

IDB WORKING PAPER SERIES N° IDB-WP-821

# Do Tests Applied to Teachers Predict their Effectiveness?

Yyannú Cruz-Aguayo  
Pablo Ibararán  
Norbert Schady

Inter-American Development Bank  
Social Sector - SCL/SCL

May, 2017

# Do Tests Applied to Teachers Predict their Effectiveness?

Yyannú Cruz-Aguayo  
Pablo Ibararán  
Norbert Schady

Cataloging-in-Publication data provided by the  
Inter-American Development Bank

Felipe Herrera Library  
Cruz Aguayo, Yyannú.

Do tests applied to teachers predict their effectiveness? / Yyannú Cruz-Aguayo, Pablo  
Ibarraran, Norbert Schady.

p. cm. — (IDB Working Paper Series ; 821)

Includes bibliographic references.

1. Teacher effectiveness-Ecuador-Evaluation. 2. Teachers-Rating of-Ecuador. 3.  
Teaching-Ecuador-Evaluation. I. Ibarraran, Pablo. II. Schady, Norbert Rüdiger, 1967-  
III. Inter-American Development Bank. Education Division. IV. Title. V. Series.  
IDB-WP-821

<http://www.iadb.org>

Copyright © 2017 Inter-American Development Bank. This work is licensed under a Creative Commons IGO 3.0 Attribution-NonCommercial-NoDerivatives (CC-IGO BY-NC-ND 3.0 IGO) license (<http://creativecommons.org/licenses/by-nc-nd/3.0/igo/legalcode>) and may be reproduced with attribution to the IDB and for any non-commercial purpose, as provided below. No derivative work is allowed.

Any dispute related to the use of the works of the IDB that cannot be settled amicably shall be submitted to arbitration pursuant to the UNCITRAL rules. The use of the IDB's name for any purpose other than for attribution, and the use of IDB's logo shall be subject to a separate written license agreement between the IDB and the user and is not authorized as part of this CC-IGO license.

Following a peer review process, and with previous written consent by the Inter-American Development Bank (IDB), a revised version of this work may also be reproduced in any academic journal, including those indexed by the American Economic Association's EconLit, provided that the IDB is credited and that the author(s) receive no income from the publication. Therefore, the restriction to receive income from such publication shall only extend to the publication's author(s). With regard to such restriction, in case of any inconsistency between the Creative Commons IGO 3.0 Attribution-NonCommercial-NoDerivatives license and these statements, the latter shall prevail.

Note that link provided above includes additional terms and conditions of the license.

The opinions expressed in this publication are those of the authors and do not necessarily reflect the views of the Inter-American Development Bank, its Board of Directors, or the countries they represent.



All three authors are with the Inter-American Development Bank.

Their e-mails are [yyannuc@iadb.org](mailto:yyannuc@iadb.org), [pibarraran@iadb.org](mailto:pibarraran@iadb.org), and [norberts@iadb.org](mailto:norberts@iadb.org).

## **Do Tests Applied to Teachers Predict their Effectiveness?**

Yyannú Cruz-Aguayo  
Pablo Ibararán  
Norbert Schady

May 8, 2017<sup>1</sup>

### **Abstract**

Teachers vary considerably in their effectiveness, but identifying teacher characteristics that predict their impact on learning outcomes has been elusive. We analyze a teacher evaluation that is used to make teacher tenure decisions in Ecuador. The evaluation includes a written test, a demonstration class, and a points system that gives higher scores to teachers with more experience, degrees, and in-service training. We find no evidence that children taught by teachers with higher scores on the evaluation learn more. Our estimates are very precise: We can rule out that teachers with one-standard deviation higher evaluation scores raise child test scores in math by 0.03 standard deviations or more, and language scores by 0.02 standard deviations or more. We conclude that the effort that is being placed by policy-makers in Latin America to design and “improve” teacher tests is unlikely to result in large improvements in child learning.

**JEL Codes:** I25, I28

**Keywords:** Teachers, evaluation, test scores.

---

<sup>1</sup> All three authors are with the Inter-American Development Bank. Their e-mails are [yyannuc@iadb.org](mailto:yyannuc@iadb.org), [pibarraran@iadb.org](mailto:pibarraran@iadb.org), and [norberts@iadb.org](mailto:norberts@iadb.org). We thank officials in the Ministry of Education in Ecuador for making data available to us, Fiorella Benedetti and Jorge Luis Castañeda for outstanding research assistance, and Pedro Carneiro for his comments.

## 1. Introduction

It is generally believed that teachers are the most important input into the production of learning within schools. Having a better teacher has long-term consequences for a variety of outcomes in adulthood, including college attendance, savings, and wages (Chetty et al. 2011; Chetty et al. 2014a). Research from the United States has credibly established that teachers—even teachers in the same school, teaching observationally equivalent students—vary a great deal in their effectiveness (Chetty et al. 2014b; Hanushek and Rivkin 2012). However, knowing that some teachers produce more learning than others provides no guidance on the attributes of effective teachers.

In this note, we analyze whether an evaluation that is used to determine which teachers receive tenure predicts student learning in Ecuador, a middle-income country in South America. In earlier work (Araujo et al. 2016), we showed that kindergarten teachers varied in the impact they have on learning: A one standard deviation increase in classroom quality increased test scores in math and language by 0.11 standard deviations.

In Ecuador, about one in three teachers in the public sector work on a contract basis. Tenured teachers are paid substantially more than contract teachers, and have more generous health insurance and pension benefits. Being converted from a contract to a tenured teacher therefore has important benefits for the individual, and implies increases in costs for the Ministry of Education.

Every year the Ministry carries out a planning exercise that takes account of the current distribution of tenured teachers, student population and the projected growth in demand, and any projected school openings, closings or merges. Based on this exercise, additional tenured slots are assigned to some schools. These slots are school-specific, but any individual who has a teaching degree from an accredited institution can apply. Among eligible applicants to a new slot, the applicant with the highest total score on an evaluation known as the *Concurso de Méritos y Oposiciones* (henceforth, *Concurso*) is awarded the tenured position.

The *Concurso* has three components, each of which receives equal weight in an aggregate score. The first component is given by a teacher's score on a test. This test has three parts: logical-verbal reasoning, pedagogical knowledge, and subject-specific knowledge. The second component of the aggregate score is a teacher's score on a demonstration class. The third component of the aggregate *Concurso* score is given by a point system (*Méritos*) which gives higher scores to

teachers with more experience, degrees (above and beyond the basic teaching certificate), and in-service training.

We use data on a sample of children in 2<sup>nd</sup>, 3<sup>rd</sup>, and 4<sup>th</sup> grades to analyze whether children taught by teachers with higher scores on the *Concurso* have higher achievement in language and math. Our earlier work (Araujo et al. 2016) was based on a sample of kindergarten children randomly assigned to teachers within schools. The data used for the analysis in this paper do not come from an experiment. We therefore have to make stronger identifying assumptions. In practice, we use a school fixed effects strategy that compares the test scores of observationally equivalent children taught by teachers with higher or lower *Concurso* scores in the same school.

## **2. Data and identification**

Our sample consists of 4,479 children taught by 480 teachers in 240 schools. Schools were selected so each would have at least two contract teachers in 2<sup>nd</sup> through 4<sup>th</sup> grades who had taken the *Concurso* evaluation. Within each classroom, 9-10 children were randomly selected for testing. Twenty-eight percent of the children in the sample are in 2<sup>nd</sup> grade, 38 percent in 3<sup>rd</sup> grade, and 35 percent in 4<sup>th</sup> grade. To ease with interpretation, we standardize teacher scores on the *Concurso* and its components to have mean zero and unit standard deviation.

Children were tested twice, midway through the school year, and shortly before the end. We used adapted versions of the Early Grade Reading Assessment (EGRA; RTI International 2009; Gove and Cvelich 2011), and Early Grade Math Assessment (EGMA; Reubens 2009). Both tests have been applied in more than 20 developing countries. Application takes about 30 minutes per child. We standardize the four tests (EGRA and EGMA, mid- and end-year) so they have mean zero and unit standard deviation in each grade.

For 3,346 children in the sample (74.7 percent), we also have data from a household survey. This survey includes information on the years of completed schooling of mothers and fathers; whether a child attended preschool; and information on a large number of household characteristics (such as the material of the floors, wall, and ceilings in the house; whether the house is connected to the electricity, water, and sewerage system), and ownership of assets (questions for 23 assets). We aggregate the information on household characteristics and assets into a wealth composite (following, among many others, Araujo et al. 2016, and Paxson and Schady 2007, 2010, who also work with data from Ecuador).

The regressions we estimate take the following form:

$$(1) Y_{itgs} = \theta_g + \alpha_s + \beta_1 T_{gs} + \beta_2 X_{itgs} + \varepsilon_{itgs}, \quad Z=1, 2, \dots N$$

where the subscripts  $i$ ,  $t$ ,  $g$ , and  $s$  refer to individuals, teachers, grades, and schools, respectively;  $\theta_g$  is a set of grade fixed effects;  $\alpha_s$  is a set of school fixed effects;  $T_{gs}$  is the score of a teacher on the *Concurso* or its components;  $X_{itgs}$  always includes controls for child age in months, gender, and classroom size and, in some specifications, the education of mothers and fathers, household wealth, and an indicator variable for whether a child attended preschool;  $\varepsilon_{itgs}$  is the regression error term; and the regressions are run for  $z$  different tests (the mid- and end-year scores on EGRA or EGMA, and the change in scores between mid- and end-year). The coefficient of interest is  $\beta_1$ , an estimate of the extent to which children taught by teachers with higher scores on the *Concurso* have better scores in math or language. Standard errors are corrected for clustering at the school level.

The main estimation concern is purposeful sorting of children to teachers. As a check on our identification strategy, we run regressions of the variables in  $T_{gs}$  on the variables in  $X_{itgs}$ . Since the variables in  $X_{itgs}$  are predetermined, any indication that teachers with higher (or lower) test scores on the *Concurso* were assigned to children with different characteristics would raise questions about our identification strategy. Table 1 finds no evidence that this is the case, regardless whether the variables in  $X_{itgs}$  are entered one at a time or jointly. We conclude that our identification strategy is reasonable.

### 3. Results

To put things in context, we first regress the change in math or language test scores between mid- and end-year on the full set of child and household controls. We take the residuals from these regressions, calculate classroom means, and the difference in means between classrooms in the same school. Within-school differences are large: The median difference is .18 standard deviations for language, and .22 standard deviations for math. Although some of these within-school, cross-classroom differences are probably driven by sampling error or idiosyncratic classroom shocks, it seems likely that a substantial proportion reflect true differences in the quality of teachers.<sup>2</sup>

---

<sup>2</sup> In Araujo et al. (2016) and Carneiro et al. (2017) we use an Empirical Bayes estimator to shrink the estimated classroom effects. On average, the corrected estimates are about one-third smaller than those that do not correct for sampling error. We do not do this in the current paper because we are not calculating classroom (or teacher) effects.

Next, we show that the characteristics in  $X_{itgs}$  tend to be strongly correlated with test scores, as expected. Table 2 shows that girls score about 0.14 standard deviations higher on language, and 0.23 standard deviations lower in math.<sup>3</sup> Children from households of higher socioeconomic status have higher test scores in both reading and math.

Our main findings are in Table 3. We report the results from specifications with the additional controls (even-numbered columns) and without (odd-numbered columns), and six different dependent variables: the mid- and end-year test scores in language and math, and the change in test scores between mid- and end-year.<sup>4</sup>

Table 3 shows no evidence that the test score or any of its components, the score on the demonstration class, points on the *Méritos* scale, or the aggregate score on the *Concurso* predict child achievement in language or math. There are 84 coefficients in the table, and about one-half (40) are positive and one-half (44) are negative. No coefficient is significant at conventional levels. The estimates in the table are quite precise. For example, in the specifications with the additional controls, we can rule out positive associations between the total *Concurso* score and child test scores larger than 0.02 standard deviations for language, and 0.03 standard deviations for math. Consistent with the results in Table 1, the coefficients with and without the additional controls are very similar.

In sum, although there are large differences across classrooms in the same school in how much children learn, there is no evidence that children taught by teachers with better *Concurso* scores are more effective. If teachers with better scores were systematically assigned students with lower learning potential, our results could be biased towards zero. We find no evidence that this is the case. We conclude that the evaluation that was used to make tenure decisions in Ecuador did not predict how effective teachers were at increasing the language and math test scores of children in elementary school.

---

<sup>3</sup> Carneiro et al. (2017) carefully analyze the gender gap in math achievement in kindergarten through 2<sup>nd</sup> grade in Ecuador. Consistent with what we report here, they find a gender gap of 0.17 standard deviations, favoring boys, by 2<sup>nd</sup> grade.

<sup>4</sup> Results are very similar if instead we use the largest possible sample of children and do not include the additional controls. These results are available from the authors upon request.



#### 4. Discussion and conclusion

Improving the quality of education is an important priority for Latin American countries. Teachers vary considerably in their effectiveness and many Latin American children are taught by low-quality teachers (Bruns and Luque 2015). The region has low productivity, much of which can be traced to the low skills of the workforce (Pages 2010). Some analysts have argued that the low learning outcomes of Latin American students (and eventually adults) can account for much of the difference in growth rates between Latin America and East Asia (Hanushek and Woessman 2012).

Research from the United States has generally been pessimistic about the ability of various tests to identify better teachers (Hanushek and Rivkin 2006; Staiger and Rockoff 2010). We carefully analyze a teacher evaluation that included a written test, a demonstration class, and points for experience, degrees, and in-service training. Earlier work (Araujo et al. 2016) shows that teachers in Ecuador vary considerably in the impact they have on learning outcomes. In this note we show, however, that the instrument that is used to decide which teachers get tenure in Ecuador—an instrument similar to that used in many other Latin American countries—does not predict how effective a teacher is at raising math and language achievement.<sup>5</sup> Our estimates are quite precise and we can rule out even very modest effects on test scores. The results in our paper, and others in the literature, lead us to be cautious about what one can, and cannot, expect from a teacher testing regime like those that have become popular in many developing countries.<sup>6</sup>

---

<sup>5</sup> Many countries in Latin America, including Chile, Colombia, El Salvador, Mexico, Peru, and some municipalities and states in Brazil, test teachers. The tests generally include a written component and a demonstration class. Performance on this teacher evaluation is an input into decisions about hiring, tenure, in-service training, pay and, in some countries (like Chile), which teachers can get fired.

<sup>6</sup> An influential report on teachers in Latin America concludes that “mandatory certification exams are the most powerful instrument for raising teacher standards” (Bruns and Luque 2015, p. 31). Our results are not consistent with this assertion.

## References

- Araujo, M. Caridad, Pedro Carneiro, Yyannu Cruz-Aguayo, and Norbert Schady. 2016. "Teacher Quality and Learning Outcomes in Kindergarten." *Quarterly Journal of Economics* 131(3): 1415-53.
- Bruns, Barbara, and Javier Luque. 2015. *Great Teachers: How to Raise Student Learning in Latin America and the Caribbean*. Washington, D.C.: The World Bank.
- Carneiro, Pedro, Yyannu Cruz-Aguayo, and Norbert Schady. 2017. "Where the Girls Are Not: Households, Teachers, and the Gender Gap in Early Math Achievement." Unpublished manuscript, Inter-American Development Bank.
- Chetty, Raj, John Friedman, Nathaniel Hilger, Emmanuel Saez, Diane Schanzenbach, and Danny Yagan. 2011. "How Does Your Kindergarten Classroom Affect Your Earnings? Evidence from Project STAR." *Quarterly Journal of Economics* 126(4): 1593-1660.
- Chetty, Raj, John Friedman, and Jonah Rockoff. 2014a "Measuring the Impact of Teachers I: Evaluating Bias in Teacher Value-Added Estimates." *American Economic Review* 104(9): 2593-632.
- . 2014b. "Measuring the Impact of Teachers II: Teacher Value-Added and Student Outcomes in Adulthood." *American Economic Review* 104(9): 2633-679.
- Gove, Amber, and Peter Cvelich. 2011. *Early Reading: Igniting Education for All. A Report by the Early Grade Learning Community of Practice* (Revised Edition). Research Triangle Park, N.C.: Research Triangle Institute.
- Hanushek, Eric, and Steven Rivkin. 2006. "Teacher Quality." In Eric Hanushek and Finis Welch, Eds., *Handbook of the Economics of Education*, Vol. 2. Elsevier.
- , 2012. "The Distribution of Teacher Quality and Implications for Policy." *Annual Review of Economics* 4: 131-57.
- Hanushek, Eric A., and Ludger Woessmann. 2012. "Schooling, Educational Achievement, and the Latin American Growth Puzzle." *Journal of Development Economics* 99(2): 497-512.
- Pages, Carmen. 2010. *The Age of Productivity*. New York: Palgrave Macmillan.
- Paxson, Christina, and Norbert Schady. 2007. "Cognitive Development among Young Children in Ecuador: The Roles of Wealth, Health, and Parenting." *Journal of Human Resources* 42(1): 49-84.
- , 2010. "Does Money Matter? The Effects of Cash Transfers on Child Health and Development in Rural Ecuador." *Economic Development and Cultural Change* 59(1): 187-229.
- Reubens, Andrea. 2009. *Early Grade Mathematics Assessment (EGMA): A Conceptual Framework Based on Mathematics Skills Development in Children*. Research Triangle

Park, NC: Research Triangle Institute.

Staiger, Douglas O., and Jonah E. Rockoff. 2010. "Searching for Effective Teachers with Imperfect Information." *Journal of Economic Perspectives* 24(3): 97-118.

**Table 1: Associations between teacher scores on the *Concurso* and predetermined child and household characteristics**

TEACHER SCORE	CHILD OR HOUSEHOLD CHARACTERISTIC							F-test Prob > F (8)
	Age (months) (1)	Female (2)	Class Size (3)	Mother Schooling (4)	Father Schooling (5)	Wealth (6)	Attended Preschool (7)	
Test score: aggregate score	.002 (.004)	-.014 (.018)	-.008 (.015)	.001 (.003)	.003 (.003)	-.026 (.024)	-.027 (.040)	0.31
Test score: verbal-logical reasoning	-.004 (.004)	.008 (.021)	-.008 (.016)	.004 (.003)	.005 (.003)	-.010 (.019)	.025 (.036)	0.50
Test-score: pedagogical practices	-.004 (.005)	-.009 (.016)	-.006 (.013)	.002 (.003)	-.001 (.003)	-.020 (.032)	.022 (.044)	0.93
Test score: content	.003 (.004)	-.027 (.023)	-.005 (.018)	-.003 (.003)	.003 (.004)	-.027 (.017)	.103** (.041)	0.07
Demonstration class	-.000 (.004)	.019 (.031)	-.004 (.013)	.009** (.004)	.006 (.004)	.033 (.024)	.023 (.043)	0.13
Points on <i>Méritos</i> scale	.002 (.003)	-.023 (.020)	-.008 (.014)	.003 (.003)	-.002 (.003)	.014 (.017)	-.011 (.039)	0.73
Total <i>Concurso</i> score	.002 (.004)	-.030 (.024)	.007 (.014)	.004 (.004)	.002 (.003)	.014 (.017)	-.010 (.033)	0.85

*Notes:* Regressions of teacher scores on characteristics in columns (1) through (7), with each teacher score regressed on each child or household characteristic (49 separate regressions). F-test in column (8) refers to an F-test from a regression in which all the characteristics in columns (1) through (7) enter jointly. All regressions include school fixed effects. Teacher scores have been transformed to have mean zero and unit standard deviation. Sample size is 3,058 in every regression. Standard errors clustered at the school level. \*, \*\*, and \*\*\* indicate significance at the 10 percent, 5 percent, and 1 percent level, respectively.

**Table 2: Associations between child and household characteristics, and child test scores**

CHILD TEST SCORE	CHILD OR HOUSEHOLD CHARACTERISTIC						
	Age (months) (1)	Female (2)	Class Size (3)	Mother Schooling (4)	Father Schooling (5)	Wealth (6)	Attended Preschool (7)
EGRA: mid-year	.008** (.003)	.141*** (.032)	-.007 (.006)	.040*** (.004)	.036*** (.005)	.134*** (.024)	.010 (.045)
EGRA: end-year	.007** (.003)	.135*** (.032)	-.003 (.006)	.045*** (.004)	.036*** (.005)	.139*** (.024)	.073* (.043)
EGMA: mid-year	-.004 (.003)	-.242*** (.034)	-.006 (.005)	.039*** (.005)	.032*** (.005)	.112*** (.027)	-.003** (.045)
EGMA: end-year	-.002 (.003)	-.223*** (.033)	-.002 (.006)	.043*** (.005)	.030*** (.005)	.132*** (.024)	.078* (.047)

*Notes:* Regressions of child test scores on characteristics in columns (1) through (7), with each child score regressed on each child or household characteristic (28 separate regressions). All regressions include school fixed effects. Child test scores have been transformed to have mean zero and unit standard deviation. Sample size is 3,322 in every regression. Standard errors clustered at the school level. \*, \*\*, and \*\*\* indicate significance at the 10 percent, 5 percent, and 1 percent level, respectively.

**Table 3: Main results**

TEACHER SCORE ON CONCURSO OR COMPONENTS														
CHILD TEST SCORE	Total test score		Verbal-logical reasoning score		Pedagogical practices score		Content score		Demonstration class		Points on <i>Méritos</i> scale		Total <i>Concurso</i> score	
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)	(13)	(14)
EGRA: mid-year	.020 (.037)	.018 (.037)	.040 (.037)	.032 (.037)	-.002 (.029)	-.002 (.028)	.014 (.042)	.014 (.043)	.002 (.023)	-.007 (.023)	-.021 (.034)	-.021 (.034)	.000 (.030)	-.005 (.029)
EGRA: end-year	.021 (.035)	.020 (.034)	.030 (.037)	.023 (.037)	.000 (.028)	.000 (.027)	.020 (.041)	.021 (.041)	.013 (.023)	.003 (.023)	-.032 (.033)	-.033 (.033)	-.004 (.029)	-.010 (.028)
Δ EGRA score	.001 (.018)	.001 (.018)	-.009 (.016)	-.010 (.016)	.003 (.016)	.002 (.016)	.006 (.015)	.008 (.014)	.011 (.009)	.010 (.009)	-.011 (.015)	-.011 (.015)	-.004 (.012)	-.004 (.013)
EGMA: mid-year	-.046 (.031)	-.048 (.031)	-.026 (.036)	-.033 (.036)	-.042 (.027)	-.042 (.027)	-.026 (.034)	-.025 (.034)	.021 (.025)	.012 (.024)	.026 (.029)	.025 (.029)	-.003 (.029)	-.008 (.027)
EGMA: end-year	-.046 (.033)	-.046 (.033)	-.056 (.036)	-.062 (.037)	-.022 (.029)	-.022 (.030)	-.027 (.034)	-.022 (.034)	.018 (.022)	.008 (.021)	.042 (.031)	.040 (.030)	-.012 (.029)	-.018 (.027)
Δ EGMA score	-.000 (.018)	.002 (.018)	-.030 (.020)	-.030 (.020)	.021 (.015)	.020 (.015)	-.001 (.017)	.003 (.017)	-.002 (.012)	-.003 (.012)	.016 (.018)	.016 (.018)	-.009 (.017)	-.009 (.017)

*Notes:* Regressions of child test scores (mid-year, end-year, or changes between mid-year and end-year) on teacher scores on the *Concurso* or its components (84 separate regressions). All regressions include school fixed effects, class size, child age in months, and child gender. Specifications in even-numbered columns also include education of mothers, fathers, household wealth, and whether the child attended preschool. Child test scores and teacher scores on the *Concurso* have been transformed to have mean zero and unit standard deviation. Sample size is 3,058 in every regression. Standard errors clustered at the school level. \*, \*\*, and \*\*\* indicate significance at the 10 percent, 5 percent, and 1 percent level, respectively.