

## Are Public Libraries Improving Quality of Education? When the Provision of Public Goods is not Enough<sup>1</sup>

### *¿Están las bibliotecas públicas mejorando la calidad de la educación? Cuando la provisión de bienes públicos no es suficiente*

Paul Rodríguez-Lesmes<sup>2</sup>

José D. Trujillo<sup>3</sup>

Daniel Valderrama<sup>4</sup>

DOI: 10.13043/DYS.74.5

### Abstract

We analyze the relation between public, education-related infrastructure and the quality of education in schools using a case-study of the construction and implementation of two large public libraries in low-income areas in Bogotá, Colombia. We assess the impact of these libraries on quality of education by comparing results in national test scores (Saber 11<sup>o</sup>) for schools close and far from these libraries before (2000–2002) and after (2003–2008) the libraries'

---

1 This project was funded by the Icfes under a research grant for graduate students in 2010. We express our gratitude to Hugo Ñopo, Andrés García, Ali Sharman and valuable comments from Icfes seminars. All findings, interpretations, and conclusions expressed in this paper are entirely those of the authors. They do not necessarily represent the views neither of Icfes nor other institutions that authors are part of.

2 PhD-Student, Department of Economics, University College London. E-mail: p.lesmes.11@ucl.ac.uk.

3 Consultant, DANE. E-mail: jdtrujillos@dane.gov.co.

4 Consultant, World Bank. E-mail: dvalderramagonza@worldbank.org.

Este artículo fue recibido el 14 de abril de 2014; revisado el 9 de julio de 2014 y, finalmente aceptado el 8 de octubre de 2014.

opening. We find non-statistically different from zero differences that could be attributed to the libraries' implementation. We also introduce Oaxaca-Blinder decomposition on Difference-in-Differences (DiD) estimates in order to assess if variation of traditional determinants of test scores for mathematics, verbal and science explain the result estimates. These results are robust to alternative specifications of DiD, a synthetic control approach and an alternative measure of distance.

*Key words:* Libraries, quality of education, school quality, public good provision.

*JEL classification:* D62, I21, H52.

## Resumen

Analizamos la relación entre infraestructura pública, orientada hacia la educación, y la calidad de la educación de los colegios tomando como caso de estudio la construcción e implementación de dos bibliotecas públicas de gran escala en áreas de ingresos bajos en Bogotá, Colombia. Para verificar el impacto de esas bibliotecas comparamos los resultados de los colegios cercanos y lejanos a ellas en las pruebas nacionales de educación media-secundaria Saber 11°, antes (2000-2002) y después (2003-2008) de la apertura de estas. No encontramos evidencia estadísticamente significativa de que las bibliotecas hayan generado una diferencia en los resultados. Para poder determinar si el impacto en los resultados de matemáticas, lenguaje y ciencias se explica por variaciones en determinantes tradicionales, desarrollamos una descomposición Oaxaca-Blinder del estimador de Diferencia en Diferencias (DiD). Estos resultados son robustos a diferentes especificaciones del DiD, a la utilización de la técnica de control sintético o a medidas alternativas de la distancia.

*Palabras clave:* bibliotecas, calidad de la educación, calidad de los colegios, provisión de bienes públicos.

*Clasificación JEL:* D62, I21, H52.

## Introduction

Facilitating public access to information, the traditional primary function of libraries, is being challenged by the information revolution. However, public libraries serve multiple functions beyond their role in disseminating materials. A big movement of public library construction undertaken in the developing world reflects these functions by emphasizing libraries as the center of social transformation in deprived slums, providing the general population, especially the less well-off, with access to meeting spaces, cultural activities, technology, and information services, among others. For example, impressive (and expensive), massive public libraries were constructed in the most impoverished areas of Medellín (Colombia), in zones with high criminal rates, and Bogotá (Colombia). These libraries are not only places where you can find books or magazines for free, but also places offering a wide range of services which are intended to motivate the general public towards culture and education and, ultimately, to change living conditions of the people.

The goal of this article is to establish the impact on the quality of education of the 2001 construction of two of these massive libraries (from here on mega-libraries) in the city of Bogotá (Colombia). Even public schools provide services to a selected group of students, thus they can be considered as private asset in a sense. Public libraries, however, are available to students from different schools. Thus, this study will tell us something about the possible effect of truly public, education-related infrastructure on quality of education. It is also possible to assess latent complementarities between public (libraries) and private (schools) educational services in enhancing quality education by estimating the effect of libraries on the returns that certain school characteristics have on education. In other words, the paper studies how public libraries affect the quality of education and to what extent this could be through the enhancement of services provided by schools.

This paper contributes a new perspective to the literature on the determinants of quality of education. This literature is generally limited to the use of private characteristics from the school and from the family to explain differences in student performance. By widening the perspective of determinants beyond the walls of the school and the house, this paper contributes to the education literature, looking towards public goods that are around the schools and which could be used to enhance the impact of schools' inputs. At the same

time, considering that the main objective of libraries is not their direct influence on quality of education in schools, this paper contributes to the urban economics literature by analyzing the existence of externalities and complementarities between this kind of public infrastructure and schools or households near to the libraries.

The causal effect of access to public libraries on student academic performance is assessed using a Difference-in-Differences (DiD) methodology, combined with propensity score matching as a robustness test of the results. The procedure takes advantage of the spatial location of the libraries with the first, *El Tunal* constructed in the grounds of a public park, and the second, *El Tintal*, in an old garbage processing plant. We compare the average results on standardized test scores at the end of secondary level studies of schools (Saber 11<sup>o</sup>) close to the libraries and those far from them from 2000 to 2008, that is, before and after the libraries' opening. This concept is implemented under both parametric and nonparametric specifications of the relationship between distance to the library and test scores. We also implement a Oaxaca-Blinder decomposition of the impact of the program on the quality of education to explore the possible improvement via the variation of traditional inputs of education quality.

Given our specification, we are considering both the direct and indirect impacts that the libraries could have on student performance. Direct impact might come from the possibility that students living close to libraries access library services and programs independently or that nearby schools deliberately take advantage of the library for their own activities. Indirect effects might come from the impact of the renovation of the public infrastructure on the area which could improve crime perceptions, the general mood of the population, or other neighborhood effects. Due to the lack of information on students' actual residences or on specific school programs which take advantage of the libraries, we cannot assess these channels separately.

Our main results show that while the relationship shows the expected positive sign, results are not statistically significant. This either tells us that the libraries are not fully exploited by schools or that the possible gains are concentrated among particular types of individuals. This opens the question of how aligned incentives are to foster cooperation between schools and public libraries in order to improve the quality of education. Perhaps it is not enough

to construct beautiful and well-equipped libraries that are near to schools; a second generation of policies might be required to enhance the coordination between these libraries with the current educational environment of neighborhood schools and households.

The remainder of this paper is organized as follows: Section I discusses the theoretical links between libraries and quality of education. Next, Section II describes the program and its context, Section III presents data on quality of education and other controls. Section IV discusses the identification strategy and decomposition of the effect, Section V presents the results and Section VI concludes.

## I. Libraries and Academic Performance

Vegas and Petrow (2008) classify determinants of education into demand-based and supply based components. Both groups include tangible and intangible inputs defined by students' access to private facilities or their environments. For instance, on the demand side, important inputs include an environment, defined by parental characteristics, that promotes study (Fertig and Schmidt, 2002; World Bank, 2005) and the availability of educational resources in the household, like books or well used internet (Blomeyer, Coneus, Laucht and Pfeiffer, 2009; Gamboa, Rodríguez-Acosta and García-Suaza, 2010; Murnane, Maynard and Ohls, 1981). On the supply side, libraries are included as physical infrastructure along with other, intangible, inputs which are generally considered more important, such as educational policy which incentivizes competence in schools and teacher quality (Hanushek and Woßmann, 2007).

Focusing on the impact of libraries on education beyond the 'infrastructure' component of schools, Lance (1994) in a largely descriptive study of improvements on school performance that are associated with libraries in Colorado, shows a relationship between the availability of libraries and specific skills such as reading, writing and critical thinking. Similar relationships are discussed in Lance and other's further research of libraries in the United States (Lance, 1994; Lance, Rodney and Hamilton-Pennell, 2000; Rodney, Lance, Hamilton-Pennell and Center, 2002) and the United Kingdom (Williams, Wavell and Coles, 2001). Lonsdale (2003) provides a review of studies linking libraries to educational outcomes, such as Smith (2001) which argues that libraries improve

by 4% school performance. However, this literature does not involve a causal analysis; it rests on correlation and qualitative analysis.

In terms of proper causal analyses, few in the literature analyze libraries themselves. The most relevant literature analyzes the impact of programs which make learning materials more available in schools on educational outcomes. These learning materials, a traditional part of library services, are: textbooks (Glewwe, Kremer and Moulin, 2009), flipcharts (Glewwe, Kremer, Moulin and Zitzewitz, 2004) and computers in schools (Barrera-Osorio and Linden, 2009). Across programs, each with its own particularities, no authors find impact of the respective learning material on the quality of education received by the average student.<sup>5</sup> However, these evaluations do not consider the joint effect derived from the interaction of these learning materials, an effect that could be captured in an analysis of public libraries given that these institutions provide learning materials simultaneously.

Borkum, He and Linden (2013) is the only study found that explores the role of libraries on educational outcomes. In an evaluation of an educational program in Bangalore, India that provides high quality libraries to public primary schools, the authors find no impact of school libraries on scores of different subjects and on dropout rates. Given that this study does not consider public libraries and, most importantly, the type of public libraries that we are considering (mega-libraries), the present study is the first that presents evidence on causality between public mega-libraries<sup>6</sup> on educational outcomes within impoverished areas in a developing country.

We propose that the production function of education quality for school  $i$ ,  $Y_i$ , in urban areas include not only the demand characteristics that it faces,  $X_{1i}$ , and private supply (in this case, schools) characteristics,  $X_{2i}$ , but also the benefit from public, education related facilities  $Z_i$  (equation 1). This additional input acts as a complement to the education provided by schools. Assuming that these institutions do have a positive impact on the skills related to test-scores of their users, the relationship between  $Z$  and  $Y$  might vary according to the

---

5 In an evaluation of the impact of textbooks on student achievement, Glewwe et al. (2009) finds a localized positive effect on those students who already had relatively high achievement

6 Mega-libraries are not just large buildings full of learning materials but represent a catalyst for redevelopment of urban zones and repositories of new public spaces.

interaction between both the demand and supply elements related to using the public, education-related facilities. In other words, the impact of public, education-related facilities on quality of education depends on the degree to which both families directly use them and schools facilitate their use.<sup>7</sup> Let us consider two examples: first, for school managers who obtain more benefits for promoting activities related to a particular public facility than others,  $Z$  might be larger; second, families living far from public facilities are less likely to benefit from them due to credit or time constraints, which will be reflected in a lower value of  $Z$  than for those who live close by.

$$Y_i = f(X_{1,i}, X_{2,i}, Z_i(X_{1,i}, X_{2,i})) \quad (1)$$

Our data is limited by only one kind of public, education-related facility (the mega-libraries) to calculate  $Z$  and as we don't have information about relation between schools-households and libraries, so we cannot disentangle the relationship between  $Z$  and  $Y$  at the level of detail just explained. Given these data restrictions, our data will use the proximity of schools to the libraries as a proxy of  $Z$ .

In order to link the relation between the schools and libraries we use as measure of intensity the distance between both. That is, we will identify the difference  $\delta$  of being close rather than far to the public facility based on assigning a discrete value of  $T = 1$ , if a school is within a close range of a library and  $T = 0$  if the school is outside of this range. Our main assumption is that if a school is far enough away from the public facility, their students do not receive any benefit from it ( $Z = 0$ , as shown in the Equation 2).

$$\begin{aligned} \delta &= E[Y_i | T = 1, X] - E[Y_i | T = 0, X] \\ &= f(X_{1,i}(T = 1), X_{2,i}(T = 1), Z(X_{1,i}(T = 1), X_{2,i}(T = 1))) \\ &\quad - f(X_{1,i}(T = 0), X_{2,i}(T = 0), Z(X_{1,i}(T = 0), X_{2,i}(T = 0))) \quad (2) \\ &= f(X_{1,i}(T = 1), X_{2,i}(T = 1), Z(X_{1,i}(T = 1), X_{2,i}(T = 1))) \\ &\quad - f(X_{1,i}(T = 0), X_{2,i}(T = 0), 0) \end{aligned}$$

---

7 Positive returns to higher levels of school quality based on facility use in Colombia are expected for families (Gamboa and Rodríguez-Lesmes, 2014). However, it is not clear that all schools have the same incentives (Gaviria and Barrientos, 2001).

## II. BibloRed Program and Colombian Schools

*BibloRed* is a program which Bogotá's local administration designed in 1998 and operationalized by the end of 2001. The idea was to allow the general population to get access to information services and reading and writing resources. However, the program also seeks to foment cultural growth and promote research. In the first stage, the operation started with 3 major libraries (*El Tunal*, *El Tintal* and the *Virgilio Barco*), 15 minors libraries and 1 *biblioteca*; almost ten years later another major library started operations (*Julio Mario Santo Domingo*). Each major library has an area of around 10,000 square meters, 150,000 volumes and 600 reader seats (Tolosa, 2012). Information services not only include books and magazines, but also children's rooms with specialized staff, programs for babies and their parents, activities for teens, workshops in literature, puppets, etc. The intention is to attract the public with these activities while integrating education into them. One of the main projects occurs over holidays, when *BibloRed* implements *Bibliovacaciones*, a program with the activities mentioned plus cost-free art, history and literature exhibitions such as theatre plays and films. In this context, it is evident that these libraries have many activities which enhance the quality of life, particularly through their integration of culture; thus, the possible effect on the educational performance of children and young people is just one of the multiple benefits that libraries bring to society.

Since it is not possible to have information on which of the test-takers actually use the libraries, we propose to use the distance of libraries to their schools as an alternative indicator for treatment status. As discussed in the previous section, this rests on the assumption that the use of libraries is likely to be higher for those living closer than for those who live far, supported by travel costs to libraries incurred by the latter which reduce students' incentives to visit them frequently. According to Table A1.1, 77% of students in Bogotá live less than 20 minutes from the school they attend. As a result, it is a fair assumption that distance from school to the library approximates the distance from the library to students' residence and, therefore, the likelihood that they live in an environment affected by libraries.

The Euclidean distance between the school and the local library is shown in Figure 1. We calculate it based on the information on the spatial location of



Figure 1. Libraries and Treatment Status Allocation: Euclidean Distance



Source: Own calculation based on C-600.

each school as specified by Bogotá's Department of Education. Alternatively, we use road-based distances as shown in Figure 2.<sup>8</sup> Figure A2.1 presents the

8 These calculations were made using ESRI ArcMap 10.2 Closest Facility Analysis. The road network was obtained from Open Street Map project (OSM).

Figure 2. Libraries and Treatment Status Allocation: Road Distance



Source: Own calculation based on C-600.

link between both distances. As expected, the road-based distances all fall above the blue line corresponding to the 45-degree line. The black dotted line is the predicted linear relationship between both measures, which captures

up 80% of total variation. As a robustness check, the main estimators are repeated using the fitted distance.<sup>9</sup>

*El Tintal* and *El Tunal* libraries are located in middle-low income zones, where most of the students attend nearby schools. Schools near to *Virgilio Barco* and *Julio Mario Santodomingo* are populated by, on average, wealthier families which are more likely to live far from school and use private transport for the daily commuting. If we include the last two libraries, our approximation of taking the distance between the library and the school to represent the treatment status will not be accurate. As a result, we decided to include only *El Tintal* and *El Tunal* libraries in this analysis.

In Colombia, schools can be classified according to four important characteristics that are closely related with the quality of education in the literature. These characteristics are: whether the school is managed by the government, the proportion of females to males attending the school, the start of the academic year and the length of the school day. In regards to the first characteristic, most of the students who would demand the services of libraries are part of the government-managed education system. Public schools are free at the primary level and have low tuition fees at the secondary level, but provide a lower quality of education than private schools (Núñez, Steiner, Cadena and Pardo, 2002).<sup>10</sup> In regards to the second characteristic, the fact that some parents may prefer specific types of education such as religious institutions or gender-specific schools could be correlated with demand side factors. With respect to the start of the academic year, schools can be calendar A or calendar B, which means they start in January or August, respectively. While calendar A is the norm, calendar B schools are typically private institutions usually designed in order to follow European or US schedules. This typically means that calendar B schools have higher test scores due to the strong selection related to the high income of students' families. Finally, schools can serve students for a full school day (12 hours) or implement double-shifts, with some students

9 More explicitly:  $AdjustedRD = \frac{RD - \hat{\beta}_0}{\hat{\beta}_1}$ , where  $\hat{\beta}$  come from the OLS regression between road distance

$$RD \text{ and Euclidean one } ED: RD = \beta_0 + \beta_1 ED + u$$

10 A small number of public schools are managed by the private sector and seem to follow a different pattern (Sarmiento, Alonso, Duncan and Garzón, 2005). None of them is close enough to our libraries.

coming in the morning and others in the afternoon.<sup>11</sup> Double-shifting is usually associated with lower academic results in the Latin American context as documented by Bonilla-Mejía (2011).

### III. Data

#### A. Quality of Education Data

Our measure of education quality is the Colombian equivalent to the SAT, the Saber 11° test administered by the Icfes (Colombian Institute for Evaluation of Education) which is part of the Ministry of Education. It includes a comprehensive evaluation of different areas of knowledge, specifically mathematics, verbal and sciences (biology, physics and chemistry). The test is carried out twice per year due to the existence of two main school calendars, and, though it is not compulsory for graduation, it is an entry requirement by universities in order to use it as a common filter for selecting their new students. In order to ensure comparability, test results are standardized by wave at the Bogotá level in each one of the described subject areas and an average is taken of the scores (called here the general result).

Tables A1.3 and A1.2 show average, standardized test scores of schools according to their characteristics including only the universe of schools used in the estimation, specifically, Bogotá schools located within a 3.5 Km range around the libraries as shown in Figure 1. Table A1.3 shows that students attending schools with a full-day schedule score higher, on average, than students attending double-shift schools. Among the latter, the students attending school in the morning score higher, on average, than those attending schools in the afternoon. This is related to the management of the school: students attending those managed by the government typically do worse than those managed by the private sector, which are normally private institutions. These relationships are stable over time and a common factor in the Colombian quality of education literature (Gaviria and Barrientos, 2001). Table A1.2 shows that there are also differences in test scores between students who attend different types of schools in terms of school size, the teacher-student ratio, the female-male

---

11 Other schools include night shifts or weekend shifts, but we will not consider them. Typically, these institutions are intended for young adults, who want to finish their secondary education after dropping out, thus the education incentives and the environment is totally different from a typical student.

student ratio, and teacher education level. These are all traditional inputs of education that we will discuss further in the next section

Table A1.4 shows a U-shape relationship between school quality and distance to the libraries. Schools close to the libraries are normally better than those at a medium-range distance (1 Km - 2.5 Km), but worse than or similar to those far away (2.5 Km - 3.5 Km). As this relationship might be driven by the allocation of inputs, our next section will analyze them in more detail.

## B. Other Variables and Data Restrictions

In order to take into account other sources of variation that might be correlated with distance to the libraries, we take into account variables that the literature has identified as key determinants of the quality of education. Variables used to control for institutional characteristics come from the C600 (a registry of students and school staff) and C100 (a registry of school infrastructure) from the Ministry of Education. Neighborhood controls are derived from the General Population Census of 2005 conducted by DANE (national statistics department). The relationship of these variables to our measures of quality of education is described in Table A1.3.

Though C100 information is only available starting from 2002, it provides valuable information on the physical infrastructure of schools. It includes data on sports facilities, the presence of a school library and a measure of the quality of educational assets, a dummy which is one if the school has simultaneously computer, physics and chemistry labs. From the C600 form we introduce several time-varying variables per school which are related to the supply-side of quality of education. First, we take into account the number of students per school in a logarithmic scale and the teacher-pupil ratio of the school. Larger schools are correlated with better results. To provide us with an idea of the overall quality of the facilities, we include the area in squared meters of classrooms and sport facilities per student. We also take into account the proportion of teachers with a graduate degree as a proxy of their human capital. As the public sector incentivizes the concentration of teachers with more qualifications, its relationship with quality seems to be negative as described by Núñez et al. (2002). Gender differences might be relevant, so we include the proportion of female students and teachers. Finally, we include some controls specific to the examined cohort: its size and the ratio of female test-takers.

This data was cleaned by removing schools with a teacher–student ratio greater than 0.5 (one teacher for every two students) or equal to 0 (no teacher to student) as these ratios indicate that the data may contain errors.

Finally, neighborhood–level controls are available at the census block level from 2005. We averaged the information of the blocks which were at least 50 meters from the school.

These controls are the average age and the share in the block of the population who are students, who have at most primary education, who immigrated from other municipalities and from rural areas during the last 5 years, who are of working age, who are working or looking for a job and who fasted for one week.

Tables A1.5 and A1.6 report for different ranges of distance respect to the library (column 1) the number of schools–students (column 2) and the number of schools–students used in the model (column 3), respectively.<sup>12</sup> The difference between columns two and three are due to information gaps either by C600. Hot Deck imputation methodology was used to minimize the number of missing, following the implementation of Báez and Buitrago (2010) based on Ñopo (2008) idea about donors and receptors.

### C. Test Scores and Distance to the Libraries

After observing the data on the relationships between some features of the campus and the quality of education, and considering the causal impact that the literature attributes to these features, it is prudent to identify whether the location of the libraries is correlated with the type of schools. Table A1.7 addresses this question by calculating the average characteristics of schools that are located in different ranges from the nearest mega–library. The main observation is that the nearest schools are more likely to be public. As public schools tend to have lower test scores (Gaviria and Barrientos, 2001; Núñez et al., 2002), the correlation between education quality and the distance of the libraries is negative. A first approach to the impact of libraries on test–scores score is to explore the score–distance relationship after deducting the impact of variation of common determinants from the score. For this, we turn

---

12 In the case of public institutions with a school is considered as the combination seat–day.

to a classic semi-parametric model. A partial linear regression allows us to see a non-linear relationship as presented in Equation 3.<sup>13</sup> In it,  $Y$  is the score,  $X$  is the controls,  $u$  is an error such that  $E[u | d, X] = 0$ . Figure A2.2 shows the estimates  $\hat{m}(d)$ , which gives the relationship between the score and the distance variation by discounting usual controls.

$$Y = m(d) + X\beta + u \quad (3)$$

We found a U-shaped relationship where the minimum is centered near 1500 meters. As a result, our analysis will be particularly focused on schools located between 750 and 2000 meters from libraries, where the impact of libraries are likely to reach. However, these graphs are used just to explore the relationship, because they include unobserved determinants  $u$ , in fact the U pattern is found both before and after 2002.

To estimate the effect we must assume that unobservable variables can vary across the distance, but the time variations of these unobservable variables are not related with distance. This restriction allows us to identify the average impact on the schools 'close' to the libraries compared to those that are 'distant' and supports the motivation to use the DiD strategy, as it will discuss in the next section.

#### IV. Empirical Strategy

The impact of libraries on quality of education is identified using the Difference in Difference (DiD) method. We define the schools 'near' to the libraries as treated, and those 'far' from it as controls. That is, we are assuming that any difference between these two groups of schools would have been preserved if no libraries were constructed (parallel trends assumption). It is important to remember that in these cases the 'libraries' refer to the entire intervention on the public infrastructure and urban planning development that occurred in those areas. Thus, the estimation is based on the provision, not the intensity of use, of libraries which is assumed to be a function of the distance of the school to the physical building.

---

13 The estimation was performed following the algorithm differences Yatchew (1997), implemented by Lokshin (2006).

The identification strategy involves two stages: the first refers to measure the magnitude and significance of the impact, and the second is to decompose it into the impact due to changes in observed inputs and to variations not linked to those inputs. The decomposition addresses the question of complementarities between libraries and traditional determinants of the quality of education, in other words, how the libraries enhance the impact of traditional inputs already present in schools.

As described before, our treatment indicator is the spatial proximity from schools. However, being 'near' or 'far' is an arbitrary definition and requires a selection rule that is part of the research question. Discrete and continuous options were considered to define exposure treatment using the distance of each school to the libraries,  $d$ .

A first alternative (*continuous approach*) is to impose a parametric restriction on the relationship between the distance to library and test scores. Given the results from the partially-linear regression, it is possible to presume that the impact decreases with the inverse of distance up to some far, arbitrary cutoff  $R1$  where we set the impact to be exactly 0, including all the schools within a fixed radius  $R2$ . Hence, we define  $T = \frac{R1}{d} - 1$  if  $d \leq R1$  and  $T = 0$  if  $d \geq R1$ . For this specification we present results for ratios  $R1 \in \{1500, 2000, 2500, 3000, 3500\}$  and  $R2 = 3500$ .

On the hand, the effect could be discontinuous (*discrete approach*). Hence, in order to avoid any assumption on the distance-scores' relation, schools within a certain ratio,  $R2$  is assigned into treated  $T = 1$  and control groups  $T = 0$  using an arbitrary distance to the library cut-off  $R1$ . This specification, henceforth *Discrete I*, is represented in Figure 1. An alternative, *Discrete II*, is to omit some schools between treatment and control zones, so the control zone starts at  $R3 \in [R1, R2]$ . Implementing different cut-offs in the analysis did not show substantial differences. We will present results using  $R2 = 3500$ ,  $R3 = 2000$  and  $R1 \in \{750, 1000, 1250, 1500, 1750, 2000\}$ .

## A. Estimation of the General Impact (DiD)

We define the average treatment effect on the treated  $\delta^\tau$ , as the impact on average test scores at year  $\tau$  for schools that are located close to the libraries



in comparison to those that are far from them. If we consider the continuous treatment scenario, the fullest impact occurs for schools that are located right next to one of the libraries. This parameter is estimated using the classic setup as presented in equation 4. Let  $Y_{it}$  be the average test scores of school  $i$  at year  $t$ ,  $T_i$  the treatment status of each school,  $A_t$  a dummy that is 1 if  $t \geq 2003$ ,  $1(\tau = t)$  is an indicator for year  $\tau$  being equal to year  $t$ , and fix effects  $Y_i$  and  $Y_t$ . For this specification, we assume that the parallel trends hold conditional on the school-level controls  $X_{it}$ .

$$Y_{it} = \sum_{\tau=2003}^{2008} \delta^\tau T_i \cdot 1(\tau = t) + \beta_1 A_t + \eta X_{it} + \gamma_i + \gamma_t + e_{it} \quad (4)$$

The identification assumption might be too strong; schools placed in different areas might follow dissimilar trends due to uncontrolled factors. For instance, migration of people with different willingness to spend on education may shape schools' investments in a way that is not captured by our current covariates. In essence, some schools might be improving while others worsening. In order to address this, we can include school-specific trends<sup>14</sup>,  $t \cdot Y_i$ , as shown in equation 5. The limitation of this approach is that trends can differ only as long as they do so in a linear fashion.

$$Y_{it} = \sum_{\tau=2003}^{2008} \delta^\tau T_i \cdot 1(\tau = t) + \beta_1 A_t + \eta X_{it} + \gamma_i + \gamma_t + \omega_i t \cdot \gamma_i + e_{it} \quad (5)$$

## B. Propensity Score Matching and Synthetic Control

One of the main concerns with the DiD method for studies with limited control units is how to choose the best control when there are few treated units, which implies high sensitivity of the estimation to the control selection, and when the unit of observation is an aggregate (eg. countries, states or schools). Abadie and Gardeazabal (2003) introduced an approach known as the 'synthetic control' to deal with these problems. The idea is to select a set of weights for the control units to construct the parallel trends between outcomes before the intervention. However, as is suggested by Abadie and Gardeazabal (2003), the synthetic control needs a long period of time prior to the intervention in order to control for structural patterns in both observables and non-observables (Abadie, Diamond and Hainmueller, 2010). Given that there are just three

14 For other applications that introduce this technique, see for instance, Besley and Burgess (2004).

years available before the implementation of the mega-libraries and that the objective is to forecast over the next six years, the synthetic control strategy might lead to misleading results. An alternative that might be more suitable is to weaken the DiD parallel trends assumption by introducing matching into the pre-treatment period (Blundell and Dias, 2009). The matching estimator relies on the minimization of a distance function which is increasingly hard to estimate with the number of included covariates. A traditional way to simplify this problem, when there is more than one treated unit, is to perform the matching based on the predicted likelihood of being a treated unit, the propensity score (Rosenbaum and Rubin, 1983).

In this paper we combine both approaches by implementing kernel propensity score matching<sup>15</sup> (Heckman, Ichimura and Todd, 1997) that includes as controls the pre-treatment evolution of test scores, which is in line to the synthetic control matching step. Once the synthetic control is constructed by re-weighting the non-treated schools, DiD specifications from equations 4 and 5 are applied.<sup>16</sup>

In doing so, the underlying identification assumption changes slightly. Once the observed covariates are taken into account, and schools close and far from the libraries follow similar time-trends or differ in a linear way, estimated impacts can be attributed to the mega-libraries. However, keep in mind that the identification will be invalid if there were events that were not considered and affected some of the schools (either close or far from libraries) and not the others.

Apart from the 2000-2002 test scores, the matching variables considered are the following: the proportion of teachers with graduate studies, pupil-teacher ratio, public school dummy, morning school day dummy, complete school day dummy, female-teacher ratio, 11th grade female-male students ratio, 11th grade students, total students, girls-students ratio, built area per student, classrooms area per student, sports area per student, and a dummy for the presence of a school library.

---

15 The procedure was implemented using `psmatch2` (Leuven and Sianesi, 2014) in Stata 12.

16 As the matching is based on discrete categories, the continuous approach cannot be implemented.

### C. DiD-OB: Decomposition of the Impact

As discussed, the construction of the libraries implied a massive urban development. As a result, it is likely the mega-libraries triggered changes in other inputs. For instance, the construction of mega-libraries could lead to emigration from the area due to changes in real estate prices, also they could change the number of private schools or the ratio of teacher-student. Thus part of the observed changes between schools close and far from libraries would be due to this channel. Hence, we would be interested on see if the program had an impact on the inputs and such variation explain part of the outcomes difference, let's call that part  $\Delta_x$  and if there is part of that impact that is not due to them,  $\Delta_0$ , instead this part of impact could be due to changes on the impact that teachers with high level of education could has with the presence of the libraries or could be due to changes in the efficiency of public schools who engage with the libraries' services. In that case,  $\Delta_0$  would be more likely to be related with the complementarity between schools and libraries. This is achieved by implementing a novel strategy, proposed in this study, that introduces the Oaxaca (1973) and Blinder (1973) decomposition into a DiD context (see the appendix for details). The conditions for the identification of the effect are the usual parallel trends of DiD but without conditioning on covariates. The decomposition is obtained by applying equation 6.

$$\begin{aligned}
 Y_{it} = & \alpha_0 + \alpha_1 T_{it} + \alpha_2 A_{it} + \alpha_3 X_{it} + \alpha_4 X_{it} \cdot T_{it} \\
 & + \alpha_5 X_{it} \cdot A_{it} + \alpha_6 T_{it} \cdot A_{it} + \alpha_7 X_{it} \cdot T_{it} \cdot A_{it} + u
 \end{aligned}
 \tag{6}$$

From this equation, we can define the impact generated by the covariates variation (induced by the program)  $\Delta_x$ , and the variation that is unrelated to them,  $\Delta_0$ :

$$\begin{aligned}
 \delta = & (E[y | T = 1, A = 1] - E[y | T = 0, A = 1]) \\
 & - (E[y | T = 1, A = 0] - E[y | T = 0, A = 0]) \\
 \delta = & \Delta_0 + \Delta_x \\
 \delta = & \alpha_6 + (\alpha_4 + \alpha_5 + \alpha_7)E[X | T = 1, A = 1] \\
 & - \alpha_5 E[X | T = 0, A = 1] - \alpha_4 E[X | T = 1, A = 0] \\
 & + \alpha_3 [(E[X | T = 1, A = 1] - E[X | T = 0, A = 1]) \\
 & - (E[X | T = 1, A = 0] - E[X | T = 0, A = 0])]
 \end{aligned}
 \tag{9}$$

Standard errors are calculated by bootstrapping due to the lack of an analytical expression for them. In order to present results by year, the strategy is implemented by comparing the pre-intervention period against each treatment-year in a separate regression.

## V. Results and Discussion

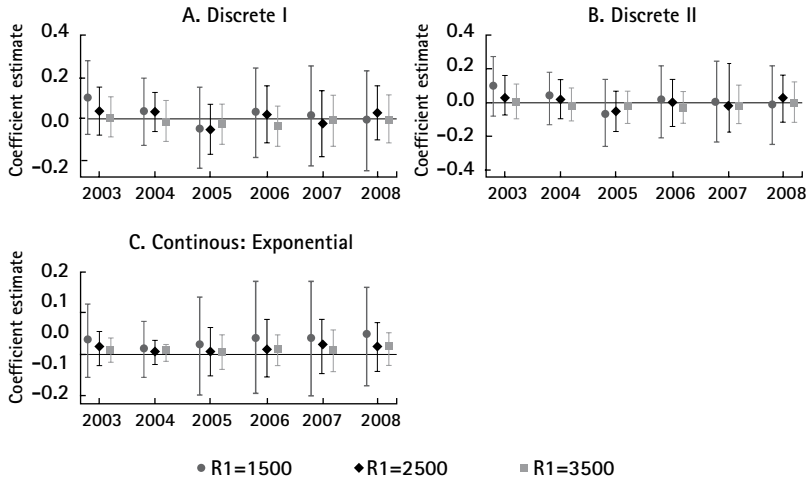
### A. Classic DiD strategy

First using the parametric approach, we compare the evolution of the treatment group in each year from 2003 to 2008 against the pre-treatment period, 2000 to 2002. In Table A1.8, we consider the intensity of treatment to be inversely proportional to the distance. It ranges from 1, the intensity received by a school in front of the library, to 0, a school that is located  $R_1$  meters or further. The general impact of being just beside the library implies an increase on average scores between 0.02 and 0.06 standard deviations ( $R_1=1500$  for 2003 and 2008, respectively). This impact is lower when we assume that there is a slower decay in the benefit received based on distance (higher  $R_1$ ), suggesting that the area of the impact is relatively small. However, those impacts are not statistically different from 0.

Table A1.9 presents the results from the discrete approach. In Panel A the treatment group are those schools between 0 and  $R_1$  meters from the libraries and the controls are those from  $R_1$  to  $R_2$  (fixed at 3.5 Km), as shown in the map from Figure 1. Estimates range between 0.21 for the lowest ratio in 2005 and -0.05 for the largest. This is consistent with the previous specification, which found that the impact is greater for the nearest schools. However, there is no evidence of impact different from 0. Similar results are found in the last specification, shown in panel B, where the controls are those schools between  $R_3 = 2000$  and  $R_2$ . That is, we are not taking into account those schools between  $R_1$  and  $R_3$  meters. These results are also presented in Figure 3, as a reference for comparison.

Equation 5 relaxed the parallel trends assumptions by allowing school-specific trends. Figure A2.3 shows that for both the discrete specification II and the continuous approaches, schools which are very close to the libraries seems to have a declining trend in the outcome. However, that pattern is still statistically non-different to zero.

Figure 3. Euclidean Distance Estimators



Confidence Intervals at 95% level.

Source: Own calculation based on C-600 and Saber 11°.

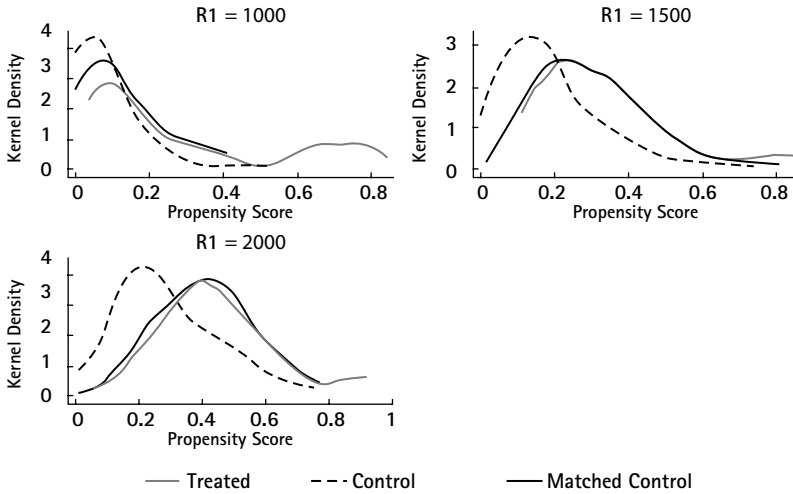
One clear concern is the measure of distance. The Euclidean approach might not capture the real cost to travel between points in certain contexts. For instance, there might be restrictions due to geographic accidents or infrastructure. However, in this urban context it might not be a bad approach. An alternative that takes these issues into account is road distance, which measures the total distance necessary to reach a mega-library while using the road infrastructure. Figure A2.4 presents the main estimates using this approach. In order to be able to compare both main and additional results, the road distance was rescaled using a linear function (see section II) as the relevant difference might come not from the absolute position of each school but from the relative one. The remainder of this paper will consider only the Euclidean measure.

## B. Synthetic Control

The next step is to introduce the matching strategy into the DiD. The main objective is to ensure that schools which are close to the libraries are compared to similar schools that are far from them. In order to achieve this, these schools were matched on the propensity score. Figures 4 and 5 show that once the matching weights are introduced, the propensity score calculated for the synthetic control group resembles the one of the treated schools (according to

the treatment definition). The purpose of this step is to ensure that by matching the score, the covariates are matched as well.

Figure 4. Propensity Score Matching at 2002: Discrete I



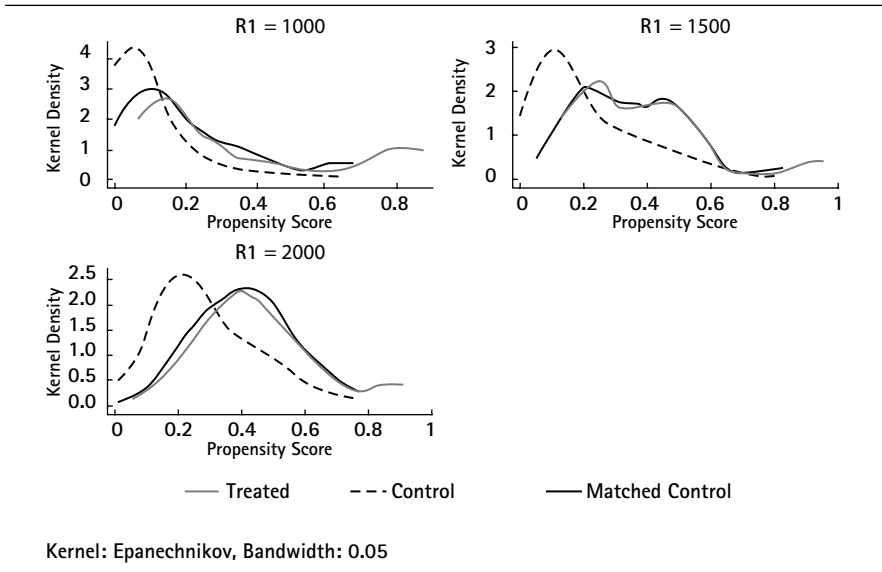
Kernel: Epanechnikov, Bandwidth: 0.05

Source: Own calculation based on C-600 and Saber 11<sup>o</sup>.

We can check the performance of the technique in Tables A1.10 and A1.11, for both discrete specification I and II respectively. For each distance definition, the tables present the difference for each match variable between treatment and control groups before (General) and after (Matched) the matching as well as the percentage reduction on the standardized bias (B.R.). Starts on the tables reflect the results of t-tests for equality of means for each difference where the null hypothesis is that the differences are equal to 0. The matched results appear balanced, and, giving that we are matching the outcome trend before the intervention, the resulting synthetic control group trend closely resembles that of the treatment. A graphic representation of this is presented in Figures A2.5 and A2.6. The only one for which the technique does not look as successful is for specification II, where the treatment seems to be following a quite different trend.

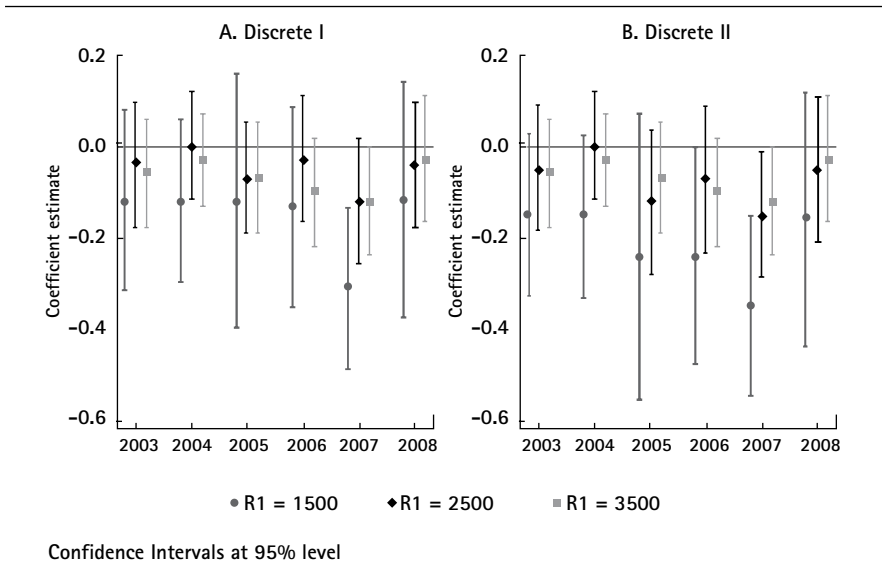
Apart from the quality of matching, Figures A2.5 and A1.6 also tell another story. It seems that schools which are closer to the libraries have a decreasing

Figure 5. Propensity Score Matching at 2002: Discrete II



Source: Own calculation based on C-600 and Saber 11°.

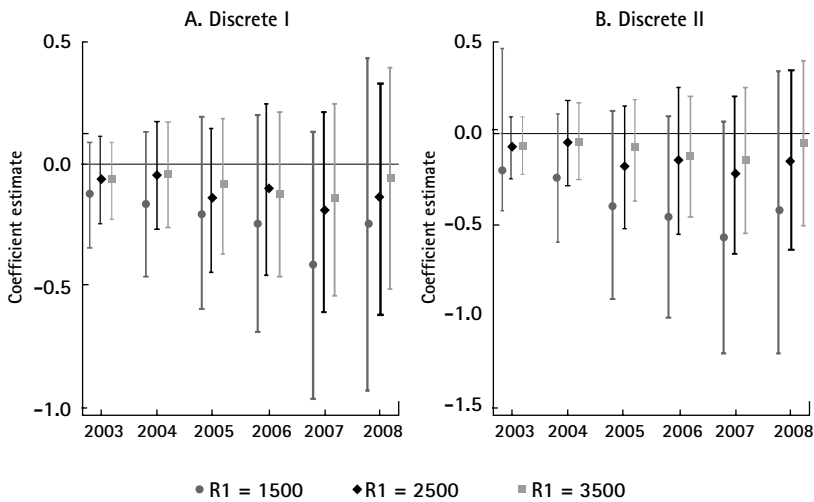
Figure 6. Matching at 2002 Estimators



Source: Own calculation based on C-600 and Saber 11°.

trend compared to distant schools which are comparable in key covariates. This is reflected in the DiD estimates in Figure 6. In contrast with Figure 3, almost all of the estimates are negative, and, for years 2006 and 2007, some of them are significant. In other words, after the libraries were constructed, schools nearby, especially those which are very close to the libraries, started to perform worse than similar ones not as close to the libraries. This means that either the libraries or the urban development in their surroundings did decrease student performance relative to their peers<sup>17</sup> or that the identification assumption is not as good as desired.

Figure 7. Matching at 2002 Estimators With School-Specific Trends



Confidence Intervals at 95% level

Source: Own calculation based on C-600 and Saber 11°.

As described before, Figure A2.6 for the 1000 meter definition according to specification II shows that the declining trend for some of these schools started prior to the construction of the libraries which was not fully controlled for by the matching. In order to assess this, performance data was de-trended by school (see Equation 5). Figure 7 and Table A1.12 present the results of this approach. Estimated coefficients are still negative but are not different from zero.

17 It might be that these schools did perform better, but not as much as to other schools in the city which is the base of our standardization.



### C. Blinder–Oaxaca Decomposition

So far it seems that there is no significant variation on the relationship between distance to the libraries and average tests scores on mathematics, science and verbal sections.<sup>18</sup> It might be the case that the urban transformation was related to changes in inputs in the quality of education production function. Table A1.13 studies this via the Oaxaca–Blinder DiD decomposition proposed before, but we should bear in mind that the identification assumptions are stronger than in the simple DiD analysis. In most of the cases, it seems that the difference between schools far and close to the libraries on test scores due to the observed inputs is negative ( $\Delta_x$ ). The direct impact of the libraries on test scores ( $\Delta_o$ ) is around 0.1 and 0.2 standard deviations for schools located between 0 and 1.5 Km from the libraries. As a reference, the difference between students with college graduated mothers and the others in the same sample (3.5 Km at most for each library) is 0.6 standard deviations. However, these results are not different from 0.

### D. Summary

The fact that estimation procedures with different sets of assumptions provide similar results gives us a good idea of the underlying relationship between the construction of mega-libraries and quality of education: there is no evidence of a positive and statistical significant impact of the libraries on average standardized scores. We can interpret these results in many different ways. First, the fact that the numbers are positive but the variance is large could be related to the small number of observations available (around 190 schools per year). If that is the case, any significant positive relationship between public libraries and schools' scores, is likely to be small. This does not mean that the libraries are useless for education: they could improve other skills that are not related with tests scores but which are important for the society, such as the availability of safe spaces and exposure to cultural activities. Current information makes it impossible to test those alternatives. Second, the high variance could be due to the positive impact of libraries only on those schools, students or teachers that decided to take advantage of the libraries and zero impact on those that did not. Heterogeneous impacts are the rule, not the exception, in

---

18 Results for each one of these scores separately are not meaningfully different from the ones presented here

the literature of educational inputs (Murnane and Ganimian, 2014).<sup>19</sup>Without further information on the selection mechanism, it is impossible to determine the impact only on those schools, students or teachers that are willing to take advantage of the public infrastructure.

In the case that some schools, students or teachers within similar distances to libraries use the libraries facilities at different rates, policy may not only be needed to construct and run these public facilities but also to impose incentive schemes that induce to use them. Glewwe and Kremer (2006) argue that the provision of resources is insufficient to improve student performance and the teachers should be instructed in order to maximize the potential advantage of the resources. Moreover, using the theoretical framework proposed by Witte and Geys (2011), the provision of most public goods, in this case the libraries, need two stages of policies: the first one for the construction of the libraries, while the second should work on how these programmatic inputs are transformed into observed and desired outputs of education. For instance, prizes for both teachers and students for projects that involve the usage of these resources might be relevant.

## VI. Conclusions

We have analyzed the impact on the quality of education, measured by mathematics, science and verbal Saber 11<sup>o</sup> scores, of the construction of two big, public libraries that involved the transformation of low-income, urban areas in Bogotá, Colombia. To do so, we measured how the construction of the libraries could change the test scores of nearby schools, controlling for observable variables that are related to students' performances. We opted for a DiD approach to analyze the evolution of the relation of distance-to-library and average test scores before and after the public libraries' introduction at the school level. This approach assumes that the effect of the libraries decays with distance and that, without the intervention, the relationship would have been unaltered over time. We also propose and implement a decomposition of the

---

19 Murnane and Ganimian (2014) remark three cases: High- and low-education parents responded very differently to initiatives to empower school councils in Niger (Beasley and Huillery, 2012); low- and high achieving students derived very different benefits from free textbooks in English in Kenya (Glewwe et al., 2009); and rural girls did not profit nearly as much as urban boys from the use of LEGO kits to teach science in Peru (Beuermann, Naslund-Hadley, Ruprah and Thompson, 2013)

effect considering the potential variations of traditional determinants of quality of education.

The libraries analyzed are public, education-related infrastructure that is progressive in a context of inequality in access to quality school education. Both libraries were built in areas populated by the less well-off and where schools have relatively poor facilities. Thus, the policy has the potential to boost the equality of opportunities in terms of quality of education. However, our findings present non-statistically different from zero impacts of the libraries on the average standardized test scores. That is, there is no evidence that schools close to the libraries are getting a clear advantage on test scores against those with similar characteristics but for their location further from the new public infrastructure.

It is important to remark that the results are correct only under the validity of the assumptions defined in the identification strategy. In general, there are two main scenarios in which the assumptions would be invalid. First, if it is the case that the intensity of the use of libraries is unrelated to the distance from them. For instance, there could be a network of teachers which take advantage of library facilities though their schools are not close to the libraries. Another reason could be that the network of medium and small libraries communicates perfectly with the more distant mega-libraries, thus there is not difference in access according to the distance. Second, it might be the case that schools close and far from the libraries were affected heterogeneously by other events which are not fully captured by observed covariates. As an example, patterns of migration or criminality in the zones that are near to the libraries which did not affect cohort sizes, gender composition, or any other observed inputs with respect to the other neighborhoods could explain those results.

These results do not necessarily mean that libraries do not improve the quality of education. On one hand, libraries might be related to skills that are not directly reflected in test scores or to these types of skills but for students in older stages of their lives, such as college students. We are unable to assess these cases via the present methodology. On the other hand, if a direct objective of these types of programs is to enhance test scores, our results imply that the policies that introduced these public facilities should be complemented with stronger programs which link and coordinate them with the already existent educational institutions. The capacity to reach the target (school-students and

teachers) is an important part of the policy which might require more attention from local governments. For instance, prizes for both teachers and students for projects that involve the usage of these resources might be relevant.

## References

1. ABADIE, A., DIAMOND, A., AND HAINMUELLER, J. (2010). "Synthetic control methods for comparative case studies: Estimating the effect of California's tobacco control program", *Journal of the American Statistical Association*, 105(490).
2. ABADIE, A., AND GARDEAZABAL, J. (2003). "The economic costs of conflict: A case study of the Basque country", *American Economic Review*, 93(1):113-132.
3. BARRERA-OSORIO, F., AND LINDEN, L. L. (2009). The use and misuse of computers in education: Evidence from a randomized experiment in Colombia (Technical Report 4836). World Bank.
4. BEASLEY, E., AND HUILLERY, E. (2012). *Empowering parents in schools: What they can (not) do*. Cambridge, MA: Abdul Latif Jameel Poverty Action Lab (J-PAL).
5. BESLEY, T., AND BURGESS, R. (2004). "Can labor regulation hinder economic performance? evidence from India", *The Quarterly Journal of Economics*, 119(1):91-134.
6. BEUERMANN, D. W., NASLUND-HADLEY, E., RUPRAH, I. J., AND THOMPSON, J. (2013). "The pedagogy of science and environment: Experimental evidence from Peru", *The Journal of Development Studies*, 49(5):719-736.
7. BÁEZ, N. A. D., AND BUITRAGO, C. F. (2010). Ingresos en el sistema de identificación de potenciales beneficiarios de programas sociales (Sisbén): tres metodologías de imputación (Archivos de Economía 006451). Departamento Nacional de Planeación.

8. BLINDER, A. (1973). "Wage discrimination: Reduced form and structural estimates", *Journal of Human resources*, 8(4):436-455.
9. BLOMEYER, D., CONEUS, K., LAUCHT, M., AND PFEIFFER, F. (2009). "Initial risk matrix, home resources, ability development, and children's achievement", *Journal of the European Economic Association*, 7(2-3):638-648.
10. BLUNDELL, R.G AND DIAS, M. (2009). "Alternative approaches to evaluation in empirical microeconomics", *Journal of Human Resources*, 44(3):565-640.
11. BONILLA-MEJÍA, L. (2011). Doble jornada escolar y calidad de la educación en Colombia (Documento de Trabajo sobre Economía Regional 143). Banco de la República.
12. BORKUM, E., HE, F., AND LINDEN, L. L. (2013). The effects of school libraries on language skills: Evidence from a randomized controlled trial in India (Discussion Papers 7267). Institute for the Study of Labor (IZA).
13. FERTIG, M., AND SCHMIDT, C. M. (2002). The role of background factors for reading literacy: Straight National Scores in the PISA 2000 Study (Discussion Papers 545). Institute for the Study of Labor (IZA).
14. FORTIN, N., LEMIEUX, T., AND FIRPO, S. (2011). "Decomposition methods in economics", *Handbook of Labor Economics*, 4:1-102.
15. GAMBOA, L. F., AND RODRÍGUEZ-LESME, P. A. (2014). Do Colombian students underestimate higher education returns? (Documentos de Trabajo 164). Universidad del Rosario, Facultad de Economía.
16. GAMBOA, L. F., RODRÍGUEZ-ACOSTA, M., AND GARCÍA-SUAZA, A. (2010). Academic achievement in sciences: The role of preferences and educative assets (Documentos de Trabajo 78). Universidad del Rosario, Facultad de Economía.
17. GAVIRIA, A., AND BARRIENTOS, J. (2001). "Características del plantel y calidad de la educación en Bogotá", *Coyuntura Social*, 25:81-98.

18. GLEWWE, P., AND KREMER, M. (2006). "Schools, teachers, and education outcomes in developing countries", *Handbook of the Economics of Education*, 2:945-1017.
19. GLEWWE, P., KREMER, M., AND MOULIN, S. (2009). "Many children left behind? textbooks and test scores in Kenya", *American Economic Journal: Applied Economics*, 1(1):112-135.
20. GLEWWE, P., KREMER, M., MOULIN, S., AND ZITZEWITZ, E. (2004). "Retrospective vs. prospective analyses of school inputs: The case of flip charts in Kenya", *Journal of development Economics*, 74(1):251-268.
21. HANUSHEK, E., AND WOBMANN, L. (2007). The role of education quality in economic growth (Technical Report). World Bank.
22. HECKMAN, J. J., ICHIMURA, H., AND TODD, P. E. (1997). "Matching as an econometric evaluation estimator: Evidence from evaluating a job training programme", *Review of Economic Studies*, 64(4):605-654.
23. LANCE, K. (1994). "The impact of school library media centers on academic achievement", *School Library Media Quarterly*, 22(3):167-170.
24. LANCE, K., RODNEY, M., AND HAMILTON-PENNELL, C. (2000). Measuring up to standards: The impact of school library programs & information literacy in Pennsylvania schools (Technical Report). Pennsylvania State Dept. of Education, Office of Commonwealth Libraries.
25. LEUVEN, E., AND SIANESI, B. (2014). *Psmatch2: Stata module to perform full mahalanobis and propensity score matching, common support graphing, and covariate imbalance testing*. Statistical Software Components.
26. LOKSHIN, M. (2006). "Semi-parametric difference-based estimation of partial linear regression models", *Stata Journal*, 6(3):377-383.G
27. LONSDALE, M. (2003). Impact of school libraries on student achievement: A review of the research (Information Analyses 70). Australian Council for Educational Research, Victoria.

28. MURNANE, R. J., AND GANIMIAN, A. J. (2014). Improving educational outcomes in developing countries (Lessons from Rigorous Evaluations 20284). National Bureau of Economic Research.
29. MURNANE, R., MAYNARD, R., AND OHLS, J. (1981). "Home resources and children's achievement", *The Review of Economics and Statistics*, 63(3):369-377.
30. ÑOPO, H. (2008). "Matching as a tool to decompose wage gaps", *The Review of Economics and Statistics*, 90(2):290-299.
31. NÚÑEZ, J., STEINER, R., CADENA, X., AND PARDO, R. (2002). "¿Cuáles colegios ofrecen mejor educación en Colombia?", *Archivos de Economía*, 193.
32. OAXACA, R. (1973). "Male-female wage differentials in urban labor markets", *International Economic Review*, 14(3):693-709.
33. RODNEY, M., LANCE, K., HAMILTON-PENNELL, C., AND CENTER, M. (2002). Make the connection: Quality school library media programs impact academic achievement in Iowa, Mississippi Bend Area Education Agency.
34. ROSENBAUM, P. R., AND RUBIN, D. B. (1983). "The central role of the propensity score in observational studies for causal effects", *Biometrika*, 70(1):41-55.
35. SARMIENTO, A., ALONSO, C. E., DUNCAN, G., AND GARZÓN, C. A. (2005). *Evaluación de la gestión de los colegios en concesión en Bogotá 2000-2003*. Bogotá: DNP.
36. SMITH, E. (2001). *Texas School Libraries: Standards, resources, services, and students performance*. EGS Research & Consulting.
37. TOLOSA, L. R. T. (2012). "Breve historia de las bibliotecas públicas en Colombia", *Códices*, 8(1):57-86.

38. VEGAS, E., AND PETROW, J. (2008). *Raising student learning in Latin America: The challenge for the 21st century*. World Bank Publications.
39. WILLIAMS, D., WAVELL, C., AND COLES, L. (2001). *Impact of school library services on achievement and learning*. Technical report, Department for Education & Skills and Resources: The Council for Museums, Archives & Libraries.
40. WITTE, K. D., AND GEYS, B. (2011). "Evaluating efficient public good provision: Theory and evidence from a generalised conditional efficiency model for public libraries", *Journal of Urban Economics*, 69(3):319-327.
41. WORLD BANK. (2005). Mexico: Determinants of learning policy note (Report 31842-MX). World Bank, Washington D. C.
42. YATCHEW, A. (1997). "An elementary estimator of the partial linear model", *Economics Letters*, 57(2):135-143.



## Annex 1. Tables

**Table A1.1** Travelling time to school

Time	Freq.	Cum.
Less than 10 min.	51%	51%
Between 10 and 20 min.	26%	77%
Between 20 y 30 min.	23%	100%

Source: DANE Population Census 2005

**Table A1.2** Average test score by institutional and environment characteristics

	Year								Total
	2000	2001	2003	2004	2005	2006	2007	2008	
<b>School day</b>									
Complete	0.040	-0.075	-0.017	-0.046	0.020	0.016	0.030	0.051	0.002
Morning	-0.065	-0.209	-0.193	-0.288	-0.251	-0.313	-0.394	-0.365	-0.263
Afternoon	-0.252	-0.451	-0.356	-0.362	-0.420	-0.474	-0.460	-0.475	-0.408
<b>Total</b>	<b>-0.082</b>	<b>-0.234</b>	<b>-0.174</b>	<b>-0.217</b>	<b>-0.197</b>	<b>-0.235</b>	<b>-0.245</b>	<b>-0.234</b>	<b>-0.204</b>
<b>Type of school</b>									
Public	-0.139	-0.312	-0.256	-0.289	-0.339	-0.410	-0.432	-0.414	-0.328
Private	-0.034	-0.162	-0.097	-0.147	-0.062	-0.068	-0.073	-0.060	-0.088
<b>Total</b>	<b>-0.082</b>	<b>-0.234</b>	<b>-0.174</b>	<b>-0.217</b>	<b>-0.197</b>	<b>-0.235</b>	<b>-0.245</b>	<b>-0.234</b>	<b>-0.204</b>

Source: Own calculations based on Saber 11<sup>o</sup> (include imputations).

Table A1.3 Average test score by infrastructure and teaching force

	Year								Total
	2000	2001	2003	2004	2005	2006	2007	2008	
<b>Students</b>									
Less than 300	-0.26	-0.52	-0.44	-0.42	-0.41	-0.33	-0.41	-0.29	-0.38
Between 300-600	-0.22	-0.36	-0.15	-0.21	-0.13	-0.19	-0.26	-0.16	-0.21
Between 600-1000	-0.02	-0.15	-0.03	-0.14	-0.19	-0.17	-0.05	-0.14	-0.12
More than 1000	0.13	-0.01	-0.12	-0.09	-0.10	-0.16	-0.16	-0.18	-0.10
<b>Total</b>	<b>-0.06</b>	<b>-0.21</b>	<b>-0.15</b>	<b>-0.18</b>	<b>-0.19</b>	<b>-0.20</b>	<b>-0.20</b>	<b>-0.19</b>	<b>-0.17</b>
<b>Teacher-student ratio</b>									
Less than .03	-0.08	-0.57	-0.32	-0.29	-0.31	-0.34	-0.25	-0.28	-0.31
Between .03-.04	-0.06	-0.16	0.01	-0.19	-0.21	-0.25	-0.20	-0.27	-0.18
Between .04-.05	-0.11	-0.21	-0.00	-0.18	-0.00	-0.04	-0.17	0.01	-0.11
Between .05-.06	0.01	-0.10	-0.46	-0.03	-0.12	-0.19	-0.19	-0.27	-0.13
More than .06	-0.21	-0.45	-0.37	-0.35	-0.30	-0.34	-0.40	-0.36	-0.36
<b>Total</b>	<b>-0.08</b>	<b>-0.23</b>	<b>-0.17</b>	<b>-0.22</b>	<b>-0.20</b>	<b>-0.23</b>	<b>-0.25</b>	<b>-0.23</b>	<b>-0.20</b>
<b>Girls-students ratio</b>									
Less than 0.15	0.15	0.41	0.29	0.11	0.13	0.11	0.03	-0.03	0.15
Between 0.15-0.43	0.10	-0.06	-0.08	-0.11	-0.11	-0.15	-0.28	-0.15	-0.11
Between 0.43-0.48	-0.11	-0.29	-0.19	-0.19	-0.13	-0.18	-0.06	-0.20	-0.17
Between 0.48-0.52	-0.20	-0.40	-0.31	-0.30	-0.36	-0.37	-0.40	-0.33	-0.34
Between 0.52-0.85	-0.22	-0.36	-0.18	-0.39	-0.23	-0.28	-0.34	-0.42	-0.30
More than 0.85	0.24	0.26	0.33	0.13	0.09	0.01	0.03	0.16	0.16
<b>Total</b>	<b>-0.08</b>	<b>-0.23</b>	<b>-0.17</b>	<b>-0.22</b>	<b>-0.20</b>	<b>-0.23</b>	<b>-0.25</b>	<b>-0.23</b>	<b>-0.20</b>
<b>Basic level teachers</b>									
Less than .25	-0.03	-0.20	-0.11	-0.16	-0.19	-0.22	-0.20	-0.21	-0.17
Between .25-.5	-0.26	-0.35	-0.43	-0.47	-0.13	-0.38	-0.36	-0.21	-0.33
Between .5-.75	0.21	-0.03		-0.22	-0.36				-0.05
More than .75	-0.31	-0.83	-0.66	-0.52	-0.76	-0.59		-0.90	-0.61
<b>Total</b>	<b>-0.06</b>	<b>-0.22</b>	<b>-0.15</b>	<b>-0.20</b>	<b>-0.19</b>	<b>-0.23</b>	<b>-0.21</b>	<b>-0.21</b>	<b>-0.19</b>
<b>Highest Level teachers</b>									
Less than .25	-0.08	-0.25	-0.19	-0.24	-0.14	-0.18	-0.17	-0.17	-0.18
Between .25-.5	-0.16	-0.37	-0.15	-0.55	-0.17	-0.18	-0.81	-0.26	-0.30
Between .5-.75	-0.12	-0.28	-0.19	-0.27	-0.33	-0.32	-0.41	-0.34	-0.29
More than .75	-0.04	-0.20	-0.08	-0.18	-0.11	-0.30	-0.24	-0.32	-0.19
<b>Total</b>	<b>-0.09</b>	<b>-0.26</b>	<b>-0.18</b>	<b>-0.24</b>	<b>-0.18</b>	<b>-0.21</b>	<b>-0.23</b>	<b>-0.22</b>	<b>-0.20</b>

Source: Own calculations based on C600 and Saber 11° (include imputations).

Table A1.4 Average test score by distance

Distance to library	Years			Total
	2000-2002	2003-2005	2006-2008	
Less than 1000	-0.141	-0.048	-0.150	-0.110
Between 1000-2500	-0.239	-0.306	-0.346	-0.306
More than 2500	-0.112	-0.140	-0.177	-0.147
<b>Total</b>	<b>-0.161</b>	<b>-0.196</b>	<b>-0.238</b>	<b>-0.204</b>

Source: Own calculations based on Saber 11<sup>a</sup> (include imputations).

Table A1.5 Schools by distance

Distance to the library (meters)	Schools	Used in the models
0-500m	5	4
500m-1000m	15	11
1000m-1500m	28	27
1500m-2000m	30	24
2000m-2500m	48	40
2500m-3000m	45	39
3000m-3500m	45	38
3500m-4000m	59	49
<b>Total</b>	<b>275</b>	<b>232</b>

Source: Own calculations.

Table A1.6 Students by distance

Distance to the library (meters)	Students	Used in the models
0-500m	237	115
500m-1000m	2996	2888
1000m-1500m	5372	5178
1500m-2000m	5229	4820
2000m-2500m	6322	5629
2500m-3000m	7634	7032
3000m-3500m	6263	6086
3500m-4000m	7685	7195
<b>Total</b>	<b>41738</b>	<b>38943</b>

Source: Own calculations.

Table A1.7 Distribution by distances and school characteristics

	Distance to the library		
	Between 0 and 1 Km	Between 1 and 2 Km	Between 2 and 4 Km
	%	%	%
<b>Type of School</b>			
Public	59.84	58.96	50.80
Private	40.16	41.04	49.20
<b>Total</b>	<b>100</b>	<b>100</b>	<b>100</b>
<b>Post-graduated teachers ratio</b>			
Less than 30%	51.18	61.32	63.19
Between 30% y 60%	25.20	16.98	19.93
More than 70%	23.62	21.70	16.88
<b>Total</b>	<b>100</b>	<b>100</b>	<b>100</b>
<b>School day</b>			
Complete	31.50	37.26	42.90
Morning	35.43	28.07	24.64
Afternoon	33.07	34.67	32.46
<b>Total</b>	<b>100</b>	<b>100</b>	<b>100</b>
<b>Student-teacher ratio</b>			
Less than 20	20.47	21.70	23.84
Between 20 and 30	59.84	58.96	54.06
More than 30	19.69	19.34	22.10
<b>Total</b>	<b>100</b>	<b>100</b>	<b>100</b>
<b>School size</b>			
More than 1000 students	39.37	50.47	27.90
Between 500 and 1000 students	37.01	25.47	38.84
Less than 500 students	23.62	24.06	33.26
<b>Total</b>	<b>100</b>	<b>100</b>	<b>100</b>
<b>Gender of the school</b>			
Boys or Girls school	0	11.79	11.67
Coeducational school	100	88.21	88.33
<b>Total</b>	<b>100</b>	<b>100</b>	<b>100</b>

Source: Own calculation based on Saber 11° and C-600.

**Table A1.8** DID Continuous Specification: Exponential

Estimated values of  $\delta^r$  from

$$Y_{it} = \sum_{\tau=2003}^{2008} \delta^r T_i \cdot 1(\tau = t) + \beta_1 A_i + \eta X_{it} + \gamma_i + \gamma_\tau + e_{it}$$


---

**Exponential Specification:** For a school of distance  $d_i$  from a library,  $T_i = \frac{R1}{d_i} - 1$  if  $d_i \leq R1$  and  $T_i = 0$  if  $d_i > R1$

Distance Def	2003	2004	2005	2006	2007	2008
R1=1500	0.03 (0.07)	-0.02 (0.05)	0.06 (0.08)	0.05 (0.09)	0.04 (0.10)	0.06 (0.09)
R1=2000	0.02 (0.04)	-0.01 (0.03)	0.03 (0.05)	0.03 (0.06)	0.02 (0.07)	0.04 (0.06)
R1=2500	0.01 (0.03)	-0.01 (0.02)	0.02 (0.04)	0.02 (0.04)	0.01 (0.05)	0.03 (0.04)
R1=3000	0.01 (0.03)	-0.00 (0.02)	0.02 (0.03)	0.01 (0.03)	0.01 (0.04)	0.03 (0.04)
R1=3500	0.01 (0.02)	-0.00 (0.02)	0.01 (0.03)	0.01 (0.03)	0.01 (0.03)	0.02 (0.03)

R2=3500. Standard errors clustered by locality in parentheses. Significance level: \* 90%, \*\* 95%, \*\*\* 99%. Source: Own calculation based on Saber 11° and C-600.

Table A1.9 DiD Discrete

Estimated values of $\delta^r$ from						
$Y_{it} = \sum_{\tau=2003}^{2008} \delta^\tau T_i \cdot 1(\tau = t) + \beta_i A_i + \eta X_{it} + \gamma_i + \gamma_t + e_{it}$						
<b>A. Specification I:</b> Schools between 0 and R1 meters are treated, $T_i = 1$ , and from R1 to R2 meters are controls, $T_i = 0$						
Distance Def	2003	2004	2005	2006	2007	2008
R1=750	0.04 (0.19)	0.02 (0.17)	0.21 (0.22)	0.12 (0.26)	0.10 (0.28)	0.15 (0.26)
R1=1000	0.10 (0.13)	0.04 (0.11)	-0.02 (0.14)	0.04 (0.16)	-0.03 (0.16)	0.02 (0.16)
R1=1250	0.10 (0.09)	0.06 (0.08)	-0.02 (0.10)	0.03 (0.11)	-0.02 (0.12)	0.04 (0.11)
R1=1500	0.03 (0.07)	0.05 (0.06)	-0.03 (0.07)	0.04 (0.07)	-0.02 (0.08)	0.05 (0.08)
R1=1750	0.02 (0.06)	-0.00 (0.06)	-0.06 (0.06)	-0.00 (0.07)	-0.04 (0.07)	-0.03 (0.07)
R1=2000	0.02 (0.06)	-0.01 (0.05)	-0.02 (0.06)	-0.05 (0.06)	-0.05 (0.06)	-0.02 (0.06)
<b>B. Specification II:</b> Schools between 0 and R1 meters are treated, $T_i = 1$ , and from R3 to R2 meters are controls, $T_i = 0$						
Distance Def	2003	2004	2005	2006	2007	2008
R1=750	0.06 (0.19)	0.02 (0.17)	0.18 (0.22)	0.09 (0.26)	0.07 (0.28)	0.14 (0.26)
R1=1000	0.10 (0.13)	0.03 (0.11)	-0.03 (0.14)	0.01 (0.15)	-0.05 (0.16)	0.02 (0.16)
R1=1250	0.10 (0.09)	0.05 (0.08)	-0.03 (0.10)	-0.00 (0.11)	-0.04 (0.12)	0.03 (0.11)
R1=1500	0.03 (0.07)	0.03 (0.06)	-0.03 (0.07)	0.01 (0.07)	-0.03 (0.08)	0.04 (0.08)
R1=1750	0.02 (0.06)	-0.01 (0.06)	-0.05 (0.06)	-0.03 (0.07)	-0.05 (0.07)	0.03 (0.07)
R1=2000	0.02 (0.06)	-0.01 (0.05)	-0.02 (0.06)	-0.05 (0.06)	-0.05 (0.06)	0.02 (0.06)

R2=3500, R3=2000. Standard errors clustered by locality in parentheses. Significance level: \* 90%, \*\* 95%, \*\*\* 99%.

Source: Own calculation based on Saber 11<sup>o</sup> and C-600.

Table A1.10 Balance Status after Matching: Discrete I

Variables	RI=750			RI=1000			RI=1250		
	General	Matched	B.R.	General	Matched	B.R.	General	Matched	B.R.
Proportion of teachers with graduate studies	-0.05	-0.02	63.3	-0.01	0.01	48.7	-0.00	-0.01	-130.0
Pupil-Teacher Ratio	0.01*	0.01	45.7	0.01	0.00	70.5	0.00	-0.00	31.6
Female-Teacher ratio	-0.07*	-0.04	50.9	-0.03	0.00	89.2	-0.02	0.00	97.3
11 Grade Girls-students ratio	-0.02	-0.03	-66.7	0.03	0.00	86.0	0.02	0.01	39.3
11 Grade Students	11.25	-25.70	-240.0	3.75	-10.38	-290.0	13.15	-4.20	68.3
Total Students	85.47	-243.67	-290.0	58.98	-57.09	-20.0	118.48	-46.13	60.9
Girls-Students ratio	-0.04	-0.03	16.2	0.02	0.01	52.0	0.01	0.02	-89.3
Public School	0.17	0.13	19.7	0.02	-0.00	95.5	-0.09	0.03	72.5
School day: morning	-0.04	0.04	-6.4	-0.01	0.12	-1400.0	0.07	0.04	49.6
School day: complete	0.19	0.09	57.2	0.04	-0.06	-41.0	-0.03	-0.01	73.5
Built area per student	-0.25	-0.27	-11.8	-0.20	-0.04	81.6	2.80**	1.18	58.6
Classrooms area per student	0.09	-0.10	-2.8	0.20	0.11	51.9	0.51	0.48	18.1
Sports area per student	0.77	-0.27	65.3	1.32	0.35	73.4	1.99**	0.70	71.1
Has a library (C100)	-0.13*	0.05	24.1	-0.11**	0.05	13.2	-0.05	0.04	-45.7
Avg Std Test Score: 2000	-0.07	-0.06	16.9	-0.11	-0.11	8.4	-0.10	-0.05	47.0
Avg Std Test Score: 2001	0.22	0.05	75.0	0.18	-0.04	79.8	0.09	-0.04	53.6
Avg Std Test Score: 2002	0.09	0.02	76.0	0.04	-0.02	37.3	-0.08	-0.01	85.1

T: Treated, C: Control, SC: Synthetic Control (weighted control group)

(Continued)

Table A1.10 Balance Status after Matching: Discrete I

Year Variables	RI=1500			RI=1750			RI=2000		
	General	Matched	B.R.	General	Matched	B.R.	General	Matched	B.R.
Proportion of teachers with graduate studies	0.02	0.00	83.7	0.04	0.00	95.6	0.05	0.00	93.0
Pupil-Teacher Ratio	-0.00	-0.00	42.8	0.00	0.00	43.7	0.00	0.00	21.4
Girls-teacher ratio	-0.01	0.00	72.5	-0.01	0.01	31.0	-0.02	0.00	88.2
11 Grade Girls-students ratio	0.01	-0.01	-14.3	-0.00	-0.01	-110.0	0.00	0.00	-8.2
11 Grade Students	23.25**	-0.83	96.7	21.97**	0.55	97.5	23.11***	-3.76	79.4
Total Students	234.79***	30.50	87.8	191.24**	28.81	84.3	208.51***	-43.47	75.0
Girls-students ratio	0.00	-0.00	-73.6	0.00	0.00	76.8	0.01	0.01	-27.0
Public School	-0.04	0.00	96.6	-0.04	-0.01	70.4	-0.04	-0.01	80.8
School day: morning	0.02	0.03	-110.0	0.00	-0.00	-33.9	0.01	0.02	-52.1
School day: complete	0.00	-0.02	-4800.0	0.03	-0.01	58.9	0.01	-0.00	83.1
Built area per student	1.01	0.65	42.1	0.56	0.74	-22.8	1.23	1.16	12.2
Classrooms area per student	0.19	0.05	76.4	0.05	0.08	-46.3	0.25	0.09	68.9
Sports area per student	0.89	0.17	82.4	0.53	0.64	-24.0	0.89	0.03	96.8
Has a library (C100)	-0.01	-0.00	95.5	-0.00	-0.02	-19000.0	0.02	-0.00	58.4
Avg Std Test Score: 2000	-0.10	-0.00	95.4	-0.09	-0.02	77.6	-0.00	0.00	-68.5
Avg Std Test Score: 2001	0.06	-0.00	93.2	0.04	-0.02	21.2	0.09	0.02	79.3
Avg Std Test Score: 2002	-0.09	0.01	90.5	-0.09	0.00	98.8	-0.01	0.01	-25.9

Significance level for t-tests for equality of means: \* 90%, \*\* 95%, \*\*\* 99%  
Source: Own calculation based on Saber 11° and C-600.



Table A1.11 Balance Status after Matching: Discreta III

Variables	RI=750			RI=1000			RI=1250		
	General	Matched	B.R.	General	Matched	B.R.	General	Matched	B.R.
Proportion of teachers with graduate studies	-0.03	-0.02	46.5	0.01	0.04	-320.0	0.02	-0.01	17.0
Pupil-Teacher Ratio	0.01*	0.00	85.4	0.01	0.00	71.5	0.00	0.00	92.7
Female-Teacher ratio	-0.08*	-0.02	75.9	-0.03	-0.00	96.2	-0.02	-0.00	92.1
11 Grade Girls-students ratio	-0.02	-0.03	-130.0	0.03	0.03	8.7	0.02	-0.00	82.3
11 Grade Students	18.99	-13.30	-20.4	11.74	-4.50	42.8	19.55*	-2.17	87.8
Total Students	156.17	-113.86	-17.0	128.49	-37.55	63.4	176.24*	-20.45	87.3
Girls-Students ratio	-0.04	-0.02	35.9	0.02	0.03	-79.4	0.01	-0.00	82.2
Public School	0.14	0.20	-43.9	0.00	-0.00	46.0	-0.10	0.02	76.5
School day: morning	-0.03	0.07	-130.0	-0.00	0.10	-3900.0	0.06	0.04	43.0
School day: complete	0.19	0.05	74.6	0.04	-0.02	43.5	-0.02	0.03	-17.4
Built area per student	0.20	-2.08	-650.0	0.26	0.31	-9.3	2.82**	0.96	66.9
Classrooms area per student	0.18	-0.43	-100.0	0.27	-0.12	65.3	0.52	0.34	45.8
Sports area per student	1.05	0.34	64.5	1.51	0.07	95.7	2.02**	0.37	85.3
Has a library (C100)	-0.12	0.04	21.3	-0.09	0.05	-9.5	-0.04	-0.02	31.1
Avg Std Test Score: 2000	-0.06	-0.05	28.0	-0.10	0.02	84.2	-0.08	-0.02	68.0
Avg Std Test Score: 2001	0.24	0.05	80.1	0.20	0.08	57.0	0.11	-0.04	56.3
Avg Std Test Score: 2002	0.08	-0.01	88.2	0.03	0.12	-270.0	-0.07	-0.01	88.4

(Continued)

Table A1.11 Balance Status after Matching: Discreta II

Year	RI=1500			RI=1750			RI=2000		
	General	Matched	B.R.	General	Matched	B.R.	General	Matched	B.R.
Proportion of teachers with graduate studies	0.03	0.02	42.8	0.05	0.01	88.3	0.05	0.00	93.0
Pupil-Teacher Ratio	-0.00	-0.00	-36.2	0.00	-0.00	3.0	0.00	0.00	21.4
Girls-teacher ratio	-0.01	0.00	99.3	-0.01	0.01	17.9	-0.02	0.00	88.2
11 Grade Girls-students ratio	0.01	-0.01	-9.0	-0.00	-0.01	-540.0	0.00	0.00	-8.2
11 Grade Students	26.19***	6.26	74.4	24.34***	-2.17	90.2	23.11***	-3.76	79.4
Total Students	255.52***	74.85	68.7	214.53***	6.53	96.6	208.51***	-43.47	75.0
Girls-students ratio	0.00	-0.00	-33.1	0.01	-0.00	49.3	0.01	0.01	-27.0
Public School	-0.05	-0.04	18.4	-0.04	-0.01	75.6	-0.04	-0.01	80.8
School day: morning	0.02	0.05	-210.0	0.01	0.03	-280.0	0.01	0.02	-52.1
School day: complete	0.01	-0.05	-870.0	0.03	-0.01	41.1	0.01	-0.00	83.1
Built area per student	1.22	0.65	53.6	0.85	0.85	12.3	1.23	1.16	12.2
Classrooms area per student	0.24	0.02	92.4	0.13	0.04	74.9	0.25	0.09	68.9
Sports area per student	1.01	0.06	94.2	0.70	0.42	41.7	0.89	0.03	96.8
Has a library (C100)	-0.00	-0.04	-2900.0	0.01	-0.01	-170.0	0.02	-0.00	58.4
Avg Std Test Score: 2000	-0.08	0.02	62.7	-0.07	-0.04	25.2	-0.00	0.00	-68.5
Avg Std Test Score: 2001	0.08	0.01	86.8	0.06	-0.03	39.4	0.09	0.02	79.3
Avg Std Test Score: 2002	-0.08	0.03	47.4	-0.07	-0.03	42.5	-0.01	0.01	-25.9

Significance level for t-tests for equality of means: \* 90%, \*\* 95%, \*\*\* 99%  
 Source: Own calculation based on Saber 11° and C-600.

Table A1.12 DiD Discrete after Matching Including School-Specific Trends

Estimated values of  $\delta^r$  from

$$Y_{it} = \sum_{\tau=2003}^{2008} \delta^r T_i \cdot 1(\tau = t) + \beta_i A_i + \eta X_{it} + \gamma_i + \gamma_t + \omega_i t \cdot \gamma_i + e_{it}$$

**A. Specification I:** Schools between 0 and R1 meters are treated,  $T_i = 1$ , and from R1 to R2 meters are controls,  $T_i = 0$

Distance Def	2003	2004	2005	2006	2007	2008
R1=750	-0.23 (0.15)	-0.22 (0.25)	0.02 (0.24)	-0.14 (0.38)	-0.29 (0.51)	-0.09 (0.61)
R1=1000	-0.12 (0.11)	-0.16 (0.15)	-0.20 (0.20)	-0.24 (0.23)	-0.41 (0.28)	-0.24 (0.35)
R1=1250	-0.08 (0.11)	-0.11 (0.16)	-0.22 (0.22)	-0.21 (0.26)	-0.33 (0.29)	-0.27 (0.35)
R1=1500	-0.06 (0.09)	-0.04 (0.11)	-0.14 (0.15)	-0.10 (0.18)	-0.19 (0.21)	-0.13 (0.24)
R1=1750	-0.06 (0.09)	-0.06 (0.11)	-0.16 (0.14)	-0.13 (0.17)	-0.19 (0.19)	-0.14 (0.23)
R1=2000	-0.06 (0.08)	-0.04 (0.11)	-0.08 (0.14)	-0.12 (0.17)	-0.14 (0.20)	-0.05 (0.23)

**B. Specification II:** Schools between 0 and R1 meters are treated,  $T_i = 1$ , and from R3 to R2 meters are controls,  $T_i = 0$

Distance Def	2003	2004	2005	2006	2007	2008
R1=750	-0.27 (0.20)	-0.15 (0.32)	0.01 (0.29)	-0.09 (0.46)	-0.24 (0.59)	-0.01 (0.71)
R1=1000	-0.20* (0.11)	-0.24 (0.18)	-0.39 (0.26)	-0.45 (0.28)	-0.56* (0.32)	-0.42 (0.39)
R1=1250	-0.05 (0.13)	-0.09 (0.17)	-0.25 (0.23)	-0.23 (0.27)	-0.35 (0.30)	-0.28 (0.34)
R1=1500	-0.07 (0.09)	-0.04 (0.12)	-0.18 (0.17)	-0.14 (0.20)	-0.22 (0.22)	-0.14 (0.25)
R1=1750	-0.07 (0.09)	-0.04 (0.11)	-0.13 (0.15)	-0.11 (0.18)	-0.18 (0.21)	-0.12 (0.24)
R1=2000	-0.06 (0.08)	-0.04 (0.11)	-0.08 (0.14)	-0.12 (0.17)	-0.14 (0.20)	-0.05 (0.23)

R2=3500, R3=2000. Standard errors clustered by locality in parentheses. Significance level: \* 90%, \*\* 95%, \*\*\* 99%. Source: Own calculation based on Saber 11° and C-600.

Table A1.13 BO-DD Discreta

Blinder-Oaxaca decomposition of the treatment effect:  $\delta = \Delta_0 + \Delta_x$

$\delta$  : Total impact

$\Delta_x$  : Impact due to variation on covariates

$\Delta_0$  : Impact due to other channels

	Treated/Controls	2003	2004	2005	2006	2007	2008
R1=750 10/182	$\delta$	0.2008 (0.3981)	0.2252 (0.3644)	0.4293 (0.4220)	0.3343 (0.5306)	0.1444 (0.4816)	0.1726 (0.4339)
	$\Delta_0$	0.0310 (0.4573)	0.2346 (0.3555)	0.5356 (0.4093)	0.4517 (0.4972)	0.3101 (0.4399)	0.3223 (0.4144)
	$\Delta_x$	0.1698 (0.1710)	-0.0094 (0.0985)	-0.1063 (0.1055)	-0.1174 (0.1137)	-0.1657 (0.1078)	-0.1497 (0.1005)
R1=1000 19/173	$\delta$	0.1626 (0.2353)	0.1460 (0.2138)	0.0851 (0.2474)	0.1524 (0.2810)	-0.0009 (0.2593)	0.0465 (0.2562)
	$\Delta_0$	0.1253 (0.2616)	0.1787 (0.2075)	0.2083 (0.2410)	0.2776 (0.2625)	0.1597 (0.2414)	0.2296 (0.2469)
	$\Delta_x$	0.0373 (0.1308)	-0.0327 (0.1149)	-0.1232 (0.1120)	-0.1252 (0.1241)	-0.1605 (0.1073)	-0.1831* (0.1080)
R1=1250 27/165	$\delta$	0.1308 (0.1595)	0.1056 (0.1489)	0.0286 (0.1789)	0.0639 (0.1891)	-0.0229 (0.1787)	0.0463 (0.1774)
	$\Delta_0$	0.1360 (0.1895)	0.1890 (0.1644)	0.1740 (0.1890)	0.2745 (0.2016)	0.2032 (0.1754)	0.2613 (0.1850)
	$\Delta_x$	-0.0051 (0.1091)	-0.0833 (0.1158)	-0.1455 (0.1079)	-0.2107 (0.1331)	-0.2262** (0.1118)	-0.2149* (0.1173)
R1=1500 44/148	$\delta$	0.0560 (0.1194)	0.0840 (0.1145)	0.0115 (0.1239)	0.0790 (0.1245)	-0.0202 (0.1430)	0.0637 (0.1270)
	$\Delta_0$	0.0535 (0.1390)	0.1459 (0.1278)	0.1349 (0.1330)	0.2530 (0.1586)	0.1201 (0.1466)	0.1831 (0.1568)
	$\Delta_x$	0.0026 (0.0990)	-0.0619 (0.1230)	-0.1235 (0.1052)	-0.1740 (0.1341)	-0.1403 (0.1155)	-0.1193 (0.1303)
R1=1750 52/140	$\delta$	0.0437 (0.1081)	0.0320 (0.1109)	-0.0172 (0.1106)	0.0379 (0.1213)	-0.0237 (0.1337)	0.0574 (0.1144)
	$\Delta_0$	0.0515 (0.1211)	0.0611 (0.1157)	0.0967 (0.1204)	0.2114 (0.1469)	0.1084 (0.1370)	0.1779 (0.1436)
	$\Delta_x$	-0.0078 (0.0906)	-0.0291 (0.1131)	-0.1139 (0.0958)	-0.1735 (0.1302)	-0.1321 (0.1204)	-0.1206 (0.1273)
R1=2000 70/122	$\delta$	0.0208 (0.1074)	-0.0132 (0.1079)	-0.0178 (0.1071)	-0.0446 (0.1175)	-0.0722 (0.1178)	0.0125 (0.1107)
	$\Delta_0$	0.0230 (0.1206)	0.0657 (0.1033)	0.0530 (0.1103)	0.1115 (0.1387)	0.0862 (0.1334)	0.1253 (0.1350)
	$\Delta_x$	-0.0021 (0.0853)	-0.0789 (0.1138)	-0.0708 (0.0952)	-0.1561 (0.1292)	-0.1584 (0.1195)	-0.1128 (0.1343)

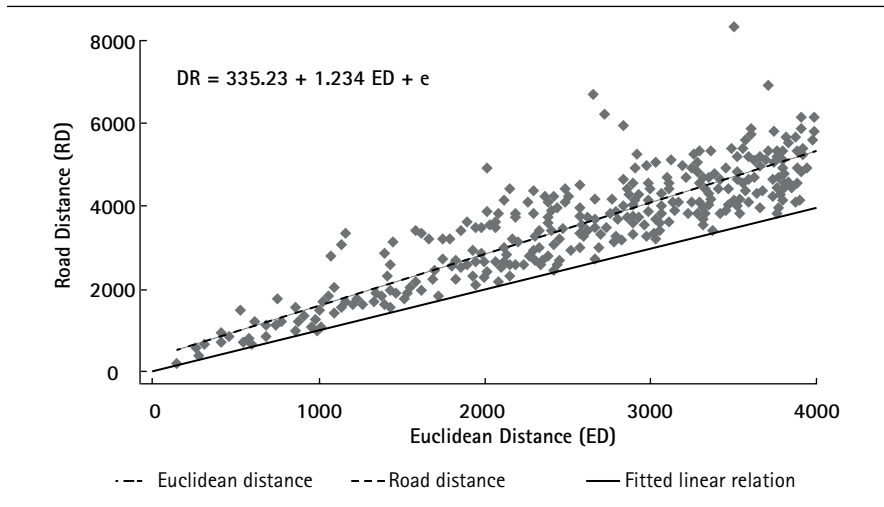
R2=3500. Clusters by locality standard errors in parentheses

\*\*\* p < 0.01, \*\* p < 0.05, \* p < 0.1

Source: Own calculation based on Saber 11° and C-600.

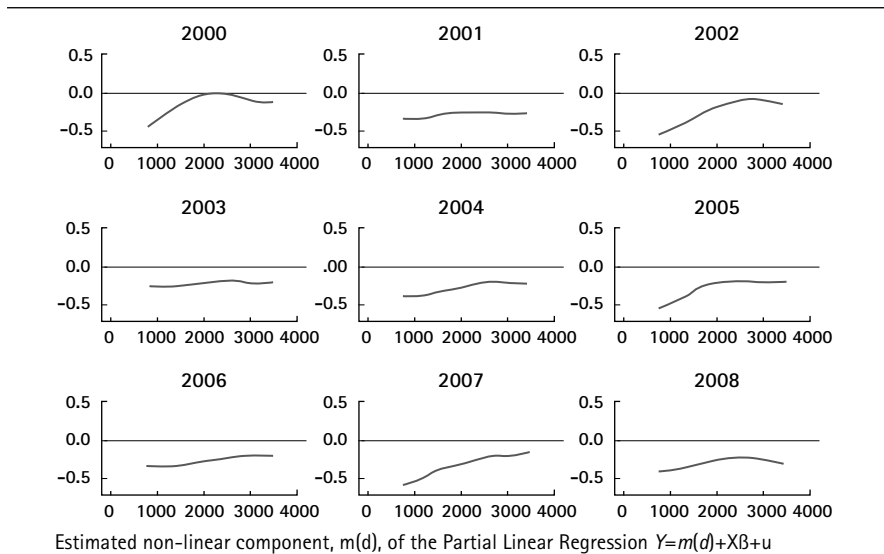
## Annex 2. Figures

Figure A2.1 Euclidean vs Road Distances



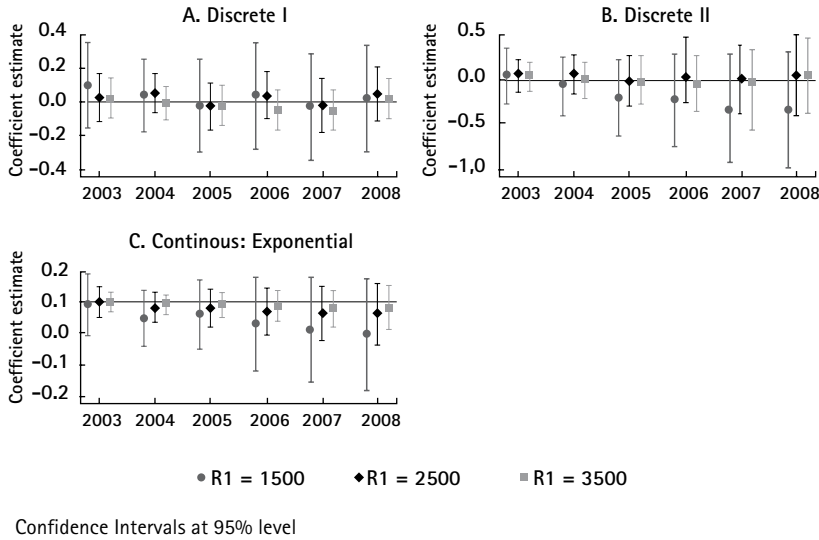
Source: Own calculations base on OSM roads network

Figure A2.2 Distance and Scores Relationship



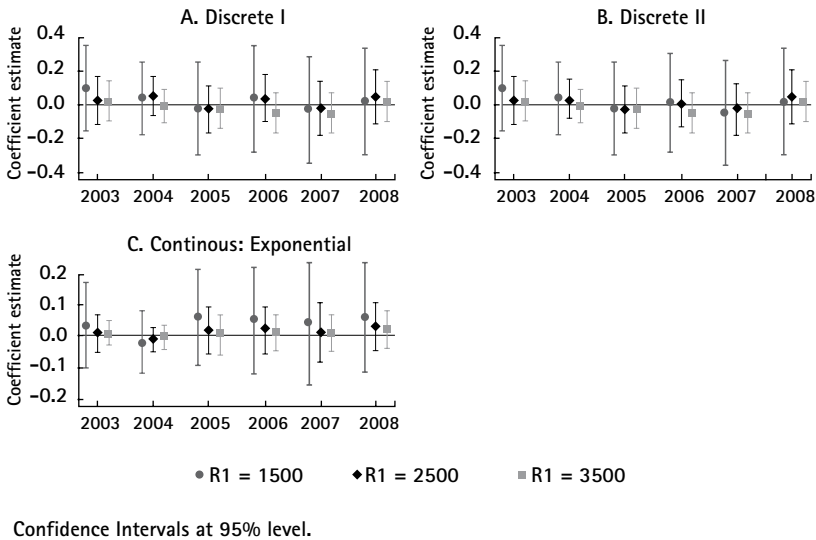
Source: Own calculation based on Saber 11°.

Figure A2.3 Euclidean Distance Estimators with School-Specific Trends



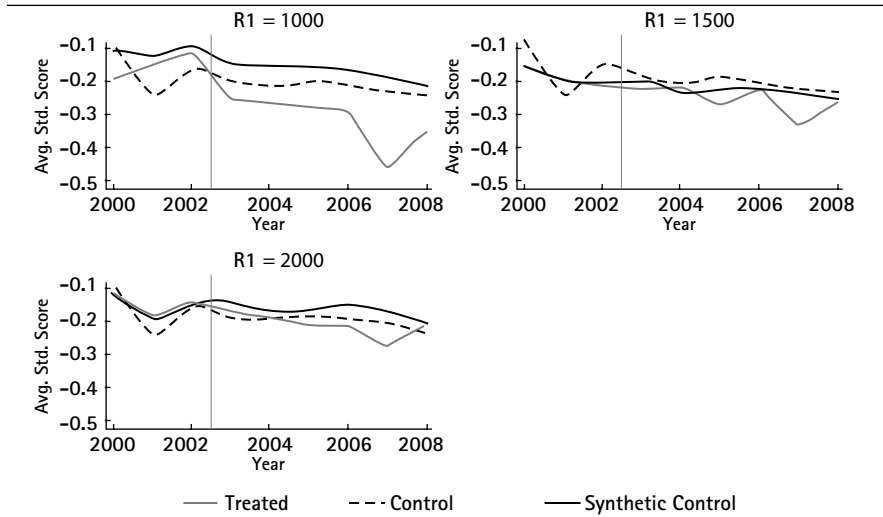
Source: Own calculation based on C-600 and Saber 11°.

Figure A2.4 Road Distance Estimators



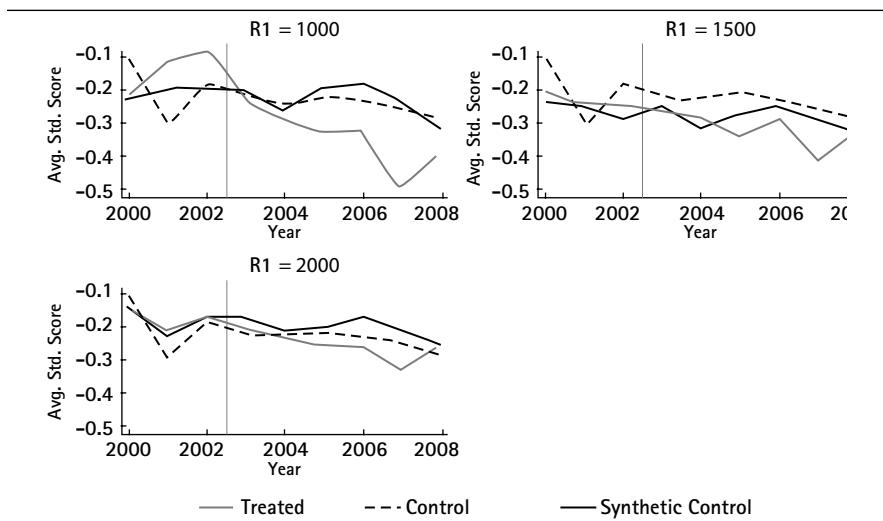
Source: Own calculation based on C-600 and Saber 11°.

Figure A2.5 Matching test Scores Evolution: Discrete I



Source: Own calculation based on C-600 and Saber 11°.

Figure A2.6 Matching test Scores Evolution: Discrete II



Source: Own calculation based on C-600 and Saber 11°.

## Appendix. Oaxaca–Blinder and DiD

Here we propose a new identification strategy that mix the advantages of Blinder Oaxaca decomposition with the DiD specification. The Blinder (1973) and Oaxaca (1973) procedure allows to decompose the difference of a variable  $y$  between two groups,  $\delta = E[y | T = 1] - E[y | T = 0]$ , by the difference on observed characteristics  $x$ ,  $\Delta_x$ , and a difference that is not related to them  $\Delta_0$ . Here we assume a linear relationship between observed characteristics  $x$  and the outcome  $y$  which can be specific to the group  $T$ .

$$y = \beta_0 + \beta_1 x + \beta_2 T + \beta_3 T \cdot x + e_2$$

If we impose  $Ee_2 | T = 1 = Ee_2 | T = 0$ , the difference  $\delta$  can be expressed on terms of the difference on  $x$  between the two groups and a remainder.

$$\begin{aligned} \delta &= E[y | T = 1] - E[y | T = 0] \\ &= [\beta_0 + \beta_2 + (\beta_1 + \beta_3)E[x | T = 1]] - [\beta_0 + \beta_1 E[x | T = 0]] \\ &= \beta_2 + (\beta_1 + \beta_3)E[x | T = 1] - \beta_1 E[x | T = 0] \\ &= \{\beta_2 + \beta_3 E[x | T = 1]\} + \{\beta_1 (E[x | T = 1] - E[x | T = 0])\} \\ &= \{\Delta_0\} + \{\Delta_x\} \end{aligned}$$

We define  $\Delta_x = \beta_1 (E[x | T = 1] - E[x | T = 0])$ , as the difference for being part of  $T = 1$  and not of  $T = 0$  on  $x$ . The other term,  $\Delta_0 = \beta_2 + \beta_3 E[x | T = 1]$ , reflects the difference on  $y$  which is not explained due to the difference on  $x$ . In empirical labour economics, these former term was usually interpreted as the 'discrimination' for being part of  $T = 1$ . Under the framework of treatment effects literature, where  $T$  is a treatment that has a heterogeneous effect according to  $x$ , so the 'unexplained' component is an average treatment on the treated (Fortin, Lemieux and Firpo 2011).

We propose a Difference-in-Differences (DiD) analogue of the decomposition, where we can understand which part of the variation is explained by the impact on an observed channel  $x$ . In the case of our program, we would



like to understand which part of the effect is due to an enhancement of the results of schools via the increase on certain inputs, and what is due to a general impact that is not related to them. To the best of our knowledge this is the first paper that implements this decomposition.

Let's assume that we can observe two periods,  $A \in \{0,1\}$ . Given it, we define the average treatment on the treated estimator:

$$\delta_x = (E[y | T = 1, A = 1] - E[y | T = 0, A = 1]) - (E[y | T = 1, A = 0] - E[y | T = 0, A = 0])$$

This is the classical DiD estimator under the usual parallel trends assumption. It could be retrieved by using the traditional specification,

$$y = \eta_0 + \eta_1 T + \eta_2 A + \delta T \cdot A + \varepsilon$$

Now, let's assume that part of this impact is due to a variation on a particular variable  $x$  that is affected by the treatment. Our decomposition is able to decompose the treatment effect of  $T$  on  $Y$  between the impact on the observed channel,  $\Delta_x$  and the impact via other channels,  $\Delta_0$ . It can be implemented using the following linear equation:

$$y = \alpha_0 + \alpha_1 T + \alpha_2 A + \alpha_3 x + \alpha_4 x \cdot A + \alpha_6 T \cdot A + \alpha_7 x \cdot T \cdot A + u$$

Given that

$$E[y | T = 0, A = 0] = \alpha_0 + \alpha_3 E[x | T = 0, A = 0]$$

$$E[y | T = 1, A = 0] = \alpha_0 + \alpha_1 (\alpha_3 + \alpha_4) E[x | T = 1, A = 0]$$

$$E[y | T = 0, A = 1] = \alpha_0 + \alpha_2 (\alpha_3 + \alpha_5) E[x | T = 0, A = 1]$$

$$E[y | T = 1, A = 1] = \alpha_0 + \alpha_1 + \alpha_2 + \alpha_6 + (\alpha_3 + \alpha_4 + \alpha_5 + \alpha_7) E[x | T = 1, A = 1]$$

The impact  $\delta$  is decomposed between the variation on  $x$  that is correlated with the treatment implementation,  $\Delta_x$ , and the variation that is explained due to other channels,  $\Delta_0$ .

$$\begin{aligned} \delta &= ((\alpha_0 + \alpha_1 + \alpha_2 + \alpha_6 + (\alpha_3 + \alpha_4 + \alpha_5 + \alpha_7)E[x | T = 1, A = 1]) \\ &\quad - (\alpha_0 + \alpha_2 + (\alpha_3 + \alpha_5)E[x | T = 1, A = 1])) \\ &\quad - ((\alpha_0 + \alpha_1 + (\alpha_3 + \alpha_4)E[x | T = 1, A = 0]) \\ &\quad - (\alpha_0 + \alpha_3 E[x | T = 0, A = 0])) \\ \delta &= \alpha_6 + (\alpha_4 + \alpha_4 + \alpha_7)E[x | T = 1, A = 1] \\ &\quad - \alpha_5 E[x | T = 0, A = 1] - \alpha_4 E[x | T = 1, A = 0] \end{aligned}$$

Hence, the impact on  $Y$  due to  $T$  that can be explained by the impact of  $T$  on  $X$  is:

$$\begin{aligned} \Delta_x &= \alpha_3 (E[x | T = 1, A = 1] - E[x | T = 0, A = 1]) \\ &\quad - (E[x | T = 1, A = 0] - E[x | T = 0, A = 0]) \end{aligned}$$

And the remainder variation

$$\begin{aligned} \Delta_0 &= \alpha_6 + (\alpha_4 + \alpha_5 + \alpha_7)E[x | T = 1, A = 1] - \alpha_5 E[x | T = 0, A = 1] \\ &\quad - \alpha_4 E[x | T = 1, A = 0] \end{aligned}$$