to their budget constraint, given by $y = px + tq$, where $y$ is the consumer's income, $p$ is the price of good $x$ and $t$ is the price of contingent good $q$, to define the optimal level of consumption of goods $x$ and $q$. From this, we find the indirect utility function, denoted by $v(p, q, y)$, whose usual properties with respect to $p$ and $y$ are satisfied. On the other hand, solving the problem of minimizing costs, subject to the constraint level of utility in state $q^0$, generates an expenditure function given by $e(p, q, u)$, (see Mas-Colell, Whinston and Green (1995)). According to Carson and Hanemann (2005), the value for the individual, in monetary terms, of the increment in utility caused by the change of state from $q^0$ to $q^1$ can be represented by two Hicksian measures: the compensatory variation and the equivalent variation. As shown by (MAS-COLELL; WHINSTON; GREEN, 1995). Formally, those measures are solutions to the following equations:

$$v^1(p, q^1, y - C) = v^0(p, q^0, y) \tag{3.1}$$

$$v^1(p, q^1, y) = v^0(p, q^0, y + E) \tag{3.2}$$

Based on these two concepts, one can define the willingness to pay in two different ways: i) as the difference between expenditure functions in the situation without contingent good and with contingent good, and, ii) as the monetary value that leaves the consumer indifferent between the *status quo* and the increase in the provision of contingent good. Following Carson and Hanemann (2005), it is possible to define the willingness to pay's function as a function to initial value $q^0$, the terminal value, $q^1$, and the values of $p$ and $y$ in which the changes in $q$ occur.

However, a common assumption for both $C(q^0, q^1, p, y)$ or $E(q^0, q^1, p, y)$ is the fact that what is measured is a discrete change between two deterministic states of nature with degenerate distribution, i.e., from initial value $q^0$ (*status quo*) with $Prob(q^0) = 1$ up to the terminal value, $q^1$ with $Prob(q^1) = 1$. The more general and interesting case of measuring willingness to pay for changes between (non-degenerate) lotteries of states of nature are still lacking a complete approach in the literature, although Cameron, DeShazo and Stiffler (2010) and Cameron and DeShazo (2013) are notably exceptions.

Although the scope of applicability of the CV method has grown considerably, many key areas traditionally approached by economists have not been thoroughly touched upon by contingent valuation. A notable example is the economics of crime. Since problems of measurement, externalities, and difficulties in assessing costs plague the area of crime and economics, it appears to us that underutilization of CV methods is hard to understand. In fact, very few papers have applied that method so far.

Ludwig and Cook (2001) estimate the benefits of reducing crime using CV methods. They focus on gun violence, in a national survey in the U.S. Using a parametric form, they found a value of US$ 24.5 billion as the worth for American society for a 30% reduction

in gun violence or US$1.2 million per injury avoided. Still in the U.S., Cohen et al. (2004) using a nationally representative sample of 1,300 U.S. residents, found that the representative American household would be willing to pay between US$ 100 and US$ 150 per year for programs that reduced specific crimes by 10% in their communities. Cohen et al. (2004) analyzed five types of crimes: burglary, serious assault, armed robbery, rape or sexual assault and murder.

In the U.K., Atkinson, Healey and Mourato (2005) valued the costs of three violent crimes: common assault (no injury), other wounding (moderate injury) and serious wounding (serious injury). Their data set contained 807 observations in Wales and in England. At the interview, respondents were told that the probability of being victims of each crime was 4% for common assault and 1% for both other wounding and serious wounding. Then each respondent was asked to express his WTP to reduce their chance of being victims of this offense by 50% over the next 12 months. The estimated values for WTP were £ 105.63, £ 154.54 and £ 178.33 for common assault, other wounding and serious wounding, respectively.

Finally, in Portugal, Soeiro and Teixeira (2010) studied the determinants of higher education students' willingness to pay for reducing the risk of being victims of violent crimes. They conducted an online survey with students from the University of Porto, which had 1,122 respondents. By means of a parametric approach, they modeled WTP as a function of demographic factors (age and gender), family-related factors (income, dimension, dependents), degree (undergraduate, master, PhD) and field of study (economics, arts, ...), crime-related factors (crime victim, crime time, physical injuries, psychological damages, fear of crime), averting behavior (locking doors), payment vehicle and policy. They found that variables such as age and family members had a negative impact in WTP, whereas variables such as gender, fear of crime, locking doors and payment vehicle had a positive impact on willingness to pay.

In Brazil, Araújo and Ramos (2009) used contingent valuation to estimate the loss of welfare associated with insecurity, by means of willingness to pay. The survey was conducted in the city of *João Pessoa (PB)*, and had 400 observations. Respondents were asked how much they would be willing to pay for a bundle of public security services, which includes: fixed police posts equipped with adequate weaponry; vehicles equipped for better care and effective police action; trained officers, with greater integration with the community and greater agility (speed) in citizen service; day and night patrols and conduction of educational programs to prevent violence and crime. They found that public security is a normal and common good and also that the estimated cost of insecurity in *João Pessoa* varies between R$ 6,524,727.01, considering the most conservative estimative, and R$ 104,864,863.52 for the highest value.

Although Ludwig and Cook (2001), Cohen et al. (2004), Atkinson, Healey and Mourato (2005), and Soeiro and Teixeira (2010) propose valuations between non-degenerate lotteries,

they stopped very far from building an econometric model that incorporates the basic tenets of choice under risk.

Given that state of affairs, say, the lack of a conceptual empirical strategy for CV among lotteries, and incipient literature on willingness to pay for crime reduction policies, our main contributions are: i) to build (and estimate) an econometric model capable of assessing the willingness to pay for first-order stochastic reductions in the risk of robbery, ii) to incorporate, in a sensible and manageable way, spatial effects to realistic mimics interactions present in a large and densely populated urban center in Brazil, and, iii) to apply our empirical strategy to real data, more specifically, to Brazilian data.

We believe to have succeeded in a satisfactory way. We make use of a unique geo-referenced sample of 4,030 households from the city of Fortaleza, CE (Brazil), containing information on socioeconomic background, experience, expectation of victimization, and willingness to pay to reduce some type of crimes (see, Carvalho (2012)). For the global model (i.e., without spatial effects), the parameters for all independent variables, except *Age*, show positive signs. Older people tend to pay less to reduce crime than young people do, men and more educated people tends to pay more for risk reductions. Finally, as to variables of perception and experience of victimization (variable *Perception of patrolling* is measured in decreasing order), the lower the perception of patrolling, the greater the willingness to pay to reduce the number of robberies, and people who were *Victims of robbery* tend to pay more to prevent such experience again. We also estimated an average willingness to pay of R$ 23.35 per month/household, a value of R$ 5.91 higher than the estimated value of the nonparametric form. Also, we estimated the implicit value of a statistical robbery approximately equal to R$ 11,969 per crime avoided. Both values are quite reasonable. As a matter of fact, our proposed specification made possible to implicitly estimate the average cost of each robbery in the city of Fortaleza. This amounts to approximately 4,15% of the income. Multiplying this value by average income, we have a value of R$ 61,38 per robbery.

The full spatial heterogeneity reveals our local model. By means of a geographically weighted regression (GWR), it is possible to allow for the estimation of local parameters rather than global parameters (FOTHERINGHAM; BRUNSDON; CHARLTON, 2003). Now, the main difference is the lack of one estimated parameter for each independent variable. Instead, for each independent variable, we have a possible different parameter for each sampled point. Overall, the estimated spatial heterogeneity brings us both expected results and surprises. The estimative mapping for variables *Gender*, *Age* and *Education* present a reasonable amount of spatial heterogeneity and, as expected, follow the very inertial city's socioeconomic spatial distribution profile. Given the geographically weighted regression, we implement a protocol to calculate a surface of willingness to pay. In order to do that, we apply Kriging techniques. The image that emerges from such empirical exercise is not difficult to rationalize: the income, age, and crime spatial distribution of

Fortaleza has an important effect on the surface of willigness to pay. Although peripheries present lower willingness to pay, as long as we go inwards, there is plenty of heterogeneity on the spatial distribution of willingness to pay for robbery reduction. It is worth noting that the highest willingness to pay is not necessarily the richest one, corroborating with a theory of crime that posits an active role for victim (costly) precautions.

Besides this introduction, we have more 5 Sections. Section (3.2) introduces the data set used in our estimatives. Section (3.3) develops a simple structural model of contingent choice between risky lotteries and frames the resulting equation as a fully parametric econometric model, although, given the type of data collected, we end up estimating it by fully maximum likelihood. In order to introduce our spatial effects, Section (3.4) deals with geographically weighted regression and how to manage that in our context. We call such model local to contrast with the previous one that neglects spatial effects. All estimatives are performed on Section (3.5), as well as their interpretations. Finally, Section (3.6) elaborates more on results and proposes futures improvements.

## 3.2   DATA SET

Our data set comes from a survey conducted in 2012 where a total of 4,030 households were sampled along 119 districts (*bairros*) from the city of Fortaleza (Brazil) during the months of October 2011 to January 2012, see Carvalho (2012). Besides information about socioeconomic background, experience and expectation of victimization, Carvalho (2012) induced respondents to express their willingness to pay to reduce certain types of crimes. A key component from the data set is due to the fact that household, work and school positions were georeferenced.

The section about contingent valuation presents respondents with a fictional scenario where there was a program to fight against criminality, more specifically the crime of robbery. The respondent was informed that the program was successful and succeeded in reducing 50% the amount of robberies. However, to maintain this program, it was necessary that the population funds it by means of fictitious future taxes. Then, respondents were asked if they were willing to pay a monthly fee to maintain that crime prevention program, and if so, how much they would be willing to pay monthly. The exact introduction and question wording were:

●Introductory Remark: *Now I would like to know how much you are willing to spend to reduce certain crimes in your town. In each case, I will ask you to answer whether you would vote ''yes'' or ''no'' for a bill that would require from you and from each household in your community a payment to prevent certain crimes. Remember that the money you agree to spend to prevent crimes is the same that you could use to buy food, clothes or other needs to you and to your family.*

•Question: *Q105 Now forget about this program that was able to reduce homicides and think about a new one. Let's suppose a new government program funded by the population of Fortaleza managed to cut the occurrence of personal robbery in the city in half. Would you be willing to pay a monthly amount to keep this program of crime prevention?*

Table 3.1 defines the variables used in this paper. Initially, we show the socioeconomic profile, the perception of security and experience of victimization of the research participants. From a total of 4,030 observations in the initial sample, 246 observations were removed due to lack of information about participation in the program and due to the difficulty in georeferencing respondents' addresses. Table 3.2 shows that 44,66% of the respondents were men. The overall age of respondents was 39,45 years old, with complete fundamental school level. As to income, its average level was R$ 1,488.70 per month, but about 50% of respondents earn R$ 817.50 or less[2].

Table 3.1 – Variables' descriptions

| Variable | Description |
| --- | --- |
| Gender | 1 if male; 0 if female |
| Age | years |
| Income | R$ |
| Education | 1.No education; 2.Incomplete fundamental school; 3.Complete fundamental school; 4.Incomplete high school; 5.Complete high school; 6.Incomplete undergraduate degree; 7.Complete undergraduate degree; 8.Graduate Program |
| Victim of robbery | 1 if you've been the victim of robbery; 0 Otherwise |
| Subject Prob. | $\in (0,1)$ |
| Perception of patrolling | 1.Always; 2.Often; 3.Sometimes ; 4.Rarely; 5.Never |
| Willingness to pay | R$/month |

Source: Elaborated by the Authors

As to victimization and perception of security, 23,24% of respondents were victims of robbery at least once in the last five years. As to perception of security, on average, respondents considered that the probability of being robbed in the next 12 months was about 49.2%, although the perception of patrolling is frequent (mean 2.05).

Table 3.4 shows that out of 3,784 respondents, 1,709 (45,16%) answered that they are willing to pay a monthly fee to fund the program to combat robberies, while 2,076 (54,86%) answered they would not pay any amount. Despite the fact that the number of people who are not willing to pay to keep the program to combat crime is fairly high, it is consistent with other studies about contingent valuation, for instance, Atkinson, Healey and Mourato (2005), which had 34,57%, and Araújo and Ramos (2009), which had a rate of 48.5 %, the last one for the city of João Pessoa (Brazil). This second group is defined in the contingent valuation literature as *protesters*. These people refuse to pay for a good

---

[2]   The minimum wage in Brazil at the time of the survey was R$ 545,00.

Table 3.2 – Sample Description - Total

| Variable | Mean | Std. dev | Min | Median | Max | NA | N |
|---|---|---|---|---|---|---|---|
| Gender | 0.4466 | 0.4972 | 0 | 0 | 1 | 0 | 3784 |
| Age | 39.4519 | 16.8278 | 16 | 37 | 94 | 80 | 3704 |
| Income | 1,488.70 | 1,524.27 | 272.50 | 817.50 | 10,900.00 | 137 | 3647 |
| Education | 3.5552 | 1.6451 | 1 | 3 | 8 | 0 | 3784 |
| Victim of robbery | 0.2324 | 0.4224 | 0 | 0 | 1 | 1 | 3783 |
| Subject Prob. | 0.4920 | 0.2990 | 0 | 0.5000 | 1 | 532 | 3252 |
| Perception of patrolling | 2.0533 | 1.2146 | 1 | 1 | 5 | 10 | 3774 |

Source: Elaborated by the Authors

either because they think they already pay many taxes or, in the case of public goods, because the provision of such goods is responsibility of the government, or simply because it is the duty of other groups to pay for the provision of that good[3]. However, it is possible that someone reports a true zero value for the reduction on the risk of being robbed or just cannot afford to pay for such amount.

Notwithstanding that, we will not enter in this debate[4], and we simply characterize them as *protesters*. The *protesters'* group, 50.48% of them were men, with an average age of 41.42 years old and with complete fundamental school level. In this group, the average income was equal to R$ 1,488.90, but 50% of them earned R$ 817.50 or less. Concerning the expectation of victimization and perception of security, 22,74% of them suffered at least one robbery in the last five years and they consider that the probability of being robbed in the next 12 month is about 48,47%, even though the perception of patrolling is frequent. In the CV's literature, the standard procedure for dealing with this group is to remove them from the sample and proceed to estimation of the maximum willingness to pay (STRAZZERA et al., 2003). However, Strazzera et al. (2003) states that this procedure is valid only when both groups are similar due to the fact that, if this is not the case, selection bias will pop up.

We compared the empirical distributions for both protesters and those who are willing to pay a positive amount of money and they are quite similar. Thus, *protesters* and those who are willing to pay are quite homogeneous, which indicates that the estimates of willingness to pay using only the second group should not be affected by selection bias (STRAZZERA et al., 2003). Thus, we remove the group of *protesters* from the sample in order to estimate the cost of robberies. Table 3.4 shows the characteristics of those who are willing to pay to maintain the crime's reduction program.

Table 3.5 shows the frequency distribution of willingness to pay. From the total of 1,708 respondents who answered that they would accept to pay some amount to reduce robberies, 70 did not know/ did not want to inform a value from those presented in the

---

[3] We also consider the fact that individuals do not report their willingness to pay for fear that, once answered a value, the research can be used to make them pay the reported amount.

[4] For details, see Jorgensen et al. (1999)

Table 3.3 – Sample description - Protesters

| Variable | Mean | Std. dev | Min | Median | Max | NA | N |
|---|---|---|---|---|---|---|---|
| Gender | 0.5048 | 0.5001 | 0 | 1 | 1 | 0 | 2076 |
| Age | 41.4205 | 17.4722 | 16 | 39 | 93 | 26 | 2050 |
| Income | 1,488.90 | 1,553.82 | 272.50 | 817.50 | 10,900.00 | 84 | 1992 |
| Education | 3.5111 | 1.6823 | 1 | 3 | 8 | 0 | 2076 |
| Victim of robbery | 0.2274 | 0.4192 | 0 | 0 | 1 | 0 | 2076 |
| Subject Prob. | 0.4847 | 0.2973 | 0 | 0.5000 | 1 | 302 | 1774 |
| Perception of patrolling | 2.0778 | 1.2273 | 1 | 2 | 5 | 6 | 2070 |

Fonte: Elaborated by the Authors

Table 3.4 – Sample description - Willing to Pay

| Variable | Mean | Std. dev | Min | Median | Max | NA | N |
|---|---|---|---|---|---|---|---|
| Gender | 0.3759 | 0.4845 | 0 | 0 | 1 | 0 | 1708 |
| Age | 37.0121 | 15.6586 | 16 | 34 | 94 | 54 | 1654 |
| Income | 1,488.45 | 1,488.40 | 272.50 | 817.50 | 10,900.00 | 53 | 1655 |
| Education | 3.6089 | 1.5976 | 1 | 4 | 8 | 0 | 1708 |
| Victim of robbery | 0.2384 | 0.4262 | 0 | 0 | 1 | 1 | 1707 |
| Subject Prob. | 0.5007 | 0.3008 | 0.0100 | 0.5000 | 1 | 230 | 1478 |
| Perception of patrolling | 2.0235 | 1.1987 | 1 | 1 | 5 | 4 | 1704 |

Source: Elaborated by the Authors

payment card. Thus, we had 1,638 observations. The columns *Cumulative frequency* and *Survival probability* in table 3.5 indicate, respectively, the number of people and the percentage of the sample that is willing to pay at least the indicated value. Thus, it can be seen that 990 people, equivalent to 60.40%, are willing to pay at least R\$ 10 for the maintenance of the combating robbery crimes program.

Table 3.5 – Willingness to pay frequency distribution

| WTP | Frequency | Cumulative frequency | Survival probability |
|---|---|---|---|
| 1 | 212 | 1638 | 1.0000 |
| 5 | 428 | 1426 | 0.8706 |
| 10 | 460 | 998 | 0.6093 |
| 15 | 159 | 538 | 0.3284 |
| 25 | 173 | 379 | 0.2314 |
| 50 | 123 | 206 | 0.1258 |
| 75 | 12 | 83 | 0.0507 |
| 100 | 38 | 71 | 0.0433 |
| 150 | 8 | 33 | 0.0201 |
| +150 | 25 | 25 | 0.0153 |

Source: Elaborated by the Authors

One can also notice that, as the value of willingness to pay increases, fewer people will be willing to pay this amount.

From this empirical distribution of willingness to pay, we estimated, nonparametrically,
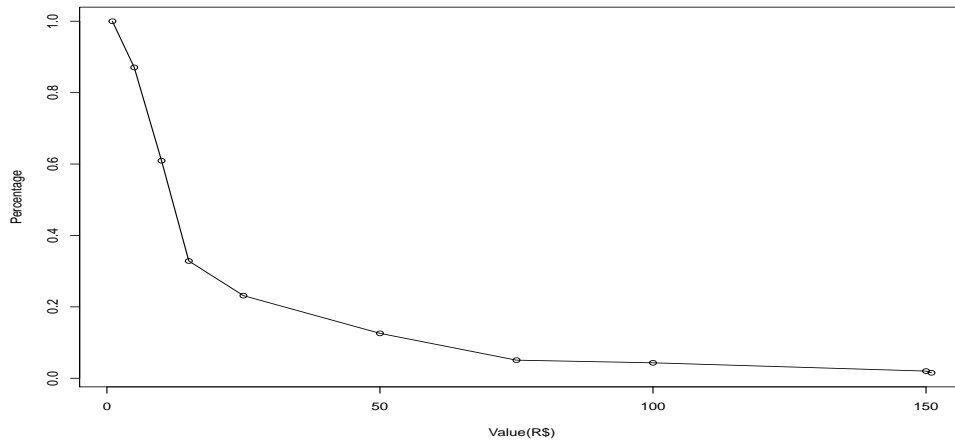
Figure 3.1 – Survival Function

maximum willingness to pay[5]. We estimated the value of R$ 17.44 as the average monthly value or R$ 173.28 per year, as the value that each household would be willing to contribute to reduce the number of robberies in the city of Fortaleza by 50%. Thus, multiplying this value by the total amount of households in Fortaleza, that according to Estatística (2012) is 709,952 households, we estimated a total value of willingness to pay of approximately R$ 123.02 million per year. Finally, considering that the number of robberies in Fortaleza in 2011 was equal to 33.240[6], we have that the implicit value of a statistical robbery[7] is equal to R$ 7,402.03, that is the cost of a robbery to society. However, this nonparametric estimation is not the ideal procedure to estimate the maximum willingness to pay, once it is expected that individual characteristics influence the amount that individuals are willing to pay. So, in the next section, a parametric model to estimate the maximum willingness to pay for the maintenance the program to reduce robberies will be presented.

## 3.3 ECONOMETRIC MODEL

Our objective is to build a contingent valuation model to assess willingness to pay for a first-order stochastic improvement on the odds of being robbed in the city of Fortaleza, Brazil, when subjective expectations about the risk are available. Since our data sets come from the same urban space, spatial effects should also be considered. The random vector $(R, M, X)$, where $R \in \{0, 1\}$, is a binary indicator if a shock did not occur or did occur, $M \in \mathrm{R}_+$ measures shock's monetary cost (tangible and intangible costs), $X \in \mathrm{R}^K$ a vector

---

[5]   We consider only who answered that would be willing to pay more than R$ 1.00 and equal or less than R$ 100.00 per month

[6]   Considering only robberies informed to the public security authorities. Source: SSPDC-CE

[7]   To obtain this value, just divide R$ 123.02 million by 16,620, the last one being the number of robberies avoided.

of individual and/or state-specific characteristics. $\theta \in \{0, 1\}$ is an indicator of *status-quo* situation or alternative status to be achieved with transfers.

We define four objective distribution functions. $P_{R,M,X}(r, m, x) \equiv Prob(R \leq r, M \leq x, X \leq x)$, the distribution function of $(R, M, X)$. Accordingly, the conditional distributions $P_{R,M|X}(r, m|X = x) \equiv Prob(M \leq m, R = r|X = x)$, $P_{M|R,X}(m|R = r, X = x) \equiv Prob(M \leq m|R = r, X = x)$, and $P_{R|X}(r|X = x) \equiv Prob(R \leq r|X = x)$ are defined. Index individuals by $i \in \{1, 2, \cdots, n\}$. We also define four subjective distribution functions, say, $P^i_{R,M,X}(r, m, x)$, $P^i_{R,M|X}(r, m|X = x)$, $P^i_{M|R,X}(m|R = r, X = x)$, $P^i_{R|X}(r|X = x)$

**Hypothesis 1.** *The values for $P_{R,M,X}(r, m, x)$, $P_{R,M|X}(r, m|X = x)$, $P_{M|R,X}(m|R = r, X = x)$, $P_{R|X}(r|X = x)$ exist and are well-defined for any $\theta \in \{0, 1\}$ and $i \in \{1, 2, \cdots, n\}$.*

**Hypothesis 2.** *$P^i_{R,M,X}(r, m, x)$, $P^i_{R,M|X}(r, m|X = x)$, $P^i_{M|R,X}(m|R = r, X = x)$, $P^i_{R|X}(r|X = x)$ exist and are well-defined for any $\theta \in \{0, 1\}$ and $i \in \{1, 2, \cdots, n\}$.*

**Hypothesis 3.** *Except for $P^i_{R|X}(r|X = x)$, the distribution functions are homogenous across individuals and equal to its respective objective distribution.*

With a slight abuse of notation, our basic random set up is described by the following vector $\left( P^\theta_{R,M,X}, P^\theta_{R,M|X}, P^\theta_{M|R,X}, P^{\theta,i}_{R|X} \right)$, for all $\theta \in \{0, 1\}$ and $i \in \{1, 2, \cdots, n\}$.

For each $\theta \in \{0, 1\}$, any individual $i \in \{1, 2, \cdots, n\}$ is endowed with an indirect utility function given by $V_{i,\theta} = V(y_i, \theta)$. Where $y_i$ is a sure amount of money and $\theta \in \{0, 1\}$ is an indicator of *status-quo* situation or alternative status to be achieved with transfers. Two quantities of interests are:

$$\mathrm{E}\left( V_{i,0} \right) = V(y_i - m, 0) Pr^0_{M|R,X} \times Pr^{0,i}_{R|X} + V(y_i, 0)(1 - Pr^{0,i}_{R|X}) \tag{3.3}$$

$$\mathrm{E}\left( V_{i,1} \right) = V(y_i - s_i - m, 1) Pr^1_{M|R,X} \times Pr^{1,i}_{R|X} + V(y_i - s_i, 0)(1 - Pr^{1,i}_{R|X}) \tag{3.4}$$

For pragmatic reasons, we assume that each individual is risk neutral and assume a linear functional form for his/her indirect utility function.

**Hypothesis 4.** *The indirect utility function for each $\theta \in \{0, 1\}$, any individual $i \in \{1, 2, \cdots, n\}$ is parametrized as $V(\widetilde{y}_i, \theta) = \beta \widetilde{y}_i + \alpha_\theta X_i + \epsilon_{i,\theta}$.*

We also assume that the distribution of shock's size is independent from the occurrence of the shock, say $R$, and observed heterogeneity, $X$. Also, for simplicity, the expected value of shock's size depends only linearly on individual income.

**Hypothesis 5.** *$Pr^0_{M|R,X} = Pr^1_{M|R,X} = P(M)$, and $\overline{M} = \tau_1 + \tau_2 Y$.*

From Equations (3.3) and (3.4), we have:

$$\beta \left(y_i - mP(M)\right) Pr_{R|X}^{0,i} + \beta y_i \left(1 - Pr_{R|X}^{0,i}\right) + \alpha_0 X_i + \epsilon_{i,0} \tag{3.5}$$

$$\beta \left((y_i - s_i - mP(M)) Pr_{R|X}^{1,i} + \beta(y_i - s_i) \left(1 - Pr_{R|X}^{1,i}\right) + \alpha_1 X_i + \epsilon_{i,1}\right) \tag{3.6}$$

Note, however, that the change in *status-quo* is a change in $P_{R|X}^{i,\theta}(r|X = x)$. In fact, it is easy to see that $P_{R|X}^{i,1}(r|X = x) \geq_{FSD} P_{R|X}^{i,0}(r|X = x)$, where $\geq_{FSD}$ means first-order stochastic dominance[8]. Note that $Pr_{R|X}^{1,i} = k \times Pr_{R|X}^{0,i}$, where $k \in (0,1)$, the payoffs are $\{-\overline{m}, 0\}$, and $Pr_{R|X}^{\theta,i} = Prob(R = -\overline{m}|X)$, for $\theta \in \{0, 1\}$. See Figure (3.2). Now we are able to develop the expression for the willingness to pay by equating Equations (3.5) to (3.6), and solving for $s_i$.



Figure 3.2 – First-Order Stochastic Dominance

$$s_i = \left(Pr_{R|X}^{0,i} - Pr_{R|X}^{1,i}\right) \overline{m} + \frac{(\alpha_1 - \alpha_0)}{\beta} X_i + \frac{(\epsilon_{i,1} - \epsilon_{i,0})}{\beta} \tag{3.7}$$

First, note that the expression for the willingness to pay $s_i$ depends on the difference in the expected value of the shock between the *status quo* and the new situation, say, $\left(Pr_{R|X}^{0,i} - Pr_{R|X}^{1,i}\right) \overline{m}$, as well as it depends on observed and unobserved heterogeneity. In fact, as expected, given a risk neutral agent, the income does not have a bite. However, in the sequel, we show that the shock value is a function of the cross-product of income and subjective probability of robbery, say, $Pr_{R|X}^{0,i} y_i$! Hence, the final expression for $s_i$ is dependent on $y_i$(individual income).

---

[8]    Remember that the counterfactual proposed by question 105 in Carvalho (2012) was phrased like *"to cut the occurrence of personal robbery in Fortaleza in half"*

Remember that $Pr_{R|X}^{1,i} = k \times Pr_{R|X}^{0,i}$, and $\overline{m} = \tau_1 + \tau_2 y_i$. Defining $\alpha \equiv \frac{(\alpha_1 - \alpha_0)}{\beta}$, and $\epsilon_i \equiv \frac{(\epsilon_{i,1} - \epsilon_{i,0})}{\beta}$, and $Z_i \equiv (1 - k)Pr_{R|X}^{0,i}$, $W_i \equiv Z_i y_i$ we get the estimable equation, where, with no loss of generality, right-hand side variable appears in logarithmic form:

$$\ln(s_i) = \ln\left(\tau_1 Z_i + \tau_2 W_i + \alpha X_i\right) + \epsilon_i \tag{3.8}$$

Where, $\epsilon_i \sim N(0, \sigma^2)$. Approximating the left-hand side of Equation (3.8) by a first order Taylor's expansion (Note that the full Taylor's approximation is $\ln(x) = (x - 1) - \frac{(x-1)^2}{2} + \frac{(x-1)^3}{3} - \frac{(x-1)^4}{4} \cdots$, as long as we re-scale monetary values $s_i$ to belong to the interval $(-1, 1]$ actually close to 1, we get:

$$\ln\left(\frac{s_i}{\theta}\right) = \ln\left(\frac{\tau_1}{\theta}Z_i + \frac{\tau_2}{\theta}W_i + \frac{\alpha}{\theta}X_i\right) + \epsilon_i \tag{3.9}$$

Where $\theta \in (\min(s_i), \max(s_i))$. Assuming that $X_i$ has an intercept whose parameter is $\frac{\alpha_{intercpt}}{\beta}$:

$$\ln(s_i) = \left(\frac{\alpha_{intercpt}}{\beta} + \ln(\theta) - 1\right) + \frac{\tau_1}{\theta}Z_i + \frac{\tau_2}{\theta}W_i + \frac{\alpha}{\theta}X_i + \epsilon_i \tag{3.10}$$

Before we proceed, it is worth noting that models which incorporate willingness to pay for first-order stochastic dominance improvements on risks is a quite new endeavor. In fact, there are only two papers we are aware of, say, Cameron, DeShazo and Stiffler (2010) and Cameron and DeShazo (2013), that build on this topic. Their approach is different from ours, however. Last, but not least, it is important to stress that individual's subjective expectations play a crucial role in our modeling strategy.

So, Cameron and Huppert (1989) propose that contingent valuation data sets obtained by means of payment cards' method can be analyzed parametrically by means of maximum likelihood models with data in intervals. They suggest that when an agent chooses a value in payment card, say $t_{ui}$, the true value of the agent's willingness to pay is greater than or equal to this value, but less than the next card value, say $t_{u+1i}$. Therefore, the probability that the agent chooses to pay the $t_{ui}$ value is equal to the probability that the true willingness to pay is in the range defined by $t_{ui}$ and $t_{u+1i}$.

$$P(t_{ui}) = P(t_{ui} \leq s < t_{u+1i}) \tag{3.11}$$

Thus, it is possible to rewrite 3.11 as:

$$P(t_{ui}) = P(log(t_{ui}) \leq log(s_i) < log(t_{u+1i})) \tag{3.12}$$

By equation 3.10, we have that $s$ has mean $\mu$ and standard deviation $\sigma$. Then, define $\mu$ as:

$$\mu = \left(\frac{\alpha_{intercpt}}{\beta} + \ln(\theta) - 1\right) + \frac{\tau_1}{\theta}Z_i + \frac{\tau_2}{\theta}W_i + \frac{\alpha}{\theta}X_i \tag{3.13}$$

we can standardize each pair of interval thresholds and state that:

$$P(t_{ui}) = P\left(\frac{log(t_{ui}) - \mu}{\sigma} \le z_i < \frac{log(t_{u+1i}) - \mu}{\sigma}\right) \qquad (3.14)$$

where $z_i$ is the standard normal random variable. The probability above can be rewritten as the difference between two standard normal cumulative densities. Then, let $\Phi$ denote the cumulative density function of a standard normal variable, it follows that:

$$P(t_{ui}) = \Phi\left(\frac{log(t_{u+1i}) - \mu}{\sigma}\right) - \Phi\left(\frac{log(t_{ui}) - \mu}{\sigma}\right) \qquad (3.15)$$

Finally, Cameron and Huppert (1989) assert that the joint probability density function for $n$ independent observation can be interpreted as a likelihood function, defined over the unknown parameters $\gamma$ e $\sigma$. Thus, the log-likelihood function takes the following form:

$$logL = \sum_{i=1}^{n} log\left[\Phi\left(\frac{log(t_{u+1i}) - \mu}{\sigma}\right) - \Phi\left(\frac{log(t_{ui}) - \mu}{\sigma}\right)\right] \qquad (3.16)$$

From the maximization of (3.16), we find the optimal values of $\gamma$ e $\sigma$, with values of $\gamma$ showing the impact of individual characteristics on the choice of the value of willingness to pay. From these estimated values of $\gamma$ and $\sigma$, it is possible to estimate mean and median WTP, as shown below:

$$\text{Mediana DAP} = exp\left(\left(\frac{\alpha_{intercpt}}{\beta} + \ln(\theta) - 1\right) + \frac{\tau_1}{\theta}Z_i + \frac{\tau_2}{\theta}W_i + \frac{\alpha}{\theta}X_i\right) \qquad (3.17)$$

$$\text{DAP Média} = exp\left(\left(\frac{\alpha_{intercpt}}{\beta} + \ln(\theta) - 1\right) + \frac{\tau_1}{\theta}Z_i + \frac{\tau_2}{\theta}W_i + \frac{\alpha}{\theta}X_i\right)exp(\sigma/2) \qquad (3.18)$$

These two measures provide what we call a global value for WTP. However, we expect spatial heterogeneity to have an important role in the relation between the choice of how much the agent wants to pay and his characteristics. This means that values of $\gamma$ can be different, which would make individual WTP values differ all over the city. A plausible explanation for this would be that individuals in different neighborhoods meet different levels of criminality, whether observed or not by police authorities[9], which would lead their willingness to pay to be different. Also, the spread of information about crimes throughout the urban fabric is not understood so far.

Thus, in order to handle this issue of spatial heterogeneity, we use the geographically weighted regression technique (GWR) to estimate a local WTP in such a way that it will be possible to identify in which regions the WTP will assume higher values. Next section presents the GWR model.

---

[9]  The security agencies only have access to the criminality level in an area from the time the citizen registers the event of a crime on, which does not always happen.

## 3.4   A "LOCAL" ECONOMETRIC MODEL

According to Almeida (2012), analyzing only the average or global response of a phenomenon may not be useful or convenient, since socioeconomic phenomena are not likely to be constant in different regions. Fotheringham, Brunsdon and Charlton (2003) refer to this situation as spatial non-stationarity and claim that any relationship that is non-stationary over space is not well represented by a global statistic and, indeed, this global value may be very misleading locally.

Fotheringham, Brunsdon and Charlton (2003) affirm that there are several reasons to expect that a relationship varies over space. Among possible explanations, we can cite sample variations, misspecification and, most importantly, there might be relationships which are intrinsically different across space. In the last case, it is suggested that there are spatial variations in peoples' attitudes or preferences or there are different administrative, political or other contextual issues that produce different responses to the same stimuli over space.

When it comes to the object of study in this article, it is a stylized fact that crime distribution is heterogenous across space. In big cities, like Fortaleza (the fifth largest city in Brazil with an area of 313 square kilometers, boasting one of the highest demographic densities in the country, say, 8,001 per $km^2$), robberies are concentrated on richer areas in the city, leading to formation of crime clusters in these areas. Due to this heterogeneous distribution, we expected that individuals' reactions to crime would also be heterogeneous. So, we expect that an individual who lives in a region with high rates of criminality has a different behavior than an individual who lives in low crime prone regions. Thus, unlike classical models of spatial dependence, here we do not expect that individuals can influence each other's willingness to pay, but we expected that different individuals have different factors that influence their willingness to pay. So, a variable that can influence the willingness to pay for individual $i$ maybe have no influence on individual $j$, or have more or less influence. In this sense, a local model is necessary to estimate this relationship.

The geographically weighted regression (GWR) is a method that extends the traditional regression framework by allowing the estimation of local parameters rather than global parameters (FOTHERINGHAM; BRUNSDON; CHARLTON, 2003). This method generates a sequence of regressions estimated for each region, using subsamples from the data, weighted by distance (ALMEIDA, 2012). Subsamples are created from the regression or calibration point, that is the reference point for the parameters estimation for region $i$. From this point, each observation belonging to the sample is weighted according to its distance to the calibration point. Close observations have a higher weight, while more distant observations have a lower weight (ALMEIDA, 2012).

The weights used for the creation of these subsamples are taken by the spatial kernel function. According to Almeida (2012), the kernel function is a real, continuous and symmetric function in which integral sums one, like a probability density function. This

function uses the distance ($d_{ij}$) between two points and a parameter of bandwidth ($b$) to determine a weight between these two regions, which is inversely related to geographic distance ($w_{ij}$).

Fotheringham, Brunsdon and Charlton (2003) classify the spatial kernel functions in two groups: the fixed kernels and the adaptive kernels. In the fixed kernels, bandwidth ($b$) is fixed, which may lead to problems of bias and efficiency. With a fixed bandwidth, the number of observations in each subsample may vary substantially. In regions where data are dense, the kernels are larger than they need to be and hence using information in excess, turning estimates biased. On the other hand, in regions where data is scarce, the kernels are smaller than they need to be to estimate the parameters' reliably. The adaptive kernels reduce both problems by making bandwidth ($b$) greater or smaller depending on data density in the area. Thus, we use the adaptive gaussian[10] kernel, defined by equation (3.19):

$$w_{ij} = \left\{ exp\left( -\frac{1}{2}\left(\frac{d_{ij}}{b_i}\right)^2 \right), \text{if } d_{ij} < b_i \right. \tag{3.19}$$

In the fixed case, only one bandwidth ($b$) is chosen for every data point, whereas in the adaptive case, one bandwidth ($b_i$) is chosen for each data point, such that each subsample has the same proportion of the data.

Due to the aforementioned problems, the choice of the bandwidth ($b$) must be made in order to try to solve the trade-off between bias and efficiency. To this end, to avoid arbitrary choices, the bandwidth is estimated using the data (ALMEIDA, 2012). There are several techniques[11] used to determine the optimal value of the bandwidth. In this paper, we use the cross-validation technique. It consists in minimizing the following function, represented by equation 3.20:

$$CV = \sum_{i=1}^{n} \left( y_i - \hat{y}_{\neq i}(b) \right)^2 \tag{3.20}$$

where $y_i$ is the dependent variable , $n$ is the number of observations, $b$ is the bandwidth and $\hat{y}_{\neq i}(b)$ is the fitted value of $y_i$ using a bandwidth of $b$ with the observations for point $i$ omitted from the calibration process (ALMEIDA, 2012). Fotheringham, Brunsdon and Charlton (2003) affirm that this approach has the desirable property of countering the wrap-around effect, since when $b$ becomes very small, the model is calibrated only on samples near to $i$ and not at $i$ itself.

---

[10]  For other types of kernel functions, see, among others, Fotheringham, Brunsdon and Charlton (2003) and Almeida (2012).

[11]  For more details, see Fotheringham, Brunsdon and Charlton (2003).

After obtaining these weights generated by the kernel function, it is possible to get the local spatial weighting diagonal matrix:

$$W(u_i, v_i) = \begin{bmatrix} w_{i1} & 0 & \cdots & 0 \\ 0 & w_{i2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & w_{in} \end{bmatrix} \tag{3.21}$$

where $w_{in}$ is the weight attributed to point $n$ in the model calibration in regression point $i$, obtained by means of spatial kernel function. Thus, from the model showed in equation (3.16), the local model can be specified by the following weighted maximum likelihood function, represented by equation (3.22):

$$logL = \sum_{i=1}^{n} W(u_i, v_i) log \left[ \Phi \left( \frac{log(t_{u+1i}) - \mu}{\sigma} \right) - \Phi \left( \frac{log(t_{ui}) - \mu}{\sigma} \right) \right] \tag{3.22}$$

From the equation above, are estimated parameter sets for each $n$ points. Next section presents results of the estimation for both global and local models.

## 3.5 RESULTS

### 3.5.1 Results from the Global Model

Table 3.6 shows the results of the estimation of the global model represented by equation (3.16). All parameters are statistically significant. The sign of estimates indicate the effect on willingness to pay. All variables, except *Age*, show positive signs. The negative sign of variable *Age* indicates that older people tend to pay less to reduce crime than young people, indicating that older people have a greater feeling of security than young people. As to variables *Gender* and *Education*, the positive sign indicates that men and more educated people tend to pay more to reduce the risk of being robbed. Finally, as to variables of perception and experience of victimization (variable *Perception of patrolling* is measured in decreasing order), the lower the perception of patrolling, the greater the willingness to pay to reduce the number of robberies, and people who were *Victims of robbery* tend to pay more to prevent such experience again.

Variable *I(Subject Prob. * Income)* deserves special attention. The positive sign of this variable shows that the higher the income and the subjective probability of being robbed, the higher the willingness to pay to reduce the risk of being robbed. On the other hand, when this variable is multiplied by $\theta$, we have the fraction, in average, of the income robbed in each robbery. This value is equal to approximately 0.0415. So, in each robbery, approximately 4,15% of the income is robbed. Multiplying this value by the average income, we have a value of R\$ 61,38 per robbery.

In table 3.7, the values of estimated WTP, as defined by equation 3.17, are presented.

Table 3.6 – Estimates - Parametric Maximum Likelihood

| Variable | Estimate | Stand. Dev. | t | p-value |
|---|---|---|---|---|
| (Intercept) | 2.3949 | 0.1024 | 23.3750 | 0.0000 |
| I(Subject Prob. * Income) | 0.0002 | 0.0000 | 4.7072 | 0.0000 |
| Gender | 0.1897 | 0.0477 | 3.9738 | 0.0000 |
| Age | -0.0034 | 0.0016 | -2.0946 | 0.0362 |
| Education | 0.0466 | 0.0160 | 2.9058 | 0.0036 |
| Perception of patrolling | 0.0520 | 0.0196 | 2.6559 | 0.0079 |
| Victim of robbery | 0.0926 | 0.0550 | 1.6818 | 0.0926 |
| $\sigma$ | 0.7299 | 0.0167 | 43.6735 | 0.0000 |

log-likelihood = -1851.53

Newton-Raphson maximization, 4 interations

Source: Elaborated by the authors

Table 3.7 – Results of WTP(R$) from the global parametric model

| Variable | Estimate | Stand. Dev. | Inter. Conf. |
|---|---|---|---|
| Mean | 23.35 | 6.12 | 22.98 - 23.72 |
| Median | 16.21 | 4.25 | 15.95 - 16.47 |

Source: Elaborated by the authors

The average WTP estimated from the global model is equal to R$ 23.35 per month/household, a value R$ 5.91, which is higher than the estimated value of the nonparametric form, which was only R$ 17.44. Thus, if the government decided to implement a monthly tax about this value, it would be possible to raise, per year, R$ 280.20 per household, which would generate an average tax revenue of about R$ 198.92 million per year, equivalent to approximately 20.63% of the amount spent on public security in the state of Ceará in 2011[12]. Assuming a worst case scenario, using the median WTP value of R$ 16.21 per month/household as a benchmark, we have a value of R$ 194.52 per year/household. In this case, the annual tax revenue in Fortaleza would be approximately R$ 138.09 million, equivalent to 14.32% of spending on public security in 2011.

Now, considering the damage of robberies to society, in the first scenario, where the WTP was estimated in R$ 23.35, we got an implicit value of a statistical robbery of approximately R$ 11,969 per robbery avoided. Considering the second scenario, where we assumed the WTP median value equal to R$ 16.21, the value of a statistical robbery was estimated approximately equal to R$ 8,310 per crime avoided. Next section presents results from the local model.

---

[12]    According to Pública (2012), the amount spent on public security in the state of Ceará in the year of 2011 was R$ 964,095,556.61.

## 3.5.2  Results from the Local Model

As discussed earlier, considering only the average or global response of a phenomenon may not be useful or convenient. So, we estimated[13] the local model specified by Equation (3.22). First, we present the estimated model with an adaptive bandwidth. The cross-validation technique pointed us a bandwidth (*b*) of 0.6149766 with a CV score of 580.1328, indicating that each sub-sample has approximately 61,5% of the sample. Table 3.8 shows the estimates under this value of *b*.

Table 3.8 – Estimates for the Local Model - GWR -Adaptive bandwidth

|  | Min. | 1st Qu. | Median | Mean | 3rd Qu. | Max. |
|---|---|---|---|---|---|---|
| (Intercept) | 2.362000 | 2.400000 | 2.413000 | 2.412000 | 2.431000 | 2.441000 |
| I(Subject Prob. * Income) | 0.000260 | 0.000267 | 0.000278 | 0.000277 | 0.000286 | 0.000290 |
| Gender | 0.164500 | 0.172800 | 0.183400 | 0.185600 | 0.197500 | 0.211100 |
| Age | -0.004934 | -0.004454 | -0.004079 | -0.003941 | -0.003369 | -0.002882 |
| Education | 0.043770 | 0.044990 | 0.048450 | 0.047500 | 0.049580 | 0.050950 |
| Perception of patrolling | 0.042800 | 0.047560 | 0.051630 | 0.052600 | 0.058130 | 0.064400 |
| Victim of robbery | 0.080970 | 0.085190 | 0.087000 | 0.091490 | 0.098220 | 0.107100 |
| $\sigma$ | 0.720800 | 0.724600 | 0.730800 | 0.732400 | 0.740500 | 0.746700 |

Estimation using Gaussian adaptive bandwidth equal to 0.6149766

Source: Elaborated by the Authors

Now, in contrast to the global model, we have a parameter distribution for each variable. In this type of model, the tabular representation is not a good deal. Although we should show 8 figures (one for each of the 8 variable appearing in Table 3.8), for pragmatic reasons we present six, say, *Subject Prob. * Income*, *Gender*, *Age*, *Education*, *Perception of patrolling* and *Victim of robbery*. So we present this result in Figures (3.3), (3.4), (3.5), (3.6), (3.7) and (3.8). For example, in the east region of the city of Fortaleza, the impact of variable *Subject Prob. * Income* is slightly greater than in west regions (the difference is in the 5th decimal place). The same pattern occurs for variables *Age*, *Perception of patrolling* and *Victim of robbery*.

---

[13]  To estimate this model, we use the R statistical software (R Core Team (2014)), more specifically packages "maxLik" (Henningsen and Toomet (2011)) and "spgwr" (Bivand and Yu (2013)).

Figure 3.3 – Estimated parameters of spatial distribution - Subject Prob. * Income



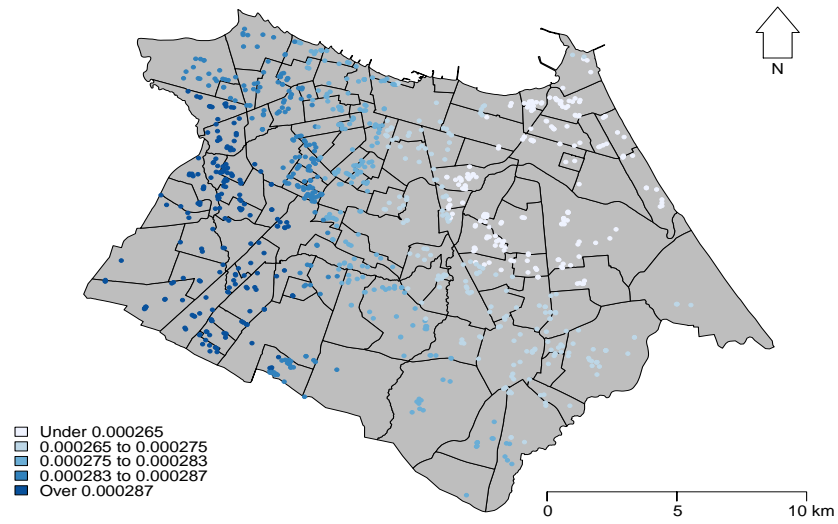Under 0.000265
0.000265 to 0.000275
0.000275 to 0.000283
0.000283 to 0.000287
Over 0.000287

Figure 3.4 – Estimated parameters of spatial distribution - Gender



Under 0.171026
0.171026 to 0.178003
0.178003 to 0.191714
0.191714 to 0.200214
Over 0.200214

Figure 3.5 – Estimated parameters of spatial distribution - Age



Under −0.004543
−0.004543 to −0.004253
−0.004253 to −0.003782
−0.003782 to −0.003254
Over −0.003254

0          5          10 km

Figure 3.6 – Estimated parameters of spatial distribution - Education



Under 0.044583
0.044583 to 0.046489
0.046489 to 0.049046
0.049046 to 0.049871
Over 0.049871

0          5          10 km

Figure 3.7 – Estimated parameters of spatial distribution - Perception of patrolling



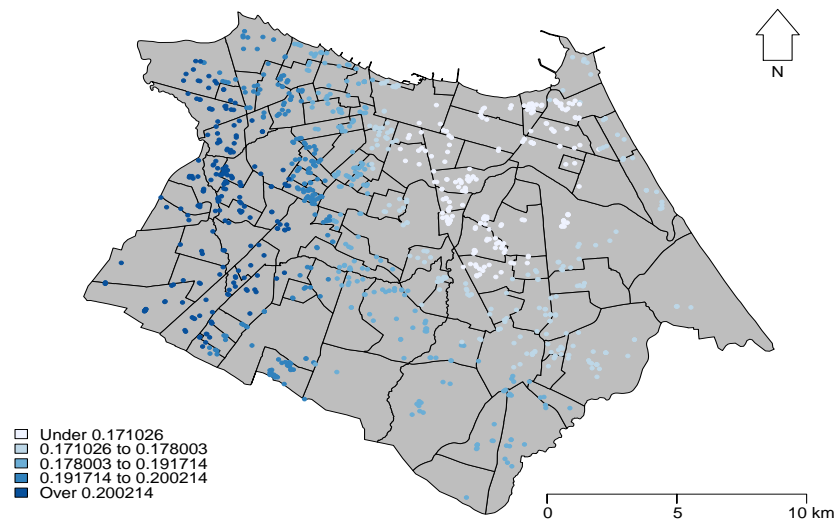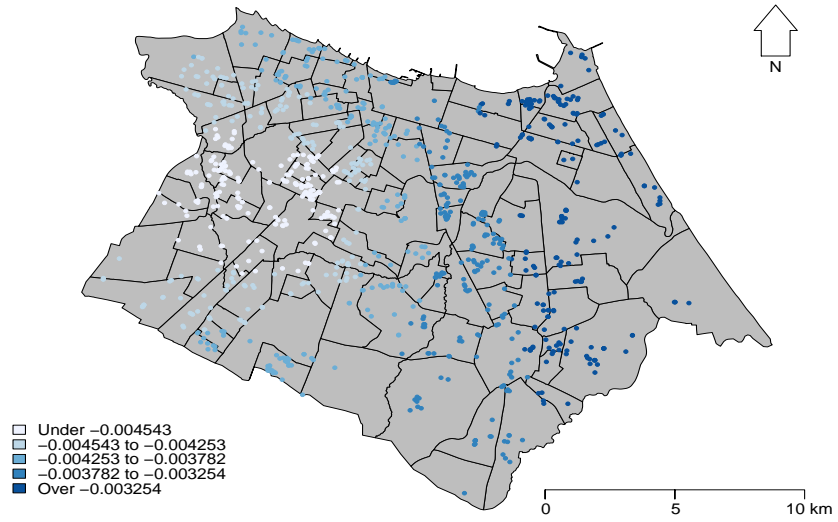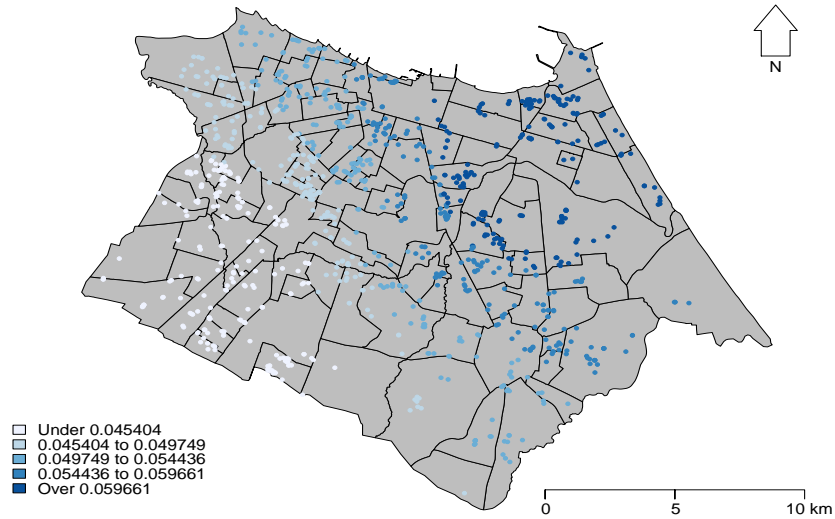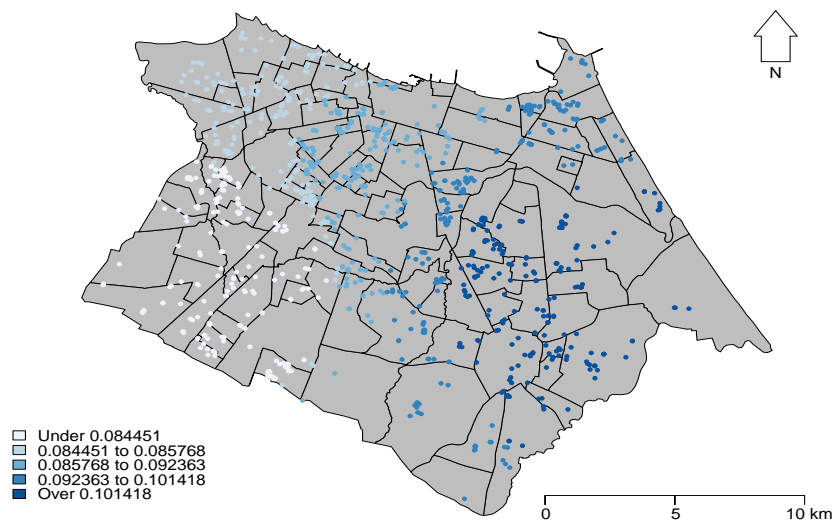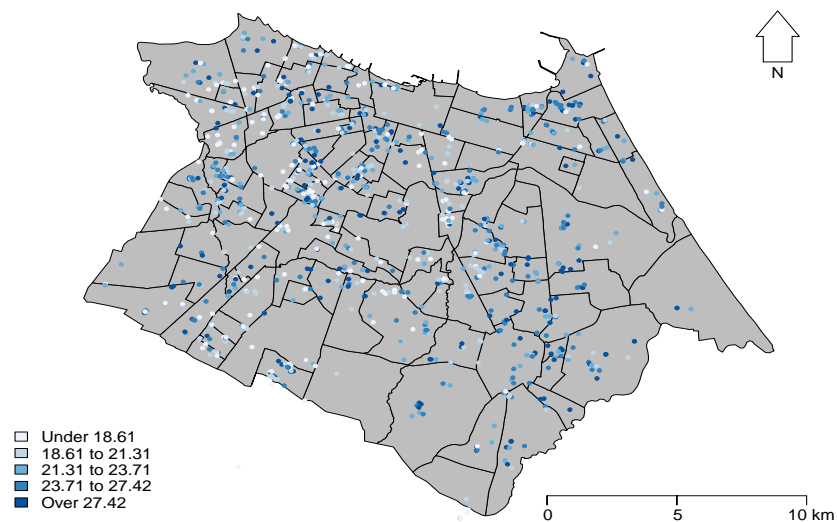Figure 3.8 – Estimated parameters of spatial distribution - Victim of robbery



For variables *Gender* and *Education*, the reverse pattern occurs. The impact is grater in eastern regions. Note that with this parameter distribution, it is possible to create a willingness to pay's distribution. In order to do that, we plug in the parameter vector

into each individual's vector of observations and calculate the expected willingness to pay and sort them into six classes. Figure 3.9 shows us the spatial distribution of willingness to pay. In this figure, we see that the highest values of estimated willingness to pay are concentrated in the central region of the city and in the southeastern region, in the prime area. This area is populated by rich people and is the area where the greatest amount of robberies in the city is concentrated, which can explain this concentration of the greatest values of willingness to pay.

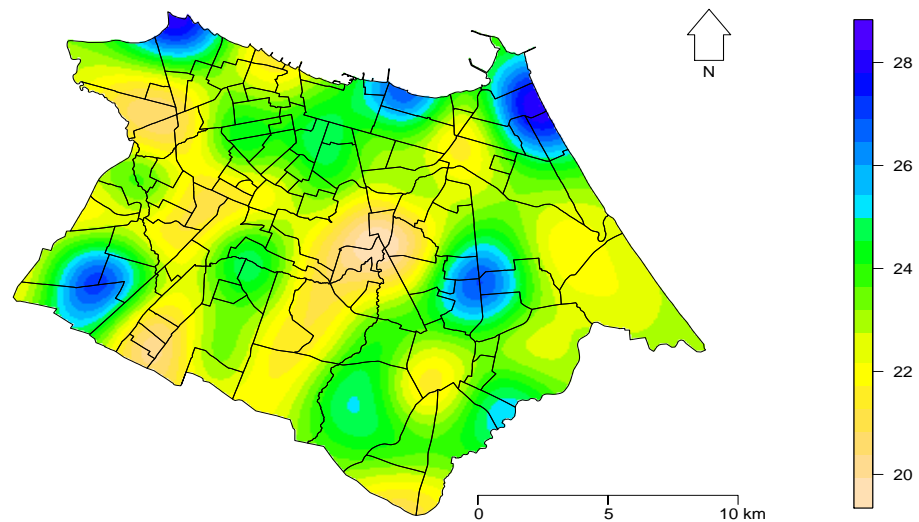Figure 3.9 – Willingness to Pay - Spatial Distribution



Before we compare these values with those from the global model, we will construct an interpolated surface to predict the willingness to pay for the entire city of Fortaleza. To do that, we will use the Ordinary Kriging technique[14]. In this map (see, Figure 3.10), we can see that in the central west and central south regions of the city have a low willingness to pay, represented by lightening colors. Although the crime rate is extremely high in this area, the types of crimes which occur are crimes against life, while robberies are less common. Furthermore, this area is populated by low-income people, who have little to be stolen. On the other hand, in dark areas, we have a high willingness to pay. In this area, the reverse pattern occurs. There is a high level of robberies and a low level of crimes against life. Moreover, this population is composed by high-income people, who have more to be stolen. The exception in this pattern occurs in dark areas in the west region, where there is a large concentration of drug trafficking activities. Thus, this heterogeneity in

---

[14]   For more details of this method, see, among others Druck et al. (2004) and Bivand, Pebesma and Gómez-Rubio (2008)

spatial distribution of crime and income may be a explanation to the spatial heterogeneity in willingness to pay distribution.

Figure 3.10 – Willingness to Pay - Kriging Surface



Comparing the values from the local model to those from the global model, we can see that the willingness to pay throughout almost the whole city is lower then the average global willingness to pay. Therefore, in case of implementing a tax, if the value set is equal to R\$ 23.35, many people lose welfare, once the value they will pay is higher than they intend to. On the other hand, there are many people who want to pay more than the potential tax, so the government will lose funds from this group. So, a flat tax to finance crime reductions is not efficient. A first degree price discrimination (see, (VARIAN, 2006)), where each unit of a good must be sold to an individual at his/her reservation price or his/her maximum willingness to pay, might be a better solution, although politically difficult. Therefore, an efficient and ideal way to go might be to determine the tax value by areas, setting it to the estimated maximum willingness to pay in that area.

## 3.6 FINAL CONSIDERATIONS

This paper sought to apply a new methodological approach to estimate willingness to pay in a large urban center in Brazil that includes spatial effects in the analysis. We constructed a theoretical model that explains the determinants of willingness to pay from the random utility model for a first-order stochastic improvement on the odds of being robbed in the city of Fortaleza, Brazil, when subjective expectations about the risk are

available. We showed that the determinant factors explaining willingness to pay in the city of Fortaleza are *Subject Prob\*Income*, *Gender*, *Age*, *Education*, *Perception of patrolling* and *Victim of robbery.*

From the global model, we estimated a mean WTP of R$ 23.35 per month/household as the value that the representative citizens of Fortaleza would be willing to pay to reduce the amount of robberies in the city in 50%. From this value, we calculated in approximately R$ 198.92 million the total cost to society, equal to 20.63% of the total amount spent on public security in the state of Ceará in 2011. We also estimated the WTP per robbery avoided equal to R$ 11,969.

Our local model, utilizing an adaptive gaussian kernel function with a bandwidth equal to approximately 0.6149, estimated a geographically weighted regression with an interval regression that, to the best of our knowledge, it is the first study to do so. We showed that in almost the whole city, the willingness to pay estimated in the local model is lower than one estimated in the global model and that there is an island where this value is greater. So, in case of tax implementation, the most efficient procedure is to discriminate the tax according to the area.

As suggestions for future studies, we believe that the construction of a new model relaxing the hypothesis of risk neutrality is a fine way to go. Finally, replicating our empirical exercise on different data sets coming from different institutional backgrounds might be something worth pursuing in order to validate our approach.

# REFERENCES

ALMEIDA, E. *Econometria espacial aplicada*. [S.l.]: Alínea, 2012.

ARAÚJO, A. F. V. d.; RAMOS, F. S. Estimação da perda de bem-estar causada pela criminalidade: o caso da cidade de João Pessoa - pb. *Economia*, v. 10, n. 3, p. 577–607, set/dez 2009.

ATKINSON, G.; HEALEY, A.; MOURATO, S. Valuing the costs of violent crime: a stated preference approach. *Oxford Economic Papers*, v. 57, n. 4, p. 559–585, 2005.

BIVAND, R.; PEBESMA, E.; GÓMEZ-RUBIO, V. *Applied spatial data analysis with R: analysis with R*. [S.l.]: Springer, 2008. (Use R!). ISBN 9780387781716.

BIVAND, R.; YU, D. *spgwr: geographically weighted regression*. [S.l.], 2013. R package version 0.6-24. Disponível em: <http://CRAN.R-project.org/package=spgwr>.

BOWEN, H. R. The interpretation of voting in the allocation of economic resources. *Quarterly Journal of Economics*, v. 58, p. 27–48, 1943.

CAMERON, T. A.; DESHAZO, J. Demand for health risk reductions. *Journal of Environmental Economics and Management*, v. 65, n. 1, p. 87–109, 2013.

CAMERON, T. A.; DESHAZO, J.; STIFFLER, P. Demand for health risk reductions: a cross-national comparison between the us and canada. *Journal of Risk and Uncertainty*, v. 41, n. 3, p. 245–273, 2010.

CAMERON, T. A.; HUPPERT, D. D. Ols versus ml estimation of non-market resource values with payment card interval data. *Journal of Environmental Economics and Management*, v. 17, n. 3, p. 230 – 246, 1989.

CARSON, R. T.; HANEMANN, W. M. Handbook of environmental economics: valuing environmental changes. In: _____. [S.l.]: Elsevier Science, 2005. (Handbook of Environmental Economics, v. 2), cap. 17.Contingent Valuation, p. 822–920. ISBN 9780080457499.

CARVALHO, J. R. Montagem de uma base de dados longitudinal de vitimização do ceará: aspectos sócio-econômicos e espaciais. Relatório Final, FUNCAP. 2012.

CIRIACY-WANTRUP, S. V. Capital returns from soil-conservation practices. *Journal of Farm Economics*, v. 29, p. 1181–1196, 1947.

COHEN, M. A. et al. Willingness-to-pay for crime control programs. *Criminology*, v. 42, n. 1, p. 89–110, 2004.

DRUCK, S. et al. *Análise espacial de dados geográficos*. EMBRAPA, 2004. Disponível em: <www.dpi.inpe.br/gilberto/livro/analise>.

ESTATÍSTICA, I. B. de Geografia e. *Censo demográfico 2010*. 2012. Http://www.censo2010.ibge.gov.br/apps/mapa/. Acessado em: 25/10/2012.

FOTHERINGHAM, A. S.; BRUNSDON, C.; CHARLTON, M. *Geographically weighted regression: the analysis of spatially varying relationships*. [S.l.]: John Wiley & Sons, 2003.

HENNINGSEN, A.; TOOMET, O. maxlik: A package for maximum likelihood estimation in R. *Computational Statistics*, v. 26, n. 3, p. 443–458, 2011. Disponível em: <http://dx.doi.org/10.1007/s00180-010-0217-1>.

HOYOS, D.; MARIEL, P. Contingent valuation: past, present and future. *Prague Economic Papers*, v. 2010, n. 4, p. 329–343, 2010.

JORGENSEN, B. et al. Protest responses in contingent valuation. *Environmental and Resource Economics*, v. 14, n. 1, p. 131–150, July 1999.

LUDWIG, J.; COOK, P. J. The benefits of reducing gun violence: evidence from contingent-valuation survey data. *Journal of Risk and Uncertainty*, v. 22, n. 3, p. 207–26, 2001.

MAS-COLELL, A.; WHINSTON, M. D.; GREEN, J. R. *Microeconomic theory*. [S.l.]: Oxford University Press, 1995.

PúBLICA, F. B. de S. *Anuário brasileiro de segurança pública*. [S.l.], 2012.

R Core Team. *R: a language and environment for statistical computing*. Vienna, Austria, 2014. Disponível em: <http://www.R-project.org/>.

SOEIRO, M.; TEIXEIRA, A. A. *Determinants of higher education students' willingness to pay for violent crime reduction: a contingent valuation study*. [S.l.], July 2010. Disponível em: <http://ideas.repec.org/p/por/fepwps/384.html>.

STRAZZERA, E. et al. The effect of protest votes on the estimates of wtp for use values of recreational sites. *Environmental and Resource Economics*, v. 25, n. 4, p. 461–476, August 2003.

VARIAN, H. *Microeconomia - Principios básicos*. [S.l.]: CAMPUS, 2006. ISBN 9788535216707.

WHITEHEAD, J. C.; BLOMQUIST, G. C. Handbook on contingent valuation. In: _____. [S.l.]: Edward Elgar Publishing Limited, 2006. cap. The Use of Contingent Valuation in Benefit Cost Analysis, p. 92–115.

# 4 PEER EFFECTS AND ACADEMIC PERFORMANCE IN HIGHER EDUCATION - A REGRESSION DISCONTINUITY DESIGN APPROACH

## 4.1 INTRODUCTION

Human beings are social creatures. This is based not only on the fact that we like company or depend on each other. Human beings are social creatures simply in the sense that our existence requires interaction with other people (GAWANDE, n.d). In the last decades, economists have devoted great attention to these interactions and its influence on individual behavior. The effects of these interactions are known in the literature as "peer effects".

Sacerdote (2011) defines peer effect as any externality, excluding those market-based or price-based, in which peers' background, current behavior, or outcomes exert an influence on a specific outcome obtained by another individual. Manski (1993) classifies this effect as endogenous, when it emanates from peers' current outcomes, and exogenous, when it is due to peers' backgrounds.

Several studies have analyzed peers' influence on criminal activity, drugs use, teenage pregnancy, educational achievement, among others (SACERDOTE, 2011). Looking specifically at the literature concerning educational achievement, peer effects have played an important role for primary and secondary education since Coleman et al. (1966) seminal work, being considered a key factor in determining children's schooling outcomes (WINSTON; ZIMMERMAN, 2004).

Even though the importance of peer effects in elementary and secondary education had been raised a long time ago, its relevance to the economics of undergraduate/graduate degrees has only recently been acknowledged (WINSTON; ZIMMERMAN, 2004). Thenceforward, this research agenda experimented an exponential growth, with several studies seeking to take a deeper look at the peer effects for higher education. So far, the empirical results bring up contradictory conclusions. Some studies find a positive effect on academic outcomes due to peers' influence, while others show negative effects or even no effect at all (EPPLE; ROMANO, 2011; SACERDOTE, 2011).

Commonly, the literature has been using roommates interaction as a standard source of peer effects. This is the case for works such as Sacerdote (2001),Zimmerman (2003) and McEwan and Soderberg (2006). On the other hand, studies like Paola and Scoppa (2010),

Androushchak, Poldin and Yudkevich (2012) and Booij, Leuven and Oosterbeek (2015) prefer to take advantage of classmates interactions. Our paper follows this guidance and uses the interactions between classmates as well. However, it contributes to the literature presenting a different group formation.

Based on this set of papers, we believe our endeavor has its own merits. Firstly, we estimate peer effects in higher education for a developing country with an institutional background which is very different from what is found in OECD members, for example. Secondly, by aggregating new methodological procedures, we can better understand the relation between peer effects and academic performance in high education.

Therefore, this paper aims to estimate peer effects of undergraduate students on academic outcomes. For this purpose, we used a micro data set conceded by the Federal University of Ceará, a public Brazilian university located in Fortaleza, the capital city of Ceará. Our data set brings several socioeconomic information and maps concerning 4 years of academic performance with respect to 2149 students enrolled in 33 undergraduate programs. Due to the entrance process specificities, we are able to estimate peer effects using a sharp regression discontinuity design. Also, in light of programs' heterogeneity, and since the assignment grade distribution pattern is different, we are able to estimate a multi-treatment effect model. For this, we classify each program according to the competition in its first and second semester classes.

We found that peer effects have a negative impact on the academic performance of our undergraduate students. The evidence suggests that low-ranked students put together with high-ranked classmates have a worse academic performance than those in a lower level class. This goes against several studies of peer effects for primary and high schools, as well as for higher education.

Notwithstanding, for a multi treatment model, we also found evidence of non-linearities as in Sacerdote (2001) and Zimmerman (2003). We found positive peer effects when both first and second semester classes are of low competition level, and negative peer effects in all other configurations, with modest magnitudes when both classes are of high competition level.

Besides this introduction and a final considerations section, this paper presents five more sections. Section 2 offers a brief literature review of peer effects on academic outcomes. Section 3 introduces the entrance process for Brazilian universities and demonstrates that it follows a sharp design. Section 4 scrutinizes our data, presenting the results of an initial exploratory analysis. Section 5 sets up the model to be estimated and a brief discussion of the estimation method. Finally, section 6 presents our results.

## 4.2 BRIEF EMPIRICAL LITERATURE REVIEW

### 4.2.1 Peer effects

The literature of peer effect on academic achievement has grown significantly in recent years. Studies trying to access its role in elementary, secondary and post-secondary educational levels are of particular interest in many countries, raising different methodological approaches to take the particularities of each educational level and backgrounds into account.

Sacerdote (2001) estimated peer effects among Dartmouth College (USA) roommates. He found that peers have an important impact on students' grades and on the decision to join social groups such as fraternities. Also, the paper attests a non-linearity in these peer effects: students whose roommates were in the top 25% of the class had higher grades. Sacerdote (2001) concluded that high-ability students had a positive effect on the academic achievement of relatively less talented colleagues, while there was no such influence for students in the middle of distribution.

Similarly, Zimmerman (2003) studied peer effects among Williams College (USA) undergraduate students. In this paper, since first year roommates were assigned randomly with respect to academic ability, the author could estimate differences in grades of high, medium, and low SAT students living with high, medium or low SAT roommates. The results indicated that a medium student tended to have worse grades if put together with a low SAT roommate, while high ability students were least influenced by peers.

McEwan and Soderberg (2006), in a study carried out at Wellesley College (USA), estimated the effects of students' background characteristics on their roommates' academic outcomes. The authors applied both a linear and a nonlinear model. Regarding the first structure, there is no evidence of peer effects on students' GPA. With respect to the nonlinear specification, the results suggest that students' SAT scores have a nonlinear effect on their roommates' achievement, yet the results are not robust. The conclusion is that there might exist roommate peer effects restricted to a small number of students. However this effect is not a key determinant for academic outcomes.

Carrell, Fullerton and West (2008) also estimated peer effects in college achievement. The paper uses data from the United States Air Force Academy, in a context in which students are exogenously assigned to peer groups. The interaction is even stronger in this case, since required activities involve both academic and non-academic duties. They find a scholarly peer influence larger than those found in previous studies relating to roommates. Furthermore, peer effect persists at a diminishing rate into sophomore, junior, and senior years, indicating long lasting ties on academic achievement.

Contreras, Badua and Adrian (2012) investigated peer effects for classroom colleagues in a Business College of a U.S. public university. The authors find a negative significant peer effect in students' performance, yet the proper direction and magnitude are sensitive

to both peers' and student's own average ability.

Besides the USA, there is a growing literature in European countries on peer effects. In Italy, Paola and Scoppa (2010) analyzed peer effects among students of Calabria University, a middle-sized public university. They found a positive and statistically significant influence, being robust to different peer group definitions and abilities measures. Also, the effect was larger than previous studies focusing on roommates. The results attest that students' ability is an important input in college education, which means that attracting high-level students is a key path to improve the overall performance by means of direct and indirect directions.

Androushchak, Poldin and Yudkevich (2012) used data about Russian undergraduate students enrolled at the Economics department of the National Research University — Higher School of Economics (HSE)— to estimate peer effects in exogenously formed groups. The evidence suggests that high-ability classmates exert a positive influence on individual academic performances. Still, the most talented ones are the greatest beneficiaries from this presence. The paper also finds that an increase in the proportion of low-performance students has an insignificant or negative influence on individual grades.

Regarding undergraduate students in Economics at the University of Amsterdam (NED), Booij, Leuven and Oosterbeek (2015) estimated peer effects from tutorial groups' ability composition. Aiming to achieve a wide range of support, the authors manipulated these compositions and assigned the students randomly. They find that low and medium ability students gain, on average, 0.2 standard deviation achievement units after switching from ability mixing to three-way tracking — a system in which each group is constituted by students of the same ability distribution, measured by GPA. They also find that high-ability students are not affected by the specific group composition, and defend that there is no evidence implying that teachers adjust their teaching to different group configurations.

## 4.2.2   The use of regression discontinuity design

In common, none of these papers use a regression discontinuity design approach to estimate peer effects in academic achievement. Actually, this approach is usually found in analysis of both remedial education effects and financial aid on academic achievement, due to particularities of post-secondary educational level.

Moss and Yeaton (2006) are an example of a sharp regression discontinuity design application in remedial education. This study analyzed the effectiveness of a developmental English program in an American university. The program offers compulsory remedial education to students of ASSET scores less or equal to 85 out of 107 points. The authors found that those students participating in the program had their English academic achievement similar to those initially out of supplemental coursework. Furthermore, the students in greatest need of the program had the major benefit from it.

Another example is the study of Butcher, McEwan and Taylor (2010). It estimated the causal effect of taking a course in quantitative reasoning on student's academic performance and classroom peer-group composition at Wellesley College (USA). The assignment rule in this program is similar to the one presented in Moss and Yeaton (2006): if the student's test score is less or equal to 9 (out of 18), then he/she is assigned to this mandatory quantitative course. The authors found that there is no impact in taking the course on student's academic outcomes. Nevertheless, they identified robust effects on classroom peer-group composition, i.e, classmates of this remedial course tend to keep studying together along other different courses.

In a similar study, Schöer and Shepherd (2013) estimated the role of taking a compulsory remedial course on students' performance in an undergraduate level microeconomics class. Using data from a South African university, they used a fuzzy regression discontinuity design and found that this program participation positively affects students' performance.

Regarding the role of financial aid on academic outcomes, the seminal paper of Klaauw (2002) analyzed the effects of universities' financial aid offers on students enrollment decisions. The author found that this recruitment resource is an effective instrument in competing with other colleges for new students. In the same line, Leeds and DesJardins (2014) found similar conclusions for the University of Iowa. In addition, the results suggest that financial aids may have strong effects on the brightest candidates.

Mealli and Rampichini (2012) analyze the relation between grants offered by an Italian university for low-income students and their dropout decision. The results suggest that, at a given threshold, the grant is an effective tool to prevent those low-income students to drop out of higher education. However, if the family income is much lower than this threshold, then the grant effect becomes smaller and not significant.

Canton and Blom (2004) analyzed the effects of financial aid on enrollment and students' performance at Mexican universities. The results indicate a positive effect for both issues. Concerning the enrollments, a strong impact is verified, since the probability of entering in higher education is raised in 24 percent. For the performances, students who receive financial aid have better academic results than those without it. A similar study was carried out by Curs and Harper (2012), attesting that students who receive financial aid have a GPA between 0.12 and 0.16 higher than students who do not.

In the studies of peer effects, the methodological framework depicted by the regression discontinuity design is more usual when dealing with elementary and high school levels. For example, Koppensteiner (2012), using Brazilian data on elementary school students, estimated the effect of being in a class of older classmates on students' achievement. The conclusion consists in a large negative impact of it. Card and Giuliano (2015) estimated the effect of being in a gifted/high achiever classroom for U.S. students. They found positive and significant effects concentrated among minorities, and found no evidence of spillovers on non-participants of the program.

The same Card and Giuliano (2015) motivation is presented in Vardardottir (2013). Refering to high school Icelander students, this paper estimated the effect of being in high-ability classes, and found a significant and sizable positive impact on the academic achievement of students around the assignment threshold. Abdulkadiroğlu, Angrist and Pathak (2014) are also concerned about peer effects at a high school level. They estimated the effect of study in a high-quality school, and the results suggest that the marked changes in peer characteristics at exam school admissions cutoffs have little causal effect on test scores or college quality.

In common, all of these studies follow a clear rule in the class group formations, which are usually not the case for U.S. and European universities. With respect to Brazil, there is a well-established rule in forming two types of undergraduate classes: the first one gets the best ranked students in the entrance process, and start academic activities in February (first semester); the second type is a place for the lowest ranked students, and runs at the beginning of August (second semester). Relying on this fact, we developed our study.

## 4.3   THE ENTRANCE PROCESS IN BRAZILIAN UNIVERSITIES

In Brazil, according to (Inep) (2014), in 2013, there were 7,3 million students enrolled in 2,391 higher education institutions. Among those institutions, 106 are public and maintained by the Brazilian federal government, counting 5.968 undergraduate programs and 1,14 million students enrolled at federal institutions. In other terms, these numbers represent 18% and 15% of Brazilian undergraduate programs and college students, respectively.

Since federal universities are free of charge and present a high teaching quality[1], they are target of many students from all social backgrounds, which translates into a high demand and, therefore, great competition for a vacancy. Thus, in order to ensure equal access, its entrance processes take place by means of a public tender.

Silva (2007) makes a historical analysis of the admission process in higher education in Brazil. In the 19th century, students who aimed to enter in higher education had to go through a series of tests after they completed the secondary education to obtain a required grade to access the higher education. These exams were called exit tests.

From 1915 on, these exams became to be called *Vestibular*, as we know nowadays, being mandatory for all students who wanted to access the higher education system. This admission process became effective, in fact, in the 1920s, when the number of candidates became higher than the number of vacancies. In this period, *Vestibular* still was an exit test.

It was after 1925 that *Vestibular* became an entry test, whose objective was to evaluate student's capability to understand studies at higher level. The exams were restricted to

---

[1]   In Brazil, public universities hold status of higher quality compared to their private counterparts.

disciplines considered pre-requisites to the undergraduate program the student intended to attend to.

In 1971, due to the pressure of the students who were unable to access higher education, new conditions of access were created. This new system established that *Vestibular* must have only one content for all programs and adopt classification criteria, in which students that obtained the greatest grades were selected. Under this system, each university became responsible to organize its own selection process, establishing the number of vacancies to be offered (SILVA, 2007). This system remained in force until 2010, when a new admission process based on ENEM and SISU arose[2].

Abreu (2013) summarizes *Vestibular's* algorithm. In a first moment, each student chooses and announces only one program of his preference. In a second moment, for each program, a preference relation is determined, utilizing the grade obtained in the exam. And finally, students are allocated based on their rankings and preferences of the chosen program. It was demonstrated that this algorithm is not stable, is not pareto efficient and is not strategy-proof.

At the time of our analysis, the Federal University which we have access to the information used *Vestibular* as its admission process. In this University, the exam consists in two stages. The access to the second stage is conditioned by the performance in the first stage. All students above a rank at the first stage exam are accepted to the second stage, which make the number of students who take the second stage exam a multiple (usually 4 sometimes 3) of the number of final available vacancies. These ranks (one for each major) define a first stage grade threshold. Similarly, second stage threshold determine who passes the exam and enters the University (CARVALHO; MAGNAC; XIONG, 2014). Based on scores achieved in the second stage, students were ranked, vacancies were filled, and the upper classified half of students for every course was assigned to start studies in the first semester of the academic year (first semester class), while the bottom half was allocated into the second semester(second semester class).

Carvalho, Magnac and Xiong (2014) demonstrated that this threshold is a Bayesian Nash equilibrium, and it is unique. This allows us to use a sharp regression discontinuity design to analyze peer effects among students allocated in first semester classes and in second semester classes.

## 4.4  DATA

Our analysis is based on a rich administrative data set, providing information on undergraduate students enrolled in 2008 at the Federal University of Ceará (UFC). This is

---

[2]   Since 2011, the new entrance process is by means of ENEM and SISU. Students take a centralized national exam (ENEM) and, in light of the obtained results, they choose any university and undergraduate program (SISU), taking their own score and the cutoff score determined by competition into account.

a public Brazilian university located in Fortaleza — the fifth largest city in Brazil with a population of 2.5 million citizens. Founded in 1954, UFC is considered one of the best universities in Brazil according to the Brazilian ministry of education, and the second best in the northeastern region. During the 2013 academic year, the university had 26,782 students enrolled in 114 undergraduate programs. With respect to graduate programs, the university had 6,061 students enrolled in 167 programs, divided into *lato sensu*, professional and academic masters, as well as doctoral programs. The Federal University of Ceará was constituted by 2,152 professors, of whom 1,436 and 543 have doctoral and master's degrees, respectively. UFC also had 3,407 administrative staff members, including the University Hospital (UFC, 2014).

The data set was collected from both the Vestibular Commission and the provost Office of Undergraduate Studies. We have information on 27 undergraduate programs with classes starting in both academic terms[3], counting 1550 students. Our data covers grades and final classification in the entrance exam for the 2008 academic year. We also bring information on the students' socioeconomic characteristics, collected at the registration stage by means of a survey held by the Vestibular commission. Concerning academic performance, students are traced from 2008 to 2011[4], which is equivalent to 8 academic semesters and to the required time before graduating. As a measure of academic achievement, we use IRA — UFC equivalent to the American GPA.

The IRA (Índice de Rendimento Acadêmico) index is a measure of student's academic performance similar to the American GPA, but in a 10-point scale. This index is used to rank students for research and teaching grants, for distinction purposes and so on. The IRA index for a student $i$ is calculated as follows:

$$IRA_i = \left(1 - 0.5\frac{T}{C}\right) \times \left(\frac{\sum_j P_i \times C_j \times N_j}{\sum_j P_j \times C_j}\right) \tag{4.1}$$

Where:

- $T$ is the sum of all withdrawn courses' workload;

- $C$ is the sum of all courses' workload, withdrawn or not;

- $C_j$ is the workload of course j;

- $N_j$ is the final grade of course j;

- $P_j$ is the period in which the course was done, obeying the following limitation: $P_j$ = min {6, semester in which the course was done}.

As shown above, the IRA index is a weighted mean. Variable $\frac{T}{C}$ measures the proportion of all withdrawn courses' workload with respect to the total amount (withdrawn or not).

---

[3]    This is equivalent to 38% of the total amount of undergraduate programs.

[4]    Or until student drops out the course.

Note that it has a negative impact on the IRA index, i.e, when this proportion becomes higher, the IRA index turns lower. So, the withdrawal of any course is penalized with a reduction of the student's IRA index.

When it comes to courses concluded, some important comments should be made. Firstly, if a course is attended more than once, which is the case for failures, the same number of times appearing in the student's transcript of records will be included in the IRA index calculation. Secondly, in case of course failure by attendance, the final grade will be zero. Table 4.1 presents the variables used in this paper.

Table 4.1 – Variables descriptions

| Variable | Description |
| --- | --- |
| IRA | Student's academic index |
| SAG | Standard student's grade obtained in the Vestibular exam |
| Age | In years |
| Gender | 1 if student is male; 0 if female |
| Log(income) | In R$ |

Source: Elaborated by the Authors

From table 4.1, variable *SAG* deserves special attention. We defined it as the difference between student's final grade in vestibular and that of the last ranked student within the same semester class — i.e. the class cutoff grade — divided by the standard deviation of the student's course. This procedure helps us to have all first semester students with positive *SAG*, while second semester students will have it negatively, with *zero* as its cutoff. In summary, all courses will have the same cutoff now.

Table 4.2 – Descriptive statistics

| Variable | Mean | Std. Dev. | N | Mean | Std. Dev. | N | Mean | Std. Dev | N |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | Full Sample | | | First Semester | | | Second Semester | |
| IRA | 7.8307 | 1.3208 | 12400 | 7.9690 | 1.2408 | 5992 | 7.7014 | 1.3790 | 6408 |
| SAG | 0.1722 | 0.9779 | 1550 | 0.9562 | 0.8191 | 749 | -0.5609 | 0.3328 | 801 |
| Age | 19.0316 | 2.9165 | 1550 | 18.8051 | 2.4943 | 749 | 19.2434 | 3.2495 | 801 |
| Gender | 0.4406 | 0.4966 | 1550 | 0.4419 | 0.4969 | 749 | 0.4395 | 0.4966 | 801 |
| Log(income) | 7.4835 | 1.0105 | 1550 | 7.5025 | 1.1091 | 749 | 7.4657 | 0.9090 | 801 |

Source: Elaborated by the Authors
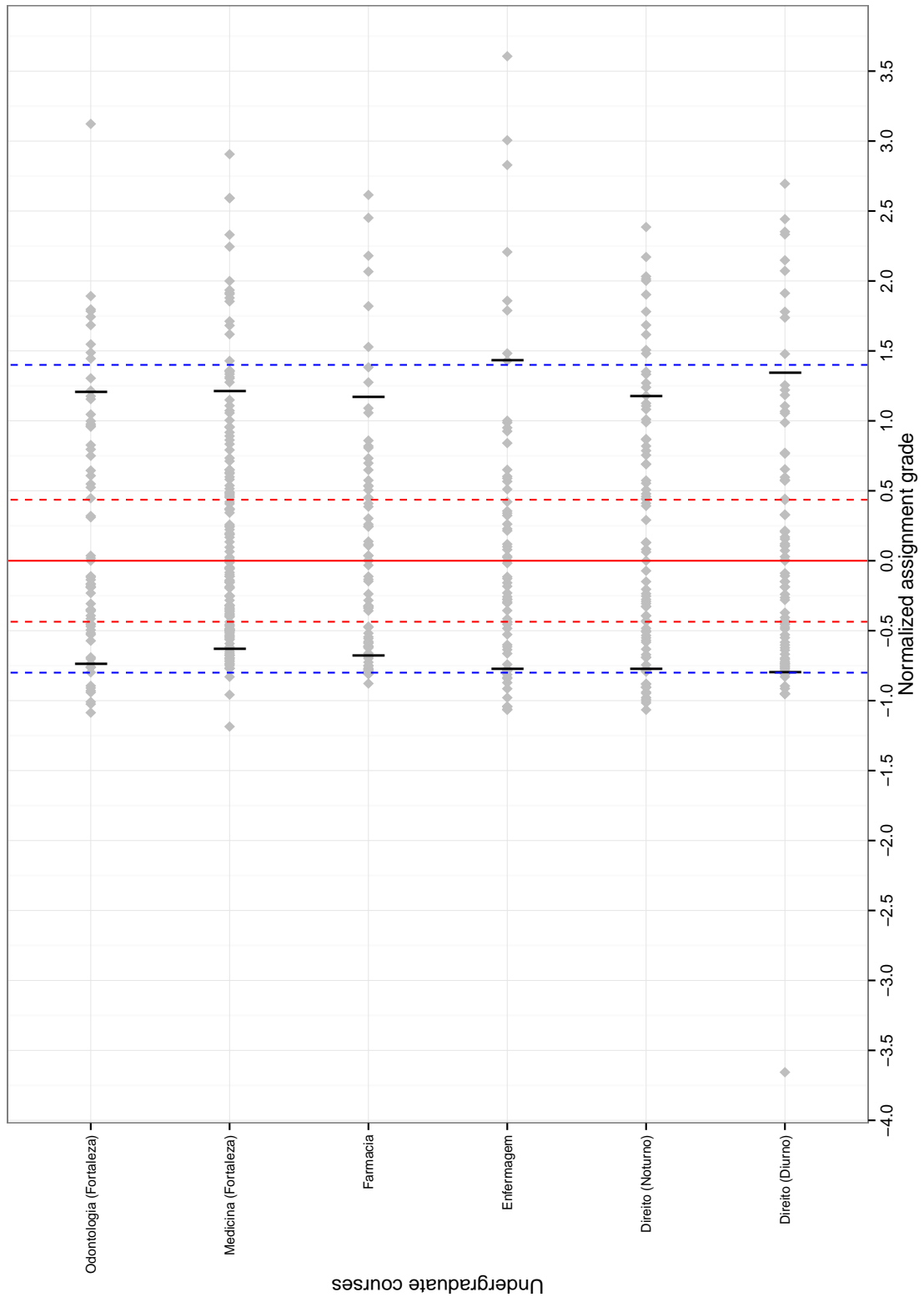Note: $N$ differs in $IRA$ because there were students who dropped out the course.

Table 4.2 shows the socioeconomic profile and academic performance of our students. The first three columns relate to the full sample. We can see that 44% of the students are males, with average age and log(income) equal to 19.03 and 7.4835, respectively. The average of students' academic performance, measured by *IRA*, equals 7.8307. Finally, the average *SAG* is 0.1722.

Now, looking at first and second semester classes, we can see that these two groups have similar socioeconomic characteristics. Note that both are composed by approximately 44% of males with 19 years old and a log(income) of 7.5 and 7.46, for the first and second

group, respectively, in average. The academic performance of the first semester group is 7.96, and 7.70 to the second.

Finally, by definition, variable $SAG$ has a positive or negative sign depending on the reference class — first or second semester, respectively. Concerning the spread of $SAG$, the standard deviation in the first group is larger than in the second, which shows that the class beginning in the last academic semester is more homogeneous. This is not a coincidence, since students ranked in the top class tend to be more prepared and achieve higher grades in vestibular. It also involves the top 1% students, whose grades are possibly too far from average. We can see this better in figures 4.1, 4.2, 4.3, 4.4 and 4.5.

Figure 4.1 – Distribution of normalized assignment grade by courses - Medical school and Law school



These figures display the $SAG$ distribution for each course. Figures 4.1, 4.2, 4.3 and

Figure 4.2 – Distribution of normalized assignment grade by courses - College of Economics, Management, Actuarial Science and Accounting

Figure 4.3 – Distribution of normalized assignment grade by courses - College of Sciences and College of Agricultural Sciences
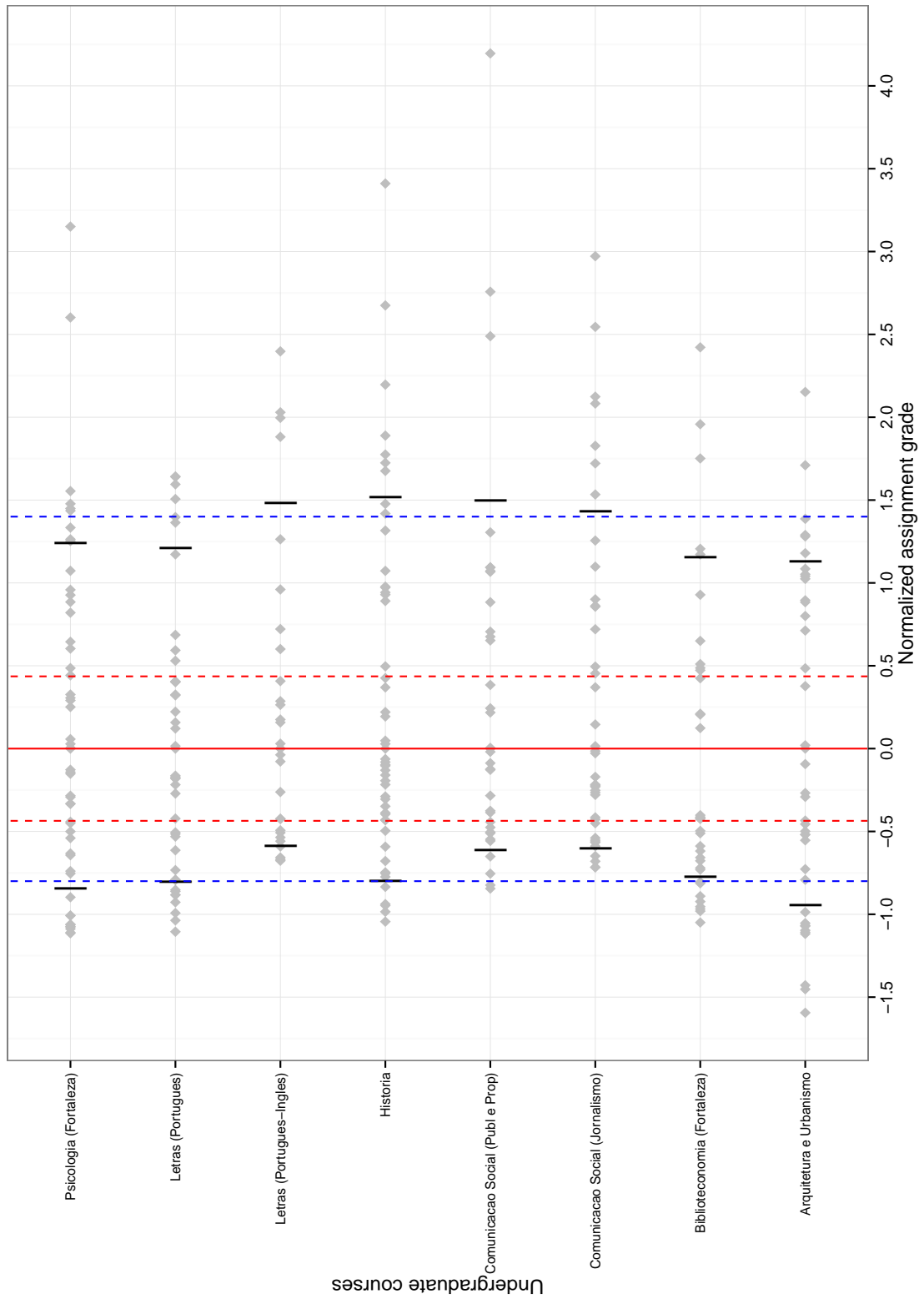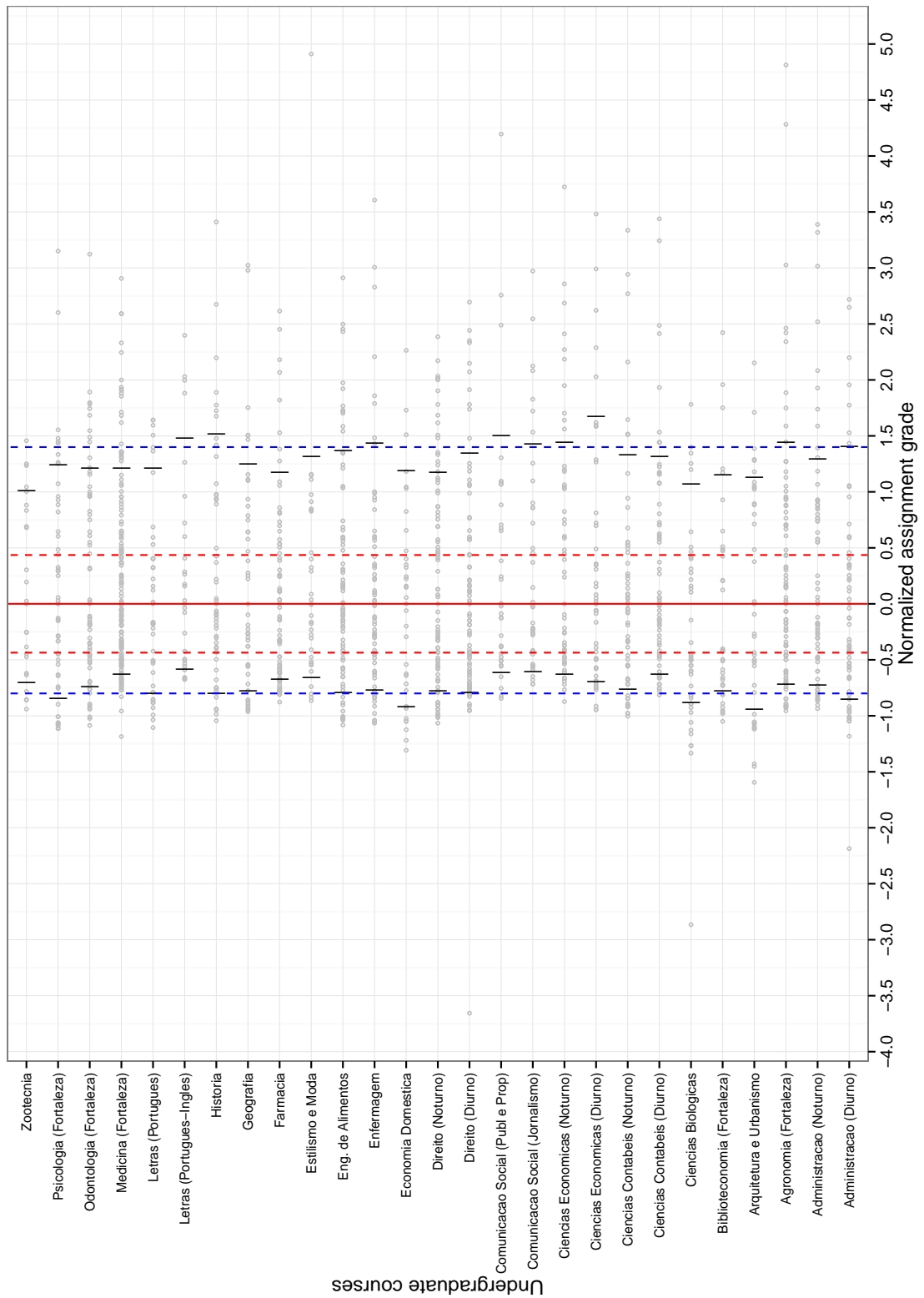
Figure 4.4 – Distribution of normalized assignment grade by courses - College of Humanities



4.4 present the courses separated by academic units, while 4.5 shows all courses. Points

Figure 4.5 – Distribution of normalized assignment grade by courses



relate to individual observations, while bars represent the average $SAG$ in each class. We can clearly see heterogeneity in $SAG$ distributions. For example, courses like *Ciencias*

*Economicas (Diurno)* (4.2) have a high *SAG* average for first semester classes, while this is not the case for courses such as *Economia Domestica* (4.3) . Another example of this heterogeneity is due to the distribution dispersion. In courses like *Medicina* (4.1) , the dispersion is very low, with the dots very close to each other, while in *Zootecnia* (4.3) this is exactly the opposite. In light of such evidence, we defined four categories or levels of treatment, according to the competition degree in each class, measured by the average *SAG*. These definitions are given in table 4.3.

Table 4.3 – Definitions of treatments groups

| Group | Control ($2^{nd}$ S) | Treatment ($1^{st}$ S) | Definition |
|-------|----------------------|------------------------|------------|
| T1 | $SAG < -0.8$ | $SAG \in [0,1.4]$ | Courses with lower competition in $2^{nd}$ S class and lower competition in $1^{st}$ S class. |
| T2 | $SAG \in [-0.8,0)$ | $SAG \in [0,1.4]$ | Courses with higher competition in $2^{nd}$ S class and lower competition in $1^{st}$ S class. |
| T3 | $SAG < -0.8$ | $SAG > 1.4$ | Courses with lower competition in $2^{nd}$ S class and higher competition in $1^{st}$ S class. |
| T4 | $SAG \in [-0.8,0)$ | $SAG > 1.4$ | Courses with higher competition in $2^{nd}$ S class and higher competition in $1^{st}$ S class. |

Source: Elaborated by the authors

Thus, we defined as high (or low) competition, classes whose average *SAG* is on the right (or left) of the dashed blue line[5], limited by the solid red line. Each class threshold values were set *ad hoc*. There are lower and upper bounds of the entire sample, representing the bottom and top 10%, respectively. Next section presents the empirical model to analyze the effect of being in the first semester class.

## 4.5   EMPIRICAL STRATEGY

Since the vestibular exam has a sharp design, we can estimate the effects of being in the first semester class on students' academic outcomes by the following model:

$$IRA_{it} = \beta_0 + \beta_1 T_i + \beta_2 SAG_i + \beta_3 T_i SAG_i + \delta X_i + \alpha_k + \alpha_t + \varepsilon_i \qquad (4.2)$$

Where $IRA_i$ is the academic performance index for student $i$, $T_i$ is a dummy variable indicating whether a student $i$ belongs to the first semester class, $SAG_i$ is the standardized assignment grade, $X_i$ is a student-specific vector of control variables such as age, gender and income. Given that our data set has information for four years (8 semesters), we are able to include fixed effects for courses and time, $\alpha_k$ and $\alpha_t$, respectively, enabling us to improve our estimators' efficiency. Finally, $\varepsilon_i$ is the error term.

---

[5]    The dashed red lines define the bandwidth used in the Sharp Regression Discontinuity Design model.

This model estimation can be done by means of parametric and nonparametric techniques. For the first case, the exercise must include high-order polynomial for the forcing variable in the model and use data distant from the cutoff. Following a nonparametric estimation, we must choose a window of width $h$ around the cutoff, and use this *local* data to perform the estimation.

Gelman and Imbens (2014) present three arguments against the use of high-order polynomials. Firstly, the implicit weights for approximations are not attractive; secondly, the results are sensitive to the polynomial approximation order; lastly, conventional inferences hold poor properties under this setting. Therefore, the authors suggest estimators based on smooth functions such as local linear or quadratic polynomials. So, following this guidance, we estimate our models by means of nonparametric techniques. Specifically, the approach here is a local linear regression.

The local linear regression is a nonparametric way to consistently estimate the treatment effect in a regression discontinuity design (LEE; LEMIEUX, 2009). This method consists in fitting linear regression functions to observations within a distance $h$ on either side of the discontinuity point. Then, treatment effects are given by the difference in intercept estimative for these two equations. Alternatively, one can estimate the average effect directly in a single regression, by solving equation 4.3 (IMBENS; LEMIEUX, 2008):

$$min_{\beta,\delta} = \sum_{i=1}^{N} 1\{c-h \leq SAG_i \leq c+h\}.(Y_i - \beta_0 - \beta_1 T_i - \beta_2 SAG_i - \beta_3 T_i SAG_i - \delta X_i)^2 \quad (4.3)$$

In this type of model, the researcher faces two important issues: selecting the kind of kernel function to be used, and, more importantly, the bandwidth determination. With respect to the first point, Imbens and Lemieux (2008) advocate that the use of a rectangular kernel, or a more sophisticated version, do not make much difference in the asymptotic bias. In this sense, if there is a difference when one varies the weights of a more sophisticated kernel, it is that the results are highly sensitive to the bandwidth. Hence, the only case in which more sophisticated kernels might be alluring is when the estimates are not much credible due to a high sensitivity in this bandwidth choice.

Even though the arguments presented in Imbens and Lemieux (2008) must be taken seriously, we will proceed with a triangular kernel. This is based on the well-known result that this kernel is an optimal choice for estimating local linear regressions at the boundary (FAN; GIJBELS, 1996). The triangular kernel function is given by the following expression:

$$K(u) = (1 - |u|) \quad \text{for} |u| \leq 1, \quad \text{where} \quad u = \frac{X_i - X_c}{h} \quad (4.4)$$

Where $X_c$ is the cutoff point and $h$ is the bandwidth.

The triangular kernel puts more weight, linearly, on observations closer to the cutoff point. So, the difference between regressions using a rectangular or triangular kernel is

that the latter involves estimating a weighted regression within a bin of width $h$, while the former is an unweighted regression (LEE; LEMIEUX, 2009).

The bandwidth determination is more intricate. According to Lee and Lemieux (2009), setting it in a nonparametric structure involves finding an optimal balance between precision and bias. If the researcher uses a larger bandwidth, more observations are available and thus he/she can obtain more precise estimates. However, the linear specification is less likely to be accurate when a larger bandwidth is used, which can bias the treatment effects estimation.

Due to the previously mentioned problems, the bandwidth ($h$) choice must be made guided by the available data to avoid arbitrary choices, and always taking its trade-off between bias and efficiency into account. To this task, we will follow Imbens and Kalyanaraman (2011)'s algorithm. This algorithm is developed to the bandwidth estimation, focusing on the local linear regression approach. The authors derived an asymptotically optimal bandwidth, conditioned on unknown data distribution functionals, and then proposed simple and consistent estimators for these functionals, obtaining a fully data-driven bandwidth algorithm.

The optimal bandwidth estimator proposed in Imbens and Kalyanaraman (2011) is given by:

$$\hat{h}_{opt} = C_K \cdot \left( \frac{\hat{\sigma}_-^2(c) + \hat{\sigma}_+^2(c)}{\hat{f}(c) \cdot ((\hat{m}_+^{(2)}(c) - \hat{m}_-^{(2)}(c))^2 + (\hat{r}_+ + \hat{r}_-))} \right)^{1/5} \cdot N^{-1/5} \tag{4.5}$$

Where the quantities $\hat{\sigma}$, $\hat{m}$, $\hat{f}(c)$ and $\hat{r}$ are, respectively, the conditional variance, conditional mean, the marginal distribution of the forcing variable X at threshold $c$, and the regularization term. Subscripts $+$ and $-$ are to identify the right or left positioning with respect to the threshold [6].

For the multi-treatment model estimation, we follow this same logic. The only difference here is that we subset our data into four competition levels, and run this model for each different category. With the estimates in hand, we are able to define two measures, commonly used in literature of multi-treatment effects: i) the incremental comparison, in which successive levels of treatment are compared; ii) the control comparison, where the different treatment levels are compared to a reference level (LEE, 2005). According to Lee (2005), assuming that the treatment effect at level $i$ is given by $\mu_i$, we have:

- Incremental effect: $\mu_i$ - $\mu_{i-1}$, $\forall i$

- Comparison with the control effect: $\mu_i$ - $\mu_0$, $\forall i$ When treatment 0 is the control

This is the methodological framework used in this paper. Next section presents the estimation results.

---

[6]  For more details, see the complete work of Imbens and Kalyanaraman (2011).

## 4.6   RESULTS

Grounded on the Imbens and Kalyanaraman (2011) algorithm, we obtained a bandwidth of 0.4360537[7]. Table 4.4 shows the estimation results of model 4.2 under this value[8].

Table 4.4 – SRD estimates for the effect of being in the first semester class on students' IRA

| Variable | Estimate | Stad. Dev. | t value | p-value |
|----------|----------|------------|---------|---------|
| Intercept | 8.5577 | 0.2138 | 40.0266 | 0.0000 |
| Tr | -0.1973 | 0.0533 | -3.6990 | 0.0002 |
| SAG | -0.0411 | 0.2220 | -0.1852 | 0.8531 |
| Tr*SAG | 0.7509 | 0.3009 | 2.4952 | 0.0126 |
| Age | -0.0201 | 0.0062 | -3.2707 | 0.0011 |
| Gender | -0.2546 | 0.0345 | -7.3877 | 0.0000 |
| Log(income) | -0.0362 | 0.0177 | -2.0401 | 0.0414 |
| Bw = 0.4360537 | $\bar{R}^2$=0.4474 | N=4256 | Nl=2376 | Nr=1880 |

Source: Elaborated by the authors

The results presented in table 4.4 show that, in fact, there is a significant difference in academic performance between students in the first semester class just above the cutoff, and those in the second semester class just below this threshold. In the first case, students have a lower academic performance when compared to those (quite similar in the vestibular results) starting their studies in the second academic semester. The magnitude of this negative effect is about 0.1973, as indicated by the coefficient of variable $Tr$. This represents a 2% decrease in $IRA$, since it is measured in a 10-point scale.

In other words, we find that, contrary to what usually happens in peer effects studies for primary and high schools (see, for example Vardardottir (2013) and Koppensteiner (2012)), belonging to a group of classmates of top students did not benefit those ranked at the bottom of first semester classes. Actually, it goes on the opposite direction, being harmful to their academic performance, *vis-a-vis* students at the top of second semester classes. Concerning high education, our results are similar to those in Contreras, Badua and Adrian (2012), which also found negative peer effects, and go against the conclusions in Paola and Scoppa (2010), Androushchak, Poldin and Yudkevich (2012) and Booij, Leuven and Oosterbeek (2015).

For comparison purposes, this result is quite similar to the financial aid effect on student's performance, as can be seen in table 4.5. In their study for Mexican universities,

---

[7] All empirical exercises in this paper were made by means of R Core Team (2014). This bandwidth was estimated using the package "rdrobust" developed by Calonico, Cattaneo and Titiunik (2015)

[8] As suggested by Doctoral Committee, we estimated this same model with the following modifications: i) undergraduate courses classified by areas of knowledge and ii) subsets, in time, of the panel. The results are presented in appendix, in tables 4.9 and 4.10, and are quite similar to the results presented here. We are grateful for these suggestions.

Canton and Blom (2004) obtained an effect equal to 0.174 on a 10-point scale, equivalent to a 2 % improvement in academic performance for students under such financial aid. For the USA, Curs and Harper (2012) studied the same impact on the first-year GPA of students enrolled at University of Oregon. They found values between 0.12 and 0.16 on a 5-point scale, which is equivalent to a GPA improvement ranging from 2.4% to 3.2%.
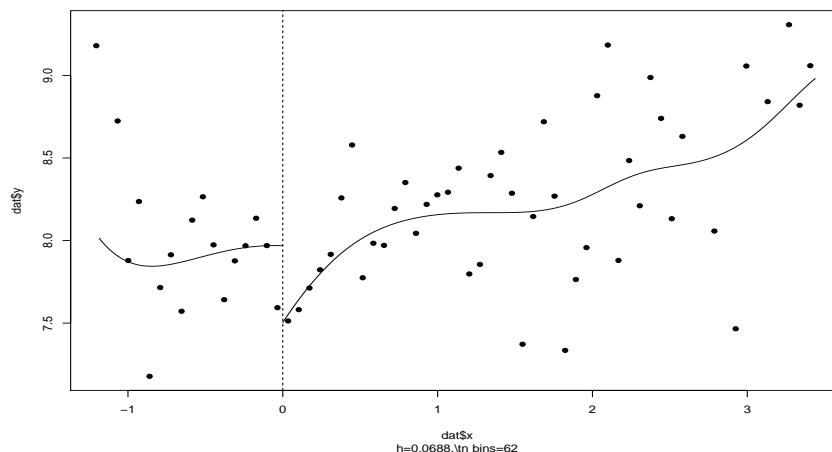
Table 4.5 – Effects of policies in some studies

| Authors | Local | Policy | Estimated value (in %) |
|---|---|---|---|
| In this paper | Brazil | Peer effects | 2% |
| Canton and Blom (2004) | Mexico | Financial aid | 2% |
| Curs and Harper (2012) | United states | Financial aid | 2.4% - 3.2% |

Source: Elaborated by the authors

Now, we turn our attention to the *running variable SAG* and investigate its effect on *IRA*. Table 4.4 shows that the $SAG$ coefficient is not significant, suggesting that it does not affect student's academic performance. Nevertheless, note that variable $Tr * SAG$ brings positive and significant results. This indicates that $SAG$ exerts an influence on the academic performance of students in first semester classes, but not in the second semesters counterparts.

The control variables are all significant and exert negative effects on $IRA$. This means that young male students with a high family income present a lower academic performance than older, female and low income colleagues. The most interesting result here is the fact that high-income students have a lower $IRA$ than those of low income. A possible line of explanation is that students from poorer backgrounds, aspiring to change their social status, could invest more efforts to obtain a higher academic performance. Graphically, the results of model 4.2 are depicted by figure 4.6.

Figure 4.6 – IRA results as a function of standard assignment grade

We verified our results' robustness to the bandwidth choice with a regression sensitivity test. This test consists in reestimating model 4.2 using several bandwidths. After that, we plotted the relation between the bandwidth and the regression discontinuity design's estimates, getting a visually powerful tool to explore the trade-off between bias and precision (JACOB et al., 2012). This is presented in figure 4.7.
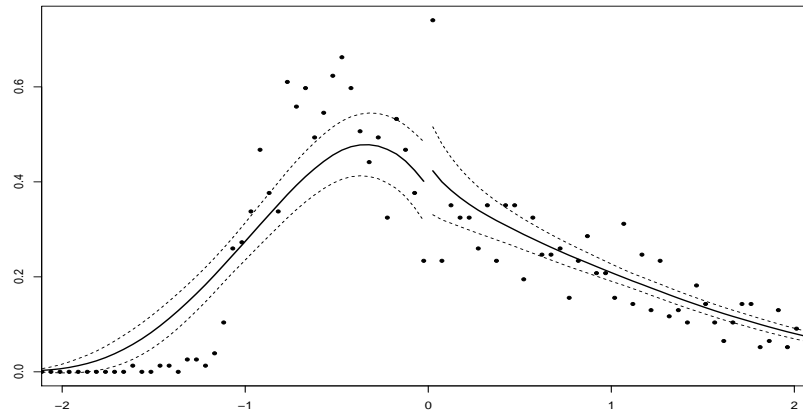
Figure 4.7 – Sensitivity test



This figure shows the treatment effects' response to a bandwidth variation ranging from 0.01 to 2.18. As expected, for small bandwidth values, precision is low and bias is high, with the treatment effect being positive. As the bandwidth assumes higher values, bias decreases and precision increases, with the treatment effect turning negative. For bandwidths larger than 0.5, the treatment effect is virtually unchanged, indicating that ours results are not much sensitive to the bandwidth choice near this value (remember that the optimal bandwidth is about 0.44).

The next step is to test the assignment variable's continuity around the cutoff. A key assumption in the regression discontinuity design approach is that agents are not able to manipulate the assignment variable. If an individual can manipulate it, then he/she can decide whether or not to receive the treatment, so that continuity assumption may not be plausible. To test the this assignment variable's continuity, we use the McCrary (2008) test.

From the McCrary (2008) test, we estimated a discontinuity around 0.10, with $z-value$ and $p-value$ equal to 0.8875 and 0.3748, respectively. In this test, the null hypothesis is that the density is continuous around the cutoff. Given the choosen $p-value$, we cannot reject this null hypothesis, hence our assignment variable is really continuous[9]. The graphical result of this test is shown in figure 4.8.

---

[9] This result is not so surprising. Remember that students do not know the cutoff point, since it is determined exogenously by competition. The only information in students' possesion is the number of vacancies. Therefore, we argue that if students do not know the cutoff, then they are not able to manipulate the assignment variable

Figure 4.8 – McCrary test



Finally, we proceeded with covariates balanced tests. In a regression discontinuity design approach, nothing else, apart from treatment status, is discontinuous in the interval under analysis (JACOB et al., 2012). So, this is equivalent to say that the treatment and control group must be similar. Table 4.6 presents the test results for equality in both means and distributions of the variables *Age*, *Gender* and *Log(income)* for first and second semester classes, around the cutoff.

Table 4.6 – Covariates balanced test

|  | Mean $2^{nd}$ S | Mean $1^{th}$ S | Difference | Statistic | p.value | Density test | p.value |
|---|---|---|---|---|---|---|---|
| Age | 19.5454 | 18.9276 | -0.6178 | -2.3614 | 0.0186 | 0.1047 | 0.1127 |
| Gender | 0.4579 | 0.3957 | -0.0618 | -1.4413 | 0.1501 | 0.0622 | 0.6911 |
| Log(income) | 7.4632 | 7.3417 | -0.1215 | -1.2732 | 0.2037 | 0.0606 | 0.7212 |

Source: Elaborated by the authors

The tests' null hypothesis is that these variables present equal means and distributions around the cutoff. The results shown in table 4.6 suggest that the only variable with a different mean for each side of the cutoff point is *Age*, since the null hypothesis is rejected (p. value = 0.0186). However, the density test shows that all variables have the same distribution in both cutoff sides. Therefore, grounded by the tests performed, we can conclude that our results are valid, since our data obeys the key assumptions in a regression discontinuity design approach.

Until now, we have analyzed our "global model". Now we turn our attention into the "multi-treatment model". The results are shown in table 4.7 [10].

---

[10]   We tried to estimate our "multi-treatment model" with others threshold, representing the bottom and top 5%, 15%, 20% and 25%. However, due problems in groups formations, we estimated only the model with bottom and top 15%, presented in table 4.11, in appendix.

Table 4.7 – SRD estimates for the effect of being in the first semester class on students' IRA — Multi-treatment

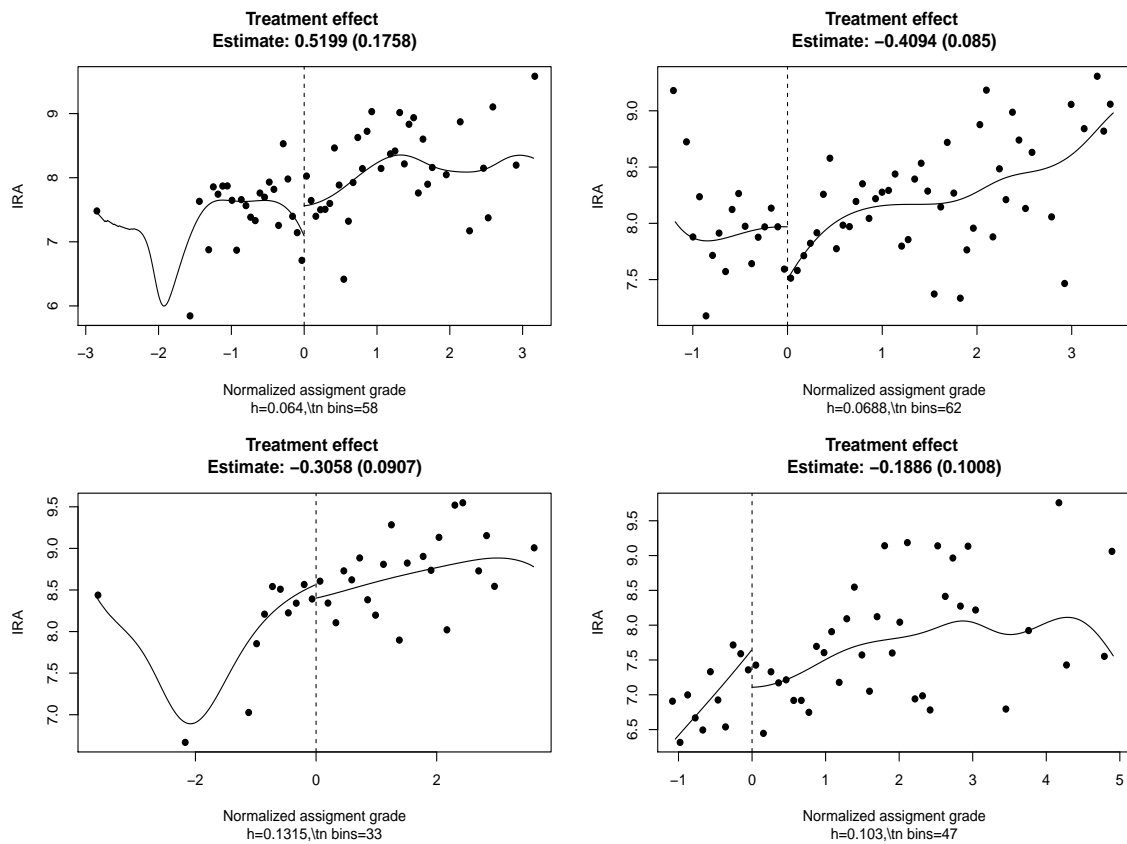| Variable | T1 | T2 | T3 | T4 |
|----------|------|------|------|------|
| Intercept | 7.8021*** | 8.2011*** | 8.9605*** | 10.3635*** |
| | (0.5267) | (0.3399) | (0.4256) | (0.3725) |
| Tr | 0.5199*** | -0.4094*** | -0.3058*** | -0.1886* |
| | (0.1758) | (0.085) | (0.0907) | (0.1008) |
| SAG | -3.9631*** | -0.2487 | 0.5584** | 0.8589** |
| | (1.0568) | (0.4372) | (0.2278) | (0.346) |
| Tr*SAG | 5.6666*** | 1.0617* | -0.2089 | -0.8586* |
| | (1.2926) | (0.5994) | (0.3178) | (0.5015) |
| Age | 0.0699*** | -0.0656*** | -0.0086 | -0.0758*** |
| | (0.0141) | (0.0121) | (0.0114) | (0.0121) |
| Gender | 0.1081 | -0.3346*** | -0.3281*** | -0.2473*** |
| | (0.1271) | (0.052) | (0.0591) | (0.0666) |
| Log(income) | -0.1708*** | 0.0418* | -0.0698* | -0.15*** |
| | (0.0578) | (0.0247) | (0.0375) | (0.0355) |
| Bandwidth | 0.3199 | 0.3441 | 0.6574 | 0.515 |
| $\bar{R}^2$ | 0.4689 | 0.3002 | 0.4131 | 0.525 |
| N | 568 | 1384 | 944 | 1360 |
| Nl | 256 | 800 | 520 | 816 |
| Nr | 312 | 584 | 424 | 544 |

Note: Standard error in parentheses
Note: Signif. codes: $p < 0.01$ "***" $p < 0.05$ "**" $p < 0.1$ "*"
Source: Elaborated by the authors

The results demonstrate a statistically significant treatment effect for all groups. It is mentioning that the treatment is positive at the $T1$ level, while in all other levels it is negative. This is just what we found in our "global model". Thus, we can conclude that there are non-linearities in our "peer effects", corroborating the results of Sacerdote (2001) and Zimmerman (2003).

At the $T1$ level, the results indicate that students at the first semester class of courses with low competition in both classes are benefited *vis-a-vis* those students in a second semester class. This diference is reflected in an $IRA$ 5% higher for the first students group. Regarding the negative treatment effect levels, $T4$ has a magnitude for the effect quite similar to that found in the global model. However, the effects in levels $T2$ and $T3$ are, respectively, 2 and 1.5 times larger, presenting students in first semester classes with $IRA$ 3% and 4% lower than second semester students. A graphical representation of our results for the multi-treatment model is given by figure 4.9.

Now, we are able to discuss the comparison with the control, as well as to analyze the incremental effects. Both measures are presented in table 4.8. We begin analyzing the comparison with the control effect, which we defined as being $T1$. The best way to

Figure 4.9 – $IRA$ results as a function of assignment grade



understand this is by thinking about a first-semester student. $T1$ has both classes with low competition and a positive treatment effect. $T2$ has the first semester class associated to a high competition and, thus, the student's loss is 0.9292, compared to $T1$. In $T3$, the first semester class continues to be of low competition, but that of the second semester is a high competition class now. In this case, the students' loss is less than in $T2$. Finally, in $T4$, both classes are of high competition, and the students' loss is 0.7085, i.e., less than in $T2$ and $T3$.

Now, we are going to analyze the treatment's incremental effect following the same logic as previously. The change from $T1$ to $T2$ is already analyzed. When the treatment of a first-semester student is initially $T2$, and changes to $T3$, his/her class becomes of low competition, and high competition is associated with the second semester class. In the case of a change from $T3$ to $T4$, both classes are of high competition now. In both changes, there is no significant effect.

Therefore, we can conclude that the peer effects are positive when both classes are of low competition, and negative in the other cases. However, note that this negative effect is lower when both classes are of high competition. In addition, the incremental effect is significant only when first and second semester classes present high and low competition, respectively.

Table 4.8 – Incremental and comparison with the control effects

| Statistic | Definition | Value | Std. Dev. | z-value |
|---|---|---|---|---|
| Comparison with control | | | | |
| E(T2 - T1) | $\beta_{T2} - \beta_{T1}$ | -0.9292** | 0.5107 | -1.8194 |
| E(T3 - T1) | $\beta_{T3} - \beta_{T1}$ | -0.8256* | 0.5162 | -1.5992 |
| E(T4 - T1) | $\beta_{T4} - \beta_{T1}$ | -0.7085* | 0.5259 | -1.3470 |
| Incremental | | | | |
| E(T2 - T1) | $\beta_{T2} - \beta_{T1}$ | -0.9292** | 0.5107 | -1.8194 |
| E(T3 - T2) | $\beta_{T3} - \beta_{T2}$ | 0.1035 | 0.4192 | 0.2470 |
| E(T4 - T3) | $\beta_{T4} - \beta_{T3}$ | 0.1171 | 0.4376 | 0.2676 |

Note: Signif. codes: $p < 0.01$ "***" $p < 0.05$ "**" $p < 0.1$ "*"
Source: Elaborated by the authors

## 4.7 FINAL CONSIDERATIONS

This paper sought to apply a well-established methodological approach in a new context: the study of peer effects in a Brazilian university. Due to specificities of the entrance process at the Federal University of Ceará until 2010 — the so called vestibular exam — we are able to use the methodological tools provided by the regression discontinuity design approach, more specifically its sharp version, to estimate peers effects among higher education students.

From this sharp regression discontinuity design approach, we estimated the effect of being in first *versus* second semester classes. We found that, in contrast to what usually happens in studies of peer effects in primary and high schools, being a classmate of high-ability students, i.e. being part of a first semester class, is harmful to a typical student. We obtained a negative effect of about 0.1973, indicating that these students have an academic performance 2% lower than those of second semester classes. For the sake of comparison, this effect is quite similar to what Canton and Blom (2004) and Curs and Harper (2012) obtained in a financial aid context.

Taking advantage that, in our data set, the undergraduate programs have heterogeneous patterns in assignment grades distributions, we classified these programs in four categories, according to the competition in first and second semester classes. After this, we estimated a model capable to assessment a multi treatment. We found, as in Sacerdote (2001) and Zimmerman (2003), that the peer effects present non-linearities. In cases which both classes are of low competition, the peer effects are positive, presenting students of first semester classes with an $IRA$ 5% higher, while in case both classes are of high competition, these students have an $IRA$ 2% lower.

As suggestions for future studies, we believe that the development of a model for the new entrance process by means of SISU could be made. We also believe that replicating our empirical exercise on different data sets, coming from different institutional backgrounds,

might be something worth pursuing to validate our approach. Finally, we believe that this should be done with ENADE's [11]score as an outcome, instead of $IRA's$. It would help us to understand this effect better.

---

[11] The National Survey of Students' Performance (ENADE) is an exam that constitutes the National System of Higher Education Assessment (Sinaes). This test aims to assess students' performance in relation to the syllabus provided in the curriculum guidelines of their undergraduate programs, and the skills and competences in their training ((INEP), 2015).

# REFERENCES

ABDULKADIROĞLU, A.; ANGRIST, J.; PATHAK, P. The elite illusion: Achievement effects at boston and new york exam schools. *Econometrica*, Wiley Online Library, v. 82, n. 1, p. 137–196, 2014.

ABREU, L. C. M. *Mecanismos de seleção Gale-Shapley dinâmicos em universidades Brasileiras: SISU, SISU$_{alpha}$, SISU$_{beta}$*. Dissertação (Dissertação de mestrado) — Universidade Federal do Ceará, 2013.

ANDROUSHCHAK, G. V.; POLDIN, O.; YUDKEVICH, M. *Peer effects in exogenously formed university student groups*. [S.l.], 2012. v. 3, n. 03/EDU/2012.

BOOIJ, A.; LEUVEN, E.; OOSTERBEEK, H. *Ability peer effects in university: Evidence from a randomized experiment*. [S.l.], January 2015.

BUTCHER, K. F.; MCEWAN, P. J.; TAYLOR, C. H. The effects of quantitative skills training on college outcomes and peers. *Economics of Education Review*, Elsevier, v. 29, n. 2, p. 187–199, 2010.

CALONICO, S.; CATTANEO, M. D.; TITIUNIK, R. *rdrobust: Robust Data-Driven Statistical Inference in Regression-Discontinuity Designs*. [S.l.], 2015. R package version 0.80. Disponível em: <http://CRAN.R-project.org/package=rdrobust>.

CANTON, E.; BLOM, A. *Do student loans improve accessibility to higher education and student performance*. [S.l.]: CPB Netherlands Bureau for Economic Policy Analysis The Hague, 2004.

CARD, d.; GIULIANO, L. Can tracking raise the test scores of high-abilitybility minority students? March 2015.

CARRELL, S. E.; FULLERTON, R. L.; WEST, J. E. *Does your cohort matter? Measuring peer effects in college achievement*. [S.l.], 2008.

CARVALHO, J.-R.; MAGNAC, T.; XIONG, Q. *College Choice Allocation Mechanisms: Structural Estimates and Counterfactuals*. [S.l.], June 2014.

CEARÁ FEDERAL UNIVERSITY. *UFC statistical yearbook - A brief edition 2013-2014*. 2014. Disponível em: <http://www.ufc.br/images/_files/a_universidade-/anuario_estatistico/statistical_yearbook_ufc_2014_2013.pdf>.

COLEMAN, J. S. et al. Equality of educational opportunity. *Washington, dc*, p. 1066–5684, 1966.

CONTRERAS, S.; BADUA, F.; ADRIAN, M. Peer effects on undergraduate business student performance. *International Review of Economic Education*, Economics Network, University of Bristol, v. 11, n. 1, p. 57–66, 2012.

CURS, B. R.; HARPER, C. E. Financial aid and first-year collegiate gpa: A regression discontinuity approach. *The Review of Higher Education*, The Johns Hopkins University Press, v. 35, n. 4, p. 627–649, 2012.

EPPLE, D.; ROMANO, R. Peer effects in education: A survey of the theory and evidence. *Handbook of social economics*, Forthcoming, v. 1, n. 11, p. 1053–1163, 2011.

FAN, J.; GIJBELS, I. *Local polynomial modelling and its applications: Monographs on statistics and applied probability 66*. Taylor and Francis, 1996. (Chapman & Hall/CRC Monographs on Statistics & Applied Probability). ISBN 9780412983214. Disponível em: <https://books.google.com.br/books?id=BM1ckQKCXP8C>.

GAWANDE, A. *BrainyQuote.com.* n.d. <http://www.brainyquote.com/quotes/quotes/a/atulgawand527239.html>. Accessed: May 22, 2015.

GELMAN, A.; IMBENS, G. *Why high-order polynomials should not be used in regression discontinuity designs*. [S.l.], 2014.

IMBENS, G.; KALYANARAMAN, K. Optimal bandwidth choice for the regression discontinuity estimator. *The Review of Economic Studies*, Oxford University Press, v. 79, n. 3, p. 933–959, 2011.

IMBENS, G. W.; LEMIEUX, T. Regression discontinuity designs: A guide to practice. *Journal of econometrics*, Elsevier, v. 142, n. 2, p. 615–635, 2008.

(INEP), I. de Estudos e P. E. A. T. *Censo da educação superior 2013*. 2014. Disponível em: <http://download.inep.gov.br/educacao_superior/censo_superior/apresentacao-/2014/coletiva_censo_superior_2013.pdf>.

(INEP), I. de Estudos e P. E. A. T. *Exame Nacional de Desempenho de Estudantes - ENADE*. May 2015. Disponível em: <http://portal.inep.gov.br/enade>.

JACOB, R. T. et al. *A practical guide to regression discontinuity*. [S.l.]: MDRC, 2012.

KLAAUW, W. Van der. Estimating the effect of financial aid offers on college enrollment: A regression discontinuity approach. *International Economic Review*, Wiley Online Library, v. 43, n. 4, p. 1249–1287, 2002.

KOPPENSTEINER, M. F. *Class assignment and peer group effects: Evidence from Brazilian primary schools*. [S.l.], 2012.

LEE, D. S.; LEMIEUX, T. *Regression discontinuity designs in economics*. [S.l.], 2009.

LEE, M.-J. *Micro-econometrics for policy, program, and treatment effects.* [S.l.]: Oxford University Press Oxford, 2005.

LEEDS, D. M.; DESJARDINS, S. L. The effect of merit aid on enrollment: A regression discontinuity analysis of iowa's national scholars award. *Research in Higher Education*, Springer, p. 1–25, 2014.

MANSKI, C. F. Identification of endogenous social effects: The reflection problem. *The review of economic studies*, Oxford University Press, v. 60, n. 3, p. 531–542, 1993.

MCCRARY, J. Manipulation of the running variable in the regression discontinuity design: A density test. *Journal of Econometrics*, Elsevier, v. 142, n. 2, p. 698–714, 2008.

MCEWAN, P. J.; SODERBERG, K. A. Roommate effects on grades: Evidence from first-year housing assignments. *Research in Higher Education*, Springer, v. 47, n. 3, p. 347–370, 2006.

MEALLI, F.; RAMPICHINI, C. Evaluating the effects of university grants by using regression discontinuity designs. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, Wiley Online Library, v. 175, n. 3, p. 775–798, 2012.

MOSS, B. G.; YEATON, W. H. Shaping policies related to developmental education: An evaluation using the regression-discontinuity design. *Educational Evaluation and Policy Analysis*, Sage Publications, v. 28, n. 3, p. 215–229, 2006.

PAOLA, M. D.; SCOPPA, V. Peer group effects on the academic performance of italian students. *Applied Economics*, Taylor & Francis, v. 42, n. 17, p. 2203–2215, 2010.

R Core Team. *R: A Language and Environment for Statistical Computing.* Vienna, Austria, 2014. Disponível em: <http://www.R-project.org/>.

SACERDOTE, B. Peer effects with random assignment: Results for dartmouth roommates. *The Quarterly Journal of Economics*, Oxford University Press (OUP), v. 116, n. 2, p. 681–704, May 2001. ISSN 1531-4650. Disponível em: <http://dx.doi.org/10.1162-/00335530151144131>.

SACERDOTE, B. Peer effects in education: How might they work, how big are they and how much do we know thus far? *Handbook of the Economics of Education*, Elsevier, v. 3, p. 249–277, 2011.

SCHÖER, V.; SHEPHERD, D. *Compulsory tutorial programmes and performance in undergraduate microeconomics: A regression discontinuity design.* [S.l.], September 2013.

SILVA, J. E. M. d. *Vestibular na UFC: Implicações para o ensino médio — Um estudo de caso a partir das mudanças realizadas de 1978 a 2004.* Dissertação (Dissertação de mestrado) — Universidade Federal do Ceará, 2007.

VARDARDOTTIR, A. Peer effects and academic achievement: A regression discontinuity approach. *Economics of Education Review*, Elsevier, v. 36, p. 108–121, 2013.

WINSTON, G.; ZIMMERMAN, D. Peer effects in higher education. In: *College choices: The economics of where to go, when to go, and how to pay for it*. [S.l.]: University of Chicago Press, 2004. p. 395–424.

ZIMMERMAN, D. J. Peer effects in academic outcomes: Evidence from a natural experiment. *Review of Economics and Statistics*, MIT Press, v. 85, n. 1, p. 9–23, 2003.

# APPENDIX

Table 4.9 – SRD estimates for the effect of being in the first semester class on students'
IRA — Classified by areas of knowledge

| Variable | CSB | CSA | CET | CLH |
|---|---|---|---|---|
| Intercept | 8.2569*** | 9.7492*** | 7.1337*** | 9.6580*** |
| | (0.2345) | (0.3197) | (0.3477) | (0.3517) |
| Tr | -0.1765*** | -0.5487*** | -0.1103 | -0.2615*** |
| | (0.0527) | (0.0785) | (0.1245) | (0.1009) |
| SAG | -0.2522*** | 0.9340*** | 0.5565 | 0.5754* |
| | (0.0862) | (0.2332) | (0.3610) | (0.3141) |
| Tr*SAG | 0.6485*** | -0.0437 | 0.3345 | -0.5638 |
| | (0.1342) | (0.3427) | (0.5172) | (0.4393) |
| Age | -0.0592*** | -0.0368*** | -0.0013 | 0.0007 |
| | (0.0074) | (0.0115) | (0.0094) | (0.0108) |
| Gender | -0.2113*** | -0.3220*** | -0.1550* | -0.1691*** |
| | (0.0320) | (0.0474) | (0.0946) | (0.0633) |
| Log(income) | 0.1101*** | -0.0712*** | -0.0894*** | -0.1047*** |
| | (0.0199) | (0.0241) | (0.0328) | (0.0358) |
| Bandwidth | 0.9530 | 0.5733 | 0.6136 | 0.5568 |
| $\bar{R}^2$ | 0.2417 | 0.5132 | 0.3872 | 0.2113 |
| N | 2464 | 1928 | 1176 | 1168 |
| Nl | 1576 | 1152 | 648 | 704 |
| Nr | 888 | 776 | 528 | 464 |

Note: Standard error in parentheses

Note: Signif. codes: $p < 0.01$ "***" $p < 0.05$ "**" $p < 0.1$ "*"

Source: Elaborated by the authors

- •CSB - Health and biological science

- •CSA - Applied social sciences

- •CET - Exact and earth sciences

- •CLH - Languages and human sciences

Table 4.10 – SRD estimates for the effect of being in the first semester class on students'
IRA — Divided by semesters

| Variable | G1 | G2 | G3 | G4 |
|---|---|---|---|---|
| Intercept | 9.6391*** | 9.0369*** | 8.7450*** | 7.8335*** |
|  | (0.2382) | (0.2244) | (0.2343) | (0.3687) |
| Tr | -0.3535*** | -0.2780*** | -0.2134*** | -0.2119** |
|  | (0.0681) | (0.0605) | (0.0591) | (0.0962) |
| SAG | 0.3031*** | 0.1774 | -0.0396 | 0.4021 |
|  | (0.0892) | (0.1639) | (0.2342) | (0.3262) |
| Tr*SAG | 0.1410 | 0.3937* | 0.7501** | 0.2652 |
|  | (0.1177) | (0.2391) | (0.3207) | (0.4595) |
| Age | -0.0220*** | -0.0175*** | -0.0203*** | -0.0164 |
|  | (0.0068) | (0.0065) | (0.0068) | (0.0108) |
| Gender | -0.2816*** | -0.3040*** | -0.2869*** | -0.1758*** |
|  | (0.0416) | (0.0390) | (0.0383) | (0.0626) |
| Log(income) | -0.0827*** | -0.0521*** | -0.0450** | -0.0201 |
|  | (0.0210) | (0.0192) | (0.0196) | (0.0316) |
| Bandwidth | 1.7633 | 0.6188 | 0.4540 | 0.5138 |
| $\bar{R}^2$ | 0.4254 | 0.4452 | 0.4526 | 0.4718 |
| N | 2866 | 3152 | 3354 | 1304 |
| Nl | 1596 | 1868 | 1872 | 752 |
| Nr | 1270 | 1284 | 1482 | 552 |

Note: Standard error in parentheses
Note: Signif. codes: $p < 0.01$ "***" $p < 0.05$ "**" $p < 0.1$ "*"
Source: Elaborated by the authors

- G1 - Only semesters 1 and 2

- G2 - From semester 1 up to semester 4

- G3 - From semester 1 up to semester 6

- G4 - Only semesters 7 and 8

Table 4.11 – SRD estimates for the effect of being in the first semester class on students'
IRA — Multi-treatment - 15%

| Variable | T1 | T2 | T3 | T4 |
|---|---|---|---|---|
| Intercept | 4.3592*** | 491.038*** | 7.6915*** | 11.0108*** |
| | (0.4381) | (132.1605) | (0.2484) | (0.2735) |
| Tr | 0.8284*** | -482.9708*** | -0.2114*** | -0.3252*** |
| | (0.2309) | (131.6681) | (0.0722) | (0.0645) |
| SAG | -2.944*** | 1899.8591*** | 0.4476* | 0.4388** |
| | (0.5825) | (518.029) | (0.2626) | (0.1859) |
| Tr*SAG | 5.3789*** | -1898.0506*** | 0.4367 | -0.1867 |
| | (0.7463) | (517.7328) | (0.3504) | (0.2845) |
| Age | 0.0079 | -0.0926 | 0.0153** | -0.0778*** |
| | (0.0094) | (0.0805) | (0.0076) | (0.0085) |
| Gender | -0.2988** | 0.0886 | -0.1168** | -0.3395*** |
| | (0.1281) | (0.1689) | (0.0483) | (0.0412) |
| Log(income) | 0.2772*** | -0.4165 | -0.0163 | -0.1138*** |
| | (0.0535) | (0.2616) | (0.0216) | (0.0228) |
| Bandwidth | 0.5789 | 0.3571 | 0.4973 | 0.5722 |
| $\bar{R}^2$ | 0.4551 | 0.5839 | 0.413 | 0.4825 |
| N | 376 | 48 | 2352 | 2672 |
| Nl | 208 | 16 | 1288 | 1688 |
| Nr | 168 | 32 | 1064 | 984 |

Note: Standard error in parentheses
Note: Signif. codes: $p < 0.01$ "***" $p < 0.05$ "**" $p < 0.1$ "*"
Source: Elaborated by the authors

- T1 - $SAG < -0.73$ and $SAG \in [0,1.18]$

- T2 - $SAG \in [-0.73,0)$ and $SAG \in [0,1.18]$

- T3 - $SAG < -0.73$ and $SAG > 1.18$

- T4 - $SAG \in [-0.73,0)$ and $SAG > 1.18$

# APPENDICES

<h1 style="text-align:center">SPRINGER LICENSE<br>TERMS AND CONDITIONS</h1>

Apr 14, 2016

This is a License Agreement between Diego André ("You") and Springer ("Springer") provided by Copyright Clearance Center ("CCC"). The license consists of your order details, the terms and conditions provided by Springer, and the payment terms and conditions.

**All payments must be made in full to CCC. For payment instructions, please see information listed at the bottom of this form.**

| | |
|---|---|
| License Number | 3847830546622 |
| License date | Apr 14, 2016 |
| Licensed content publisher | Springer |
| Licensed content publication | Water Resources Management |
| Licensed content title | Spatial Determinants of Urban Residential Water Demand in Fortaleza, Brazil |
| Licensed content author | Diego de Maria André |
| Licensed content date | Jan 1, 2014 |
| Volume number | 28 |
| Issue number | 9 |
| Type of Use | Thesis/Dissertation |
| Portion | Full text |
| Number of copies | 1 |
| Author of this Springer article | Yes and you are the sole author of the new work |
| Order reference number | None |
| Title of your thesis / dissertation | THREE ESSAYS ON APPLIED MICROECONOMETRICS WITH SPATIAL EFFECTS |
| Expected completion date | May 2016 |
| Estimated size(pages) | 84 |
| Total | 0.00 USD |
| Terms and Conditions | |

Introduction

The publisher for this copyrighted material is Springer. By clicking "accept" in connection with completing this licensing transaction, you agree that the following terms and conditions apply to this transaction (along with the Billing and Payment terms and conditions established by Copyright Clearance Center, Inc. ("CCC"), at the time that you opened your Rightslink account and that are available at any time at http://myaccount.copyright.com).

Limited License

With reference to your request to reuse material on which Springer controls the copyright, permission is granted for the use indicated in your enquiry under the following conditions:

- Licenses are for one-time use only with a maximum distribution equal to the number stated in your request.

- Springer material represents original material which does not carry references to other sources. If the material in question appears with a credit to another source, this permission is

not valid and authorization has to be obtained from the original copyright holder.
- This permission
• is non-exclusive
• is only valid if no personal rights, trademarks, or competitive products are infringed.
• explicitly excludes the right for derivatives.
- Springer does not supply original artwork or content.
- According to the format which you have selected, the following conditions apply accordingly:
• **Print and Electronic:** This License include use in electronic form provided it is password protected, on intranet, or CD-Rom/DVD or E-book/E-journal. It may not be republished in electronic open access.
• **Print:** This License excludes use in electronic form.
• **Electronic:** This License only pertains to use in electronic form provided it is password protected, on intranet, or CD-Rom/DVD or E-book/E-journal. It may not be republished in electronic open access.
For any electronic use not mentioned, please contact Springer at permissions.springer@spi-global.com.
- Although Springer controls the copyright to the material and is entitled to negotiate on rights, this license is only valid subject to courtesy information to the author (address is given in the article/chapter).
- If you are an STM Signatory or your work will be published by an STM Signatory and you are requesting to reuse figures/tables/illustrations or single text extracts, permission is granted according to STM Permissions Guidelines: http://www.stm-assoc.org/permissions-guidelines/
For any electronic use not mentioned in the Guidelines, please contact Springer at permissions.springer@spi-global.com. If you request to reuse more content than stipulated in the STM Permissions Guidelines, you will be charged a permission fee for the excess content.
Permission is valid upon payment of the fee as indicated in the licensing process. If permission is granted free of charge on this occasion, that does not prejudice any rights we might have to charge for reproduction of our copyrighted material in the future.
-If your request is for reuse in a Thesis, permission is granted free of charge under the following conditions:
This license is valid for one-time use only for the purpose of defending your thesis and with a maximum of 100 extra copies in paper. If the thesis is going to be published, permission needs to be reobtained.
- includes use in an electronic form, provided it is an author-created version of the thesis on his/her own website and his/her university's repository, including UMI (according to the definition on the Sherpa website: http://www.sherpa.ac.uk/romeo/);
- is subject to courtesy information to the co-author or corresponding author.
Geographic Rights: Scope
Licenses may be exercised anywhere in the world.
Altering/Modifying Material: Not Permitted
Figures, tables, and illustrations may be altered minimally to serve your work. You may not alter or modify text in any manner. Abbreviations, additions, deletions and/or any other alterations shall be made only with prior written authorization of the author(s).
Reservation of Rights
Springer reserves all rights not specifically granted in the combination of (i) the license details provided by you and accepted in the course of this licensing transaction and (ii) these terms and conditions and (iii) CCC's Billing and Payment terms and conditions.
License Contingent on Payment
While you may exercise the rights licensed immediately upon issuance of the license at the end of the licensing process for the transaction, provided that you have disclosed complete and accurate details of your proposed use, no license is finally effective unless and until full

payment is received from you (either by Springer or by CCC) as provided in CCC's Billing and Payment terms and conditions. If full payment is not received by the date due, then any license preliminarily granted shall be deemed automatically revoked and shall be void as if never granted. Further, in the event that you breach any of these terms and conditions or any of CCC's Billing and Payment terms and conditions, the license is automatically revoked and shall be void as if never granted. Use of materials as described in a revoked license, as well as any use of the materials beyond the scope of an unrevoked license, may constitute copyright infringement and Springer reserves the right to take any and all action to protect its copyright in the materials.

Copyright Notice: Disclaimer

You must include the following copyright and permission notice in connection with any reproduction of the licensed material:

"Springer book/journal title, chapter/article title, volume, year of publication, page, name(s) of author(s), (original copyright notice as given in the publication in which the material was originally published) "With permission of Springer"

In case of use of a graph or illustration, the caption of the graph or illustration must be included, as it is indicated in the original publication.

Warranties: None

Springer makes no representations or warranties with respect to the licensed material and adopts on its own behalf the limitations and disclaimers established by CCC on its behalf in its Billing and Payment terms and conditions for this licensing transaction.

Indemnity

You hereby indemnify and agree to hold harmless Springer and CCC, and their respective officers, directors, employees and agents, from and against any and all claims arising out of your use of the licensed material other than as specifically authorized pursuant to this license.

No Transfer of License

This license is personal to you and may not be sublicensed, assigned, or transferred by you without Springer's written permission.

No Amendment Except in Writing

This license may not be amended except in a writing signed by both parties (or, in the case of Springer, by CCC on Springer's behalf).

Objection to Contrary Terms

Springer hereby objects to any terms contained in any purchase order, acknowledgment, check endorsement or other writing prepared by you, which terms are inconsistent with these terms and conditions or CCC's Billing and Payment terms and conditions. These terms and conditions, together with CCC's Billing and Payment terms and conditions (which are incorporated herein), comprise the entire agreement between you and Springer (and CCC) concerning this licensing transaction. In the event of any conflict between your obligations established by these terms and conditions and those established by CCC's Billing and Payment terms and conditions, these terms and conditions shall control.

Jurisdiction

All disputes that may arise in connection with this present License, or the breach thereof, shall be settled exclusively by arbitration, to be held in the Federal Republic of Germany, in accordance with German law.

**Other conditions:**

V 12AUG2015

**Questions? customercare@copyright.com or +1-855-239-3415 (toll free in the US) or +1-978-646-2777.**

# DECLARAÇÃO

Atesto, para os devidos fins, que os capítulos intitulados "Spatial willingness to pay for a first-order stochastic reduction on the risk of robbery" e "Peer effects and academic performance in higher education – a regression discontinuity design approach", da tese de doutorado de autoria de Diego de Maria André, foram devidamente revisados. O material está em consonância com a gramática normativa da língua inglesa.

Fortaleza, 03 de maio de 2016.

*Ananda Badaró de A. Prata.*

**Ananda Badaró de Athayde Prata**
Tradutora/revisora de texto
Graduada em Letras Português – Inglês – Licenciatura Plena pela Universidade Federal do Ceará (UFC)
Especialista em Tradução pelo Programa de Pós-Graduação em Linguística Aplicada da Universidade Estadual do Ceará (PosLa/Uece)