

“Rely (only) on the rigorous evidence” is bad advice

Lant Pritchett 

London School of Economics, London, UK

Correspondence

Lant Pritchett, London School of Economics, London, UK.

Email: lant.hks@gmail.com

Abstract

A popular interpretation of “evidence-based” decision-making is “rely (only) on the rigorous evidence” (RORE) via “systematic” reviews that: use objective protocols to generating the potentially relevant papers from the literature; then filter those to retain only the small subset that provide impact estimates regarded as “rigorous”; and summarize only those estimates. I use two sets of cross-country impact estimates—on wage gains for migrants and private school learning gains—to illustrate this seemingly attractive approach is both empirically and conceptually unsound. *First*, the cross-country variation in the rigorous estimates of impact is very large, which implies the average(s) from a systematic review is of little predictive use. In both empirical examples the “systematic review of the rigorous estimates” approach leads to *worse* predictions of impact across countries than the naïve use of country-specific ordinary least squares estimates. *Second*, I contrast a systematic review—RORE approach with an “understanding” approach—which seeks to encompass all of the available evidence into coherent understandings in forming judgments. In both examples the notion that the impact effects are constant across countries—“external

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial-NoDerivs](https://creativecommons.org/licenses/by-nc-nd/4.0/) License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

© 2023 The Author. *Review of Development Economics* published by John Wiley & Sons Ltd.

validity”—is easily rejected. Insisting on privileged reliance on “rigorous” estimates in making context-specific decisions is logically incoherent and deeply anti-scientific.

KEYWORDS

external validity, RCT

JEL CLASSIFICATION

O1, O, C9, C, C18, C1

1 | INTRODUCTION

A simple example helps set the stage. Suppose men generally self-report they are taller than they really are and hence their self-reported height is a biased estimate. An objective measurement of the true height of a relevant sample of men would create rigorous evidence for a relevant population, a distribution of unbiased estimates. In the practical task of assigning pants sizes to a group of men there are three options. One is to rely only on the rigorous evidence, ignore the biased evidence altogether, not even ask men about their height, and just assign to every man the size of pants that fit the average man. Two, one could ignore the bias altogether and just give each man the pants corresponding to their self-report. Three, one could take self-reported height and scale it back based on a rigorous estimate of the average self-report bias. Which of these produces the smallest mismatch in pant sizes depends on three empirically contingent facts: (i) the true variability across men in height, (ii) the average magnitude of the bias from self-report, and (iii) the true variability across men in the self-report bias. As the standard deviation of men's height in the USA is about 3 in., if the average self-report bias is just 1 or 2 in. (a man who is 5'9" saying he is 6 ft tall is implausible) then relying on the non-rigorous, biased, self-reported height will produce *better* results in fit across all men than using the mean of the rigorous evidence about the population. Moreover, if the variability in self-report bias variability is small, just knocking off each man's self-reported height the average self-report bias might be the best.

In development circles the rhetoric of making “evidence-based” decisions and in particular the slogan to “rely (only) on the rigorous evidence” (RORE) has apparently been persuasive. Resources have flowed into (i) dramatically expanding the number of randomized control trial (RCT)—and other “rigorous”—estimates of the causal impact or treatment effect (TE) of potential “interventions” (policies, programs, projects) and (ii) carrying out “systematic reviews” which privilege “rigorous” studies. This leads to new research papers and reviews that ignore most of the previous literature—even when this literature is voluminous and of high quality—and pretend to contribute relevant “evidence” on the basis of very few studies and without any clear theory or model. This approach actually contributes little to the understanding necessary to make informed and truly evidence-based policy/program/project decisions. Let me give two examples, one from the peer-reviewed economics literature and one from high-profile policy literature.

Burde and Linden (2013) shows that when community-based schools were made available in 13 randomly chosen villages in northwest Afghanistan the attendance of children was very much higher in those treatment than in control villages, with attendance higher by 35 percentage points for boys and 52 percentage points for girls. I do not claim this paper is typical or representative but this paper does, in an “existence proof” way, illustrate two points.

The first point is that Burde and Linden (2013) cite the empirical results about the impact of proximity on attendance from just two papers—and those in a footnote. But there are hundreds (if not thousands) of estimates of the empirical connection between school proximity and child school attendance and all recent studies are cognizant of the possible bias from endogenous placement of schools. Filmer (2007), for instance, provides estimates of the impact of proximity on attendance in rural areas for 21 countries, compares the results to 10 other “recent” studies, and uses a number of methods to “sign and bound” the potential bias from endogenous placement. Burde and Linden (2013) was published in the *American Economic Journal: Applied Economics* which shows, again, at least as one case, that it is possible to publish the results of an RCT on a commonly researched question in a top tier economics journal without even any pretense of a literature review that would allow a comparison of their empirical estimates to the body of previous work. Moreover, the paper lacked theory or model. Publishing an empirical paper without literature review or theory/model (either formal or informal) is consistent with an attitude that one should rely only on rigorous evidence, on the premise that (i) the problem of “internal validity” is so severe, its magnitude so unknown, and/or the methods for accommodating this potential bias so unreliable that all previous studies can be ignored and (ii) that new “rigorous evidence” speaks for itself and does not need to encompass the previous knowledge nor do the new estimates need contribute to an ongoing understanding of the phenomena of school attendance which might guide the interpretation and how these particular findings from 13 villages should affect our beliefs about “the evidence.”

This paper illustrates in a more indirect way a second important point about the currently proposed uses of “rigorous” evidence. Suppose one did a “systematic review” of the “rigorous” evidence about school proximity and attendance using the current protocols and filters that accept only papers using rigorous methods that provide unbiased estimates. Burde and Linden (2013) would almost certain be included and Filmer (2007) would likely be excluded. But would this produce better decisions across developing countries? There are three points, which motivate this paper's contributions.

One, by estimating the relationship across 21 countries (and 3 countries having estimates from two different periods), Filmer (2007) provides an average estimated marginal effect of 1 km greater distance on enrollment and can also provide a standard deviation of those estimates. The average is that enrollment is lower by 1.18 percentage points per km of distance. The cross-national standard deviation of those marginal effect estimates is 1.04. This reveals large cross-national heterogeneity in the estimates, with a one standard deviation bound running from 2.2 to 0.16—which implies the estimated impact from changing from the nearest school 10 km away to 0 would be anywhere from 22 percentage points (nearing the 32 percentage points reported for males in Burde and Linden) to only 1.6 percentage points. So, as with the super simple example of assigning pants to men, even the best possible average of “rigorous” estimates would produce large prediction errors as the variation in the true impact is likely large.

Two, Filmer (2007) estimates the marginal effect of proximity to primary and secondary schools for Morocco (1992), the Philippines (1993), and India (with two estimates, for 1992–1993, 1998–1999). He finds in all four cases that estimates of the impact of distance to schools (primary and secondary) is empirically very small (0.1 percentage point gain per km reduction in distance, or less) and (obviously given these tiny magnitudes) statistically insignificant. And yet a review of the “rigorous” evidence would put weight on the evidence from Afghanistan (by being included in the meta-analysis) and zero weight on the actual observational estimates from those countries. But whether this reliance on “rigorous” evidence from other contexts (Afghanistan) improves the prediction of the TE of building schools in Morocco or the

Philippines or India depends on assumptions about the relative magnitudes of the errors from weak external validity (high variability across countries) versus the bias from internal validity. Again, as in the pants example, there is zero *ex ante* or “evidence-based” reason to believe—and, given the huge variance in observational estimates good reasons to not believe—that using an average of rigorous estimates will produce better decisions than relying on biased self-reports. Whether applying the “rigorous evidence” across contexts improves decision-making is entirely an empirically contingent, factual, question.

Three, One reaction to the above is that this is a straw man, and of course, “context” should be taken into account. Perhaps the large impact on enrollment from a school in an Afghan village is consistent with very small effects in Morocco once one takes into account the very different contexts. But “use rigorous evidence” and “take context into account” are mutually incompatible, for two reasons. (1) What the relevant “context” is, that, say, allows findings from Afghanistan to apply to Morocco or the Philippines can only come from a theory. Burde and Linden (2013) illustrate the belief that one does not need theory, just perturbing the world with an experiment and reporting what happens is, in and of itself, “rigorous” evidence. However, without a theory, any claims about the applicability of evidence from one context to another are *ad hoc* and the opposite of “rigorous.” (2) Once one acknowledges that understanding the heterogeneity in estimates requires “context” this implies that correct judgments are likely to emerge from an understanding of the relevant phenomena that encompasses all of the relevant evidence. This implies encompassing the “rigorous” evidence into understanding of the relevant phenomena with many other sources of evidence.

A second example illustrates that the slogan of “rely on the rigorous evidence” leads to reliance on small numbers of studies for generalizations. The recent Global Education Evidence Advisory Panel report (World Bank, 2020) gave (what they regarded as) globally relevant advice about “best buys” in education. The report’s main figure reports estimates of the *average* cost-effectiveness of various (classes of) potential interventions as gains in “learning-adjusted years of schooling” (LAYS) per \$100. The intervention with the highest *average* cost-effectiveness is: “Giving information on education quality, costs, and benefits.” But as the review considers only “rigorous” evidence that “average” impact is based on just *two* studies. Of those two studies one showed very low cost-effectiveness and one showed astronomically high cost-effectiveness, as a relatively large improvement was obtained at very low cost. It would seem the correct conclusion is “anything can happen” (reflecting the variance) not the report’s interpretation: “this is a best buy.” The report also reports average cost-effectiveness estimates for “Teacher accountability and incentive reforms” and “Giving merit scholarships to disadvantaged children and youth,” based on just *three* studies for each. Vivalt (2020) shows there is massive variability in the rigorous TE estimates, both across and within studies. Angrist and Meager (2022) show that the effect size estimates of the impact of a single class of pedagogical intervention, “teaching at the right level,” differ across available rigorous studies by an order of magnitude. The use of the average (central tendency) of a very small number of studies with high variance seems unlikely to produce good predictive outcomes in practice, versus the use of all of the available information.

In previous work, Pritchett and Sandefur (2014, 2015) provide empirical examples in which the root-mean-squared error (RMSE) of predicting TEs was smaller using the context-specific ordinary least squares (OLS) estimates than using the “rely on the rigorous evidence” approach of using average of the rigorous (RCT) estimates. In this paper, I extend this previous work with two additional empirical examples, which provide much larger cross-national samples: 42 countries for estimating the impact of migration on wages and 29 countries for estimating the impact

on measures of learning from private sector schooling. Using cross-national data from these two empirical examples, I illustrate four points (elaborated more fully and technically below).

One, there is no single obvious or plausible interpretation of “rely on the rigorous evidence.” The most common practice in systematic reviews is to focus exclusively on the average of the estimates of impact, but this is a completely arbitrary choice as one could just as easily take “rely on the rigorous evidence” to imply the use of rigorous estimates of bias to adjust observational estimates.

Two, in both empirical examples, predicting the impact in each country with the average of the TE produces *worse* decisions across countries in RMSE than either (i) just using each country's OLS estimate or (ii) adjusting each country's OLS estimate for the average estimated bias.

Three, in both examples, standard economic models predict heterogeneity in the true TEs across countries in ways strongly confirmed by the data, explicitly rejecting the assumption of external validity of TEs, which is needed to make the case for applying evidence from one context to another as “rigorous” evidence.

Four, in both examples, there is also evidence for systematic heterogeneity across contexts in the magnitude of OLS selectivity bias on unobservables (SBU).

The widely practiced approach of doing meta-analysis or systematic reviews that filter out nearly all of the relevant evidence and then acting as the average of the resulting TE (or causal impact) estimates is “the” evidence or “rigorous evidence” for “evidence-based” decisions is both empirically wrong and conceptually “not even wrong” (in the sense of Wolfgang Pauli). This approach relies on conceptually and empirically indefensible assumptions about external validity.

2 | THE HORSE RACE: METHODS

Here is a pretty simple empirical question: “Would using the average of rigorously estimated TEs improve the predictive accuracy of estimating TEs across countries compared to the alternative of just using each country's own OLS estimate?” Pritchett and Sandefur (2015) use the RCT results of estimating the impact of microcredit across six countries from Banerjee et al. (2015) and use the raw data from these same studies to estimate OLS estimates and show the answer is “only sometimes.” For the “reported profits” variable OLS always outperforms the (non-context specific) average RCT predictions. For the “consumption” variable OLS outperforms when the sample of RCTs is small but when all RCTs are used the results are slightly better for RCTs. I feel these results have been underappreciated as this was for a small sample (six countries) for a single intervention (microcredit) and so are ignored or treated as possibly just an anomaly. Replication is difficult because the suitable data for these calculations are scarce as it requires both an OLS estimate and a consistent estimate of the TE across a number of countries.

This paper uses cross-national empirical results that have raw, OLS and a consistent estimate of the (lower bound of the) TE for two phenomena. One, estimates of the wage gains from migration for a specific worker moving from one of 42 countries to the USA, using the ratio of PPP wages. Two, for the learning increment of math from enrollment in private school I have OLS and TE estimates for 29 countries.

2.1 | Selectivity and bias in estimating TEs

In a standard setup (e.g., Altonji et al., 2005) suppose that an outcome Y for individual i in context C depends on whether or not individual i gets “treatment” X (in the examples X is discrete,

either a migrant or nonmigrant or enrolled in private school or not enrolled in private school) and also on other determinants of the outcome, divided into those determinants observed by the econometrician and those unobserved (Equation 1):

$$Y_i^C = \beta^C X_i^C + W_{i,\text{observed}}^C + W_{i,\text{unobserved}}^C. \quad (1)$$

The difficulties of recovering a consistent estimate of the TE in this situation have long been well-known (e.g., the classic treatment in Leamer, 1983 drawing on earlier literature). Just comparing the raw average scores of, say, students in public versus private schools is likely to overstate the TE on learning of private schools as students select into private schools on characteristics of their household that also have a direct causal impact on learning (such as parental education, household wealth/income, socioeconomic status). This selectivity bias can be reduced with estimation methods, say simple multivariate OLS, that include a range of observed characteristics (W) of the student and their HH and hence $\hat{\beta}_{\text{OLS}(W_{\text{observed}})}^C$ is an estimate of the outcome difference for “observationally equivalent on W_{observed} ” individuals.

However, selection into treatment status is plausibly based on characteristics unobserved by the econometrician. This implies that any given $\hat{\beta}_{\text{OLS}(W_{\text{observed}})}^C$ suffers from omitted variables bias to the extent there are $W_{\text{unobserved}}$ which are correlated with selection into treatment. Even conditioning on all observed variables, students with, say, more unobserved grit or ambitious parents, are likely to both have higher measured learning outcomes and to be enrolled in a private school.

One can *define* for any given country and any set of observables the OLS SBU as the gap between the OLS estimate and the true TE (or, by extension, a consistent estimate of SBU is the gap between OLS and a consistent estimate of the TE).

$$\widehat{\text{SBU}}^c \equiv \hat{\beta}_{\text{OLS}(W_{\text{observed}})}^c - \beta_{\text{TE}}^c. \quad (2)$$

Oster (2019) shows that a consistent estimate of the TE of X in Equation (1), $\tilde{\beta}$, can be recovered from observational data and some assumptions via Equation (3).

$$\tilde{\beta} = \hat{\beta} - \delta \left(\hat{\beta} - \hat{\beta} \right) \frac{\bar{R} - \hat{R}}{\hat{R} - \hat{R}}. \quad (3)$$

Oster estimates require two empirical estimates and two assumptions. The two empirical estimates are: (i) the difference in the estimated β with and without W_{observed} , $\hat{\beta} - \hat{\beta}$, which, for a discrete variable X is just the raw difference in averages less the OLS estimate on X , and (ii) the difference in the regression R-squared without and with W_{observed} , $\hat{R} - \bar{R}$, how much higher the OLS R-squared is when the cofounders W_{observed} are included.

In addition to the estimated quantities Equation (3) requires two assumptions: (i) an assumption about a proportionality parameter, δ , between selectivity on the observables and unobservables and (ii) an assumption about the R-squared of Equation (1) with the unobservables included. That is, \bar{R} is the R-squared if both W_{observed} and $W_{\text{unobserved}}$ were included in the regression, and this is usually parameterized as $\bar{R} = \Pi \hat{R}$.

Obviously neither the proportionality parameter δ or Π can be estimated from the data as they depend on “unobserved” variables. Oster (2019) does a review of the literature, comparing

estimates of $\tilde{\beta}(\delta, \Pi)$ to estimates of TEs from other methods, like RCTs. Based on comparisons from the existing literature, she shows the assumptions of $\delta = 1$ and $\Pi = 1.3$ are quite conservative, in that these assumptions would produce TE estimates lower, not higher, than would result from consistent estimation methods. The proportionality assumption of $\delta = 1$ implies that there is as much selectivity into treatment from the unobserved variables as from the unobserved. The assumption that $\Pi = 1.3$ implies the inclusion of $W_{\text{unobserved}}$ would raise the R-squared by 30%. These values have become quite widely adopted.

A key innovation of this paper is to use estimates of TEs from the Oster (2019) method as consistent estimates (of lower bounds) of TEs as this allows, for the first time, the comparisons of large numbers of country estimates. Our two empirical examples, of wage gains from migration to the USA and of learning gains from private school both report Oster estimates with these values. I am going to treat the Oster (2019) estimates with those values as consistent estimates of TEs, which is likely to be conservative in that the “true” TE is larger (in absolute value) as the absence of $W_{\text{unobserved}}$ from the estimation likely creates less SBU than the assumptions of $\delta = 1, \Pi = 1.3$ imply.

2.2 | Horse race for predictive accuracy

If we accept the $\tilde{\beta}_{\text{Oster}}^c(\delta = 1, \Pi = 1.3)$, or $\text{TE}(\text{O})^c$ for short, estimates as consistent estimates of the true TE for each country the RMSE or average absolute deviation (AAD) of prediction errors can be calculated for three “horse race” possibilities.

One, use the average of the $\text{TE}(\text{O})^c$ estimates across all countries as the prediction for each country. This prediction using $\text{SR}(\overline{\text{TE}})$ mimics the “systematic review report of the average of the TEs” approach.

$$\text{RMSE}(\overline{\text{TE}(\text{O})}^c) = \text{sqrt} \left(\frac{\sum_{c=1}^{c=N} (\text{TE}(\text{O})^c - \overline{\text{TE}(\text{O})})^2}{N} \right). \quad (4)$$

Two, naïve OLS predicts each country's TE is its OLS estimate.

$$\text{RMSE}(\beta_{\text{OLS}}^c) = \text{sqrt} \left(\frac{\sum_{c=1}^{c=N} (\text{TE}(\text{O})^c - \beta_{\text{OLS}}^c)^2}{N} \right). \quad (5)$$

Three, the estimate for SBU is, for each country, defined to be equal to the gap between the estimate of the TE, $\text{TE}(\text{O})$, and the OLS estimate, which conditions on observed.

$$\text{RMSE}(\overline{\text{SBU}(\text{O})}^c) = \text{sqrt} \left(\frac{\sum_{c=1}^{c=N} (\text{TE}(\text{O})^c - (\beta_{\text{OLS}}^c - \overline{\text{SBU}(\text{O})}))^2}{N} \right). \quad (6)$$

These formulae for the AAD, which is a more robust measure of predictive accuracy as it does not heavily penalize large prediction errors, are self-explanatory.

3 | FIRST EMPIRICAL EXAMPLE: ESTIMATES OF GAINS FROM LABOR MOBILITY

Suppose you were a developing country government (say, Guatemala) who was initiating a bilateral agreement to increase labor mobility with another country (say, the USA) and you wanted an estimate of the earnings gain to an incremental mover with specific characteristics (e.g., a given level of schooling). You would understand both that observational methods comparing the wages of Guatemalans in the USA to Guatemalans in Guatemala would fail to account for selectivity on unobserved covariates. At the same time, you would understand that rigorous, RCT, estimates of TEs from other pairs of countries might not have external validity for your country. Which would be better, just to rely on estimates using OLS on observational data about Guatemalans or to rely on rigorous evidence from other contexts?

3.1 | Estimates of gains to migrants: Raw, OLS, TE(Oster)

Clemens et al. (2019) use US Census data and labor market surveys from 42 other countries, which jointly allow the comparison of the earnings differences (in PPP dollars, so “real” consumption units) between, say, people born in Guatemala and educated in Guatemala (inferred from the USA census questions about a person’s age at migration) working in the USA (migrants) versus those born and educated in Guatemala and working in Guatemala (nonmigrants). For 42 countries CMP (2019) provide estimates of (i) the raw wage ratio of migrants and nonmigrants, (ii) a standard OLS wage regression in the USA and in the sending country with observables (e.g., age, sex, education, sector, and urban residence to estimate the earnings wage ratio for “observationally equivalent” migrants and nonmigrants (at specific values of the covariates), and (iii) an Oster (2019) lower bound, $TE(O)^c$.

Figure 1 shows four results.

One, the $TE(O)$ estimates of the wage gains are large. Averaged across the 42 countries, a randomly selected low-school worker would make 4.87 times (median 4.1) times higher earnings (in PPP) in the USA than in their home country. This is an average (labor force aged population weighted) wage gain of P\$13,715 (in 2001 dollars). These estimates are consistent with a variety of other approaches to estimating the causal wage gains to a low skill worker moving from poor to rich country (Pritchett & Hani, 2020).

Two, the variation in the $TE(O)$ estimates across countries is substantial: some countries have very high estimates (Egypt at 12.1) while other countries have low estimates (Dominican Republic at 1.9). The 25th percentile (about Uruguay) is 2.6 and the 75th percentile (about Indonesia) is 5.8, more than twice as high. The standard deviation is 3.3.

Three, there typically is quite strong positive migrant selectivity on observables. Positive migrant selectivity on the observed variables leads to OLS estimated wage ratios substantially lower than the raw ratio (5.71 vs. 6.74). This in turn (via Equation 3) produces

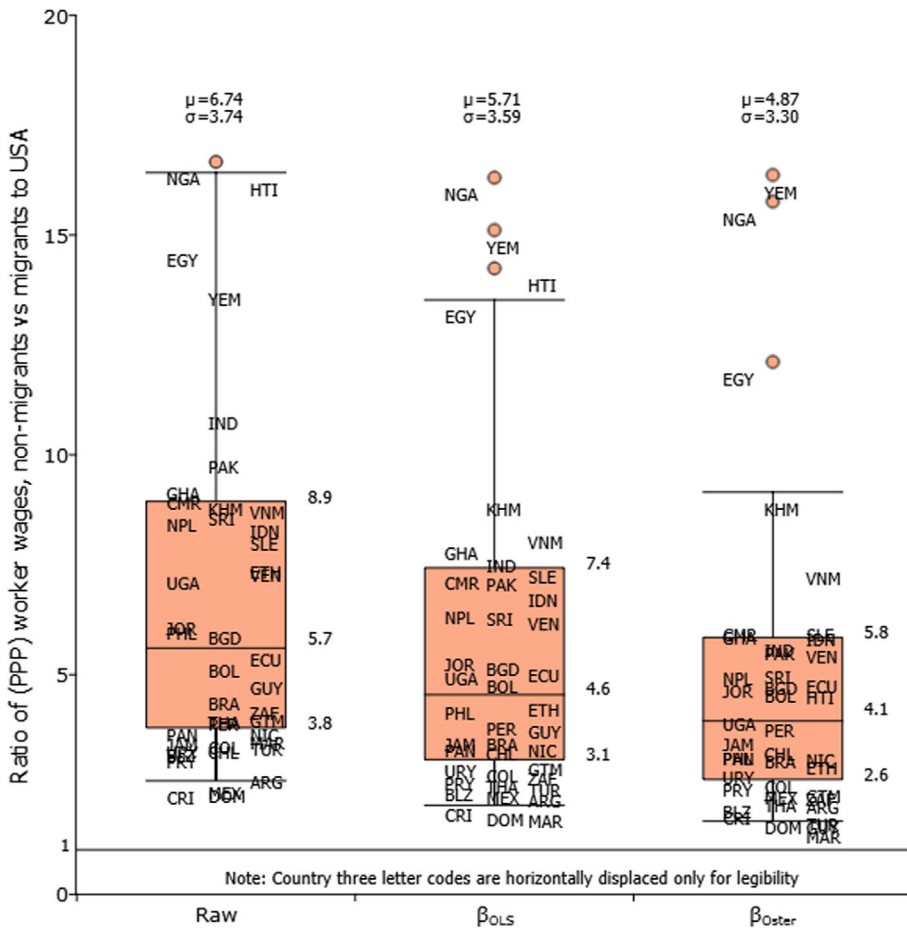


FIGURE 1 Comparison of Raw, OLS, and Oster estimates of ratios of PPP wages for migrants to home country workers from 42 different countries to the USA. Source: Author’s graph based on Clemens et al. (2019) estimates. [Colour figure can be viewed at wileyonlinelibrary.com]

$\tilde{\beta}_{Oster}(1, 1.3)$ estimates lower than OLS, an average of 4.87 versus 5.71. The average OLS SBU of 0.84 ($= 5.71 - 4.87$) is substantially less than the cross-national standard deviation of $TE(O)$ of 3.3.

Fourth, (and this has to be inferred from Figure 1 by comparing the results for specific countries across the box plots) the extent of selectivity on observables varies substantially across countries. Some countries have very small differences between the raw and OLS wage ratio estimates: the difference in Jamaica is 0, in Mexico is 0.09, and Peru is 0.10. In other countries there are very large differences (the difference for India is 3.26, for Ethiopia it is 3.14, and for Morocco is 1.81). Migrant selectivity bias on unobserved productivity in estimating wage gains is often raised in the context of highly skilled professions (e.g., doctors, engineers, academics) and economic “superstars” (e.g., CEOs, entrepreneurs) but the CMP (2019) estimates are for workers with less than high school completed (9–12 years of schooling) for which massive “long-tailed” selectivity wage gains are likely less common.

3.2 | There is no single interpretation of “rely on the rigorous evidence”

There is a well-identified estimate of the wage gains to migrants from a program in New Zealand that allowed Tongan workers to migrate on a temporary basis for agricultural work (McKenzie et al., 2010). Given an oversubscription of visa applicants, recipients were chosen randomly from the applicants, which allow researchers to correct for the potential bias from self-selection of applicants on unobservable characteristics. Their estimated TE was a wage gain ratio of 3.63. This same study also reported an OLS estimate of the migrant/home wage ratio of 4.83 so that the estimated magnitude of the OLS SBU was 1.2 ($= 4.83 - 3.63$).

This study (as pretty much any RCT study could) produced not one, but two, pieces of rigorous evidence: an estimate of the TE of 3.63 (ratio of wage gain) and an estimate of the SBU of 1.2. This implies there are two equally valid interpretations of “rely on the rigorous evidence.” One is to use $SR(\overline{TE})$ the average estimate of the rigorous estimates of TEs. The other possibility is $SR(\overline{SBU})$, use the average estimates of the SBU. The conceptual problem is that, given the large heterogeneity across countries in the OLS evidence about wage gain ratio in Figure 1, these two, equally plausible, approaches to the rigorous evidence: (i) will give will give contradictory advice about how to adjust the OLS evidence to produce an TE estimate and (ii) either interpretation will generate completely implausible implications.

As a simple example, suppose that at first the only evidence we had about wage gains were the OLS regression results for Guatemala (GTM in Figure 1, $\hat{\beta}_{OLS} = 3.2$) and for Bangladesh (BGD in Figure 1, $\hat{\beta}_{OLS} = 5.5$) and these informed our priors for those countries. Then the McKenzie et al. (2010) study was done for Tonga-NZ, which on the assumption this was, at the time, the only “rigorous” study by the filter of a systematic review, implies $SR(\overline{TE}) = 3.63$ and $SR(\overline{SBU}) = 1.2$. What would “rely on the rigorous evidence” mean?

The $SR(\overline{TE})$ interpretation would imply that we should revise our estimate of the wage gains in Guatemala *upward* from 3.2 to 3.63. The $SR(\overline{SBU})$ interpretation would imply we should revise our estimate of the wage gains in Guatemala *downward* from 3.2 to 2.0. Moreover, the $SR(\overline{TE})$ interpretation suggests we should revise our estimate of the wages gains in Bangladesh *downward* from 5.5 to 3.7. This *necessarily* means the $SR(\overline{TE})$ interpretation of “rely on the rigorous evidence” requires us to believe, given the identity in Equation (2), and that the SBU for Tonga-NZ is 1.2, the SBU for Guatemala is *negative* 0.5 (implying migrants are negatively selected and hence that OLS is too low relative to the true TE) and yet also believe that the SBU for Bangladesh is 1.8.

Things do not get any better from this simple example if there are multiple rigorous studies, as the two potential interpretations of “rely on the rigorous evidence” still contradict each other and either interpretation has obviously counterfactual implications. Suppose we treat the $TE(O)^c$ as the rigorous evidence a systematic review would be based on. If we assume TEs have external validity then all countries should believe their country’s wage gain ratio should be $SR(\overline{TE}) = 4.87$ (Equation 7a). However, if we assume estimates of the OLS SBU have external validity then each country should believe that the wage gain ratio for their country should be the OLS estimate less the average estimated SBU (Equation 7b).

$$\beta_{c,USA}^{True} = SR(\overline{TE}) = 4.87, \quad (7a)$$

$$\beta_{c,USA}^{\text{True}} = \hat{\beta}_{c,USA}^{\text{OLS}(W_{\text{obs}})} - \text{SR}(\overline{\text{SBU}}), \text{ where } \text{SR}(\overline{\text{SBU}}) = \text{mean}\left(\beta_{c,USA}^{\text{OLS}} - \tilde{\beta}_{c,USA}^{\text{Oster}}\right) = 0.84. \quad (7b)$$

The identity in Equation (2) linking OLS, SBU, and the true TE implies Equation (7c)

$$\beta_{c,USA}^{\text{OLS}} \equiv \beta_{c,USA}^{\text{True}} + \text{SBU}_{c,USA}. \quad (7c)$$

The OLS estimate for each country c , $\hat{\beta}_{c,USA}^{\text{OLS}(W_{\text{obs}})}$ is just an empirical fact (the determinate outcome of applying a given statistical procedure to a given dataset) and cannot be freely chosen.

Table 1 illustrates five serious logical and empirical problems with assuming external validity.

First, the gaps between the $\text{TE}(O)^c$ estimates and $\text{SR}(\overline{\text{TE}}) = 4.87$ are large and arbitrary and there is no theoretical or empirical justification for believing there is external validity and the country-specific $\text{TE}(O)^c$ estimates are just mistaken.

Second, assuming external validity of TE and adopting that each country's estimate should $\text{SR}(\overline{\text{TE}})$ implies the cross-national standard deviation of the “true” TE is zero (Column V). But the cross-national standard deviation of the $\text{TE}(O)^c$ estimates is 3.30. Hence assuming external validity, $\text{SR}(\overline{\text{TE}})$, implies strongly counterfactual beliefs about the cross-national variation in the true TEs. The same problem arises with assuming external validity about the bias, taking $\text{SR}(\overline{\text{SBU}})$ as the rigorous estimate of the SBU for all countries, as the estimated standard deviation of SBU is 1.5.

Third, Column VII shows the estimate of the OLS SBU for each country is implied by the identity in Equation (2) and $\text{SR}(\overline{\text{TE}})$. Given India's OLS estimate of 7.86 the implied SBU is 2.99 (7.86–4.87) which implies that Indian (low schooled) migrants to the USA are strongly positively selected on unobservables, more strongly than India's Oster estimated SBU estimate of 1.93 (Column IV). Conversely, the OLS estimate for the Dominican Republic is 2.08, which produces an OLS SBU estimate implied by $\text{SR}(\overline{\text{TE}})$ of -2.79 ($= 2.08 - 4.87$) which implies (low skill) workers from the Dominican Republic are massively *negatively* selected on unobservables—even though the Oster estimate for the Dominican Republic suggests migrants are modestly *positively* selected on unobservables, at 0.18 ($= 2.08 - 1.90$).

Moreover, the combination of assuming external validity of TEs, which implies zero variation in the true TE across countries and the actual variation of the OLS estimates (which are an empirical fact) implies that all of the variation in OLS versus $\text{SR}(\overline{\text{TE}})$ must be due to variation in the country SBU, which implies a variation in the SBU of 3.59, much higher than its estimated value of 1.5.

Fourth, the estimates of the OLS SBU implied by the assumption of external validity of $\text{SR}(\overline{\text{TE}})$ bear no relationship to the actual country-specific empirical estimates of the SBU, $\hat{\beta}_{c,USA}^{\text{OLS}(W_{\text{obs}})} - \tilde{\beta}_{c,USA}^{\text{Oster}}$, or common sense, or the existing literature. Excluding one outlier country, Haiti, the correlation between the Oster SBU estimates in Column IV and the $\text{SR}(\overline{\text{TE}})$ estimates in Column VII is modestly *negative*, at -0.11 . One implication of Column VII estimates of the OLS SBU is that 22 of 42 countries have *negative* selectivity on unobservables in spite of *positive* selection on observables, which would be extremely odd. Moreover, empirical estimates of migrant selectivity across a large number of countries suggest positive selectivity on both observables and unobservables (Clemens & Mendola, 2020).

TABLE 1 Thought experiment comparing $SR(\overline{TE})$ and $SR(\overline{SBU})$ with actual country-specific OLS and TE estimates.

Country Column number	Actual country estimates										
	$SR(\overline{TE})$					$SR(\overline{SBU})$					
	Raw OLS	Oster OLS	SBU	Average of Oster estimates	Gap with $SR(\overline{TE})$	Implied OLS SBU	Estimate	Gap with Oster estimate	Gap $SR(\overline{SBU})$ and $SR(\overline{TE})$		
I	II	III	IV (= II - II)	V	VI (= V - III)	VII (= II - V)	VIII (= II - 0.84)	IX (= VIII - III)	X (= 0.84 - IV)		
Highest 10 wage ratio countries by Oster estimate											
Yemen	13.92	15.11	16.37	-1.25	4.87	11.50	10.24	14.28	2.09	2.09	9.41
Nigeria	16.67	16.31	15.76	0.54	4.87	10.89	11.44	15.47	0.29	0.29	10.60
Egypt	14.82	13.53	12.12	1.41	4.87	7.25	8.66	12.69	-0.57	-0.57	7.82
Cambodia	9.13	9.14	9.15	-0.01	4.87	4.28	4.27	8.30	0.85	0.85	3.43
Vietnam	9.06	8.40	7.55	0.84	4.87	2.68	3.52	7.56	0.00	0.00	2.69
Cameroon	9.27	7.48	6.29	1.19	4.87	1.42	2.61	6.64	-0.35	-0.35	1.77
Sierra Leone	8.35	7.61	6.27	1.34	4.87	1.40	2.74	6.77	-0.50	-0.50	1.90
Ghana	9.51	8.16	6.23	1.93	4.87	1.36	3.29	7.32	-1.09	-1.09	2.45
Indonesia	8.64	7.07	6.19	0.88	4.87	1.32	2.20	6.23	-0.04	-0.04	1.36
India	11.12	7.86	5.93	1.93	4.87	1.06	2.99	7.02	-1.09	-1.09	2.15
Average	6.74	5.71	4.87	0.84	4.87	0.00	0.84	4.87	0.00	0.00	0.00
Std. Dev.	3.74	3.59	3.30	1.50	4.87	3.30	3.59	3.59	1.50	1.50	3.59
Smallest 10 wage ratio countries by Oster estimate											
Mexico	2.68	2.59	2.56	0.03	4.87	-2.31	-2.28	1.75	0.81	0.81	-3.12
South Africa	4.49	2.99	2.52	0.46	4.87	-2.35	-1.89	2.15	0.38	0.38	-2.72
Thailand	4.30	2.83	2.40	0.43	4.87	-2.47	-2.04	1.99	0.41	0.41	-2.88
Argentina	2.93	2.49	2.36	0.12	4.87	-2.51	-2.38	1.65	0.72	0.72	-3.22
Belize	3.52	2.63	2.25	0.39	4.87	-2.62	-2.24	1.80	0.45	0.45	-3.07
Costa Rica	2.58	2.19	2.10	0.10	4.87	-2.77	-2.68	1.36	0.74	0.74	-3.51

TABLE 1 (Continued)

Country Column number	Actual country estimates					SR(TE)	SR(SBU)			
	Raw OLS		Oster estimated OLS SBU		Average of Oster estimates	Gap with SR(TE)	Implied OLS SBU	Estimate	Gap with Oster estimate	Gap SR(SBU) and SR(TE)
	I	II	III	IV (= II - II)	V	VI (= V - III)	VII (= II - V)	VIII (= II - 0.84)	IX (= VIII - III)	X (= 0.84 - IV)
Turkey	3.68	2.74	1.95	0.79	4.87	-2.92	-2.14	1.90	0.05	-2.97
Guyana	5.08	4.07	1.90	2.17	4.87	-2.97	-0.80	3.23	-1.33	-1.64
Dom. Rep.	2.62	2.08	1.90	0.19	4.87	-2.97	-2.79	1.25	0.65	-3.62
Morocco	3.84	2.03	1.67	0.36	4.87	-3.21	-2.84	1.19	0.48	-3.68
Root-mean-square error						3.26				
Average absolute deviation						2.21				

Source: Author's calculations with estimates from CMP (2019, table 2).

Fifth, the convention that “rely on the rigorous evidence” implies some assumptions about external validity (as otherwise it is obvious an RCT evidence from one country is not rigorous evidence at all for any other country), but there is no logical or theoretical reason to believe that TEs have external validity versus that selection bias has external validity—and both cannot have external validity. For instance, the OLS estimate for Thailand is 2.83 (Column II). With $SR(\overline{TE})$ (external validity of TEs) the estimate would be 4.87 (Column V), much higher, with $SR(\overline{SBU})$ the estimate would be 1.99 (Column VIII), which is much lower. Does “rely on the rigorous evidence” mean $SR(\overline{TE})$ or $SR(\overline{SBU})$? It cannot mean both and there is no rational reason to prefer one over the other.

These five problems are quite general, in three senses.

One, one could use any single value of TE or SBU derived in whatever way from any set of rigorous estimates, that is, the “systematic review” could be filtered in any way, and still have exactly the same five issues. So many methodological issues, like the use of meta-analysis approaches to give weights in forming the “optimal” average across estimates, are irrelevant to addressing these concerns.

Two, one could modify Equations (7a)–(7d) so the country-specific prediction was a weighted average of the OLS and $SR(\overline{TE})$ with any α on $SR(\overline{TE})$ (7a with $\alpha = 1$ is a special case). One still has all the same five problems, just moderated somewhat (Pritchett & Sandefur, 2014)—and one loses the rhetorical appeal to “rigorous” as neither the OLS estimates nor the weight α can be considered “rigorous.”

$$\beta_{c,USA}^{True} = (1 - \alpha) \times \beta_{c,USA}^{OLS(W)} + \alpha \times \beta_{c,USA}^{SR(\overline{TE})}. \quad (7d)$$

Three, one could combine Equations (7a) and (7b) so that the estimate of the “true” TE in country c was a weighted combination of the average TE and the OLS adjusted for the average selection bias, Equation (7e). This however implies that “rely on the rigorous evidence” can mean pretty much anything, depending on the choice of the weight θ . For instance, the TE for Belize, estimated ratio of wages in the USA to wages in Belize for an equal productivity mover, could be anywhere from 4.87 ($\theta = 1$, Column V) to 1.80 ($\theta = 0$, Column VIII), including producing as the “rigorous” estimate the OLS estimate of 2.63 using $\theta = 0.27$ or the Oster estimate of 2.25 with $\theta = 0.15$, or with the arbitrary but focal point of equal weights, $\theta = 0.5$, one could estimate the true gain as 3.33.

$$\beta_{c,USA}^{True} = (1 - \theta) \times \left(\hat{\beta}_{c,USA}^{OLS(W_{obs})} - SR(\overline{SBU}) \right) + \theta \times SR(\overline{TE}). \quad (7e)$$

The slogan “rely on the rigorous evidence as summarized by systematic reviews” is vacuous, but any specific interpretation of that faces enormous challenges in even achieving logical coherence and even minimal empirical plausibility. A standard approach is for systematic reviews to completely ignore the OLS evidence, so that $\alpha = 1$, and completely ignore the estimates of the SBU, so that $\theta = 1$, and hence the slogan is the special case, $RORE(1,1)$ which implies the best prediction for each country is $SR(\overline{TE})$. This, however, implies the assertion that TEs have external validity but estimates of SBU have no external validity, and worse, the estimates of SBU implied by $SR(\overline{TE})$ have to take on country by country values that are an arbitrary set of measure zero.

The existing systematic reviews never acknowledge these conceptual challenges and avoid them by “feigned ignorance” (Pritchett, 2020), which just (i) pretends the OLS estimates do not

exist and (ii) ignores that studies which can produce a rigorous estimate of TE also (can) produce OLS and hence via an identity, produce rigorous estimates of the SBU of OLS.

3.3 | $SR(\overline{TE})$ produces worse cross-national predictions than OLS or RORE(bias)

$SR(\overline{TE})$ produces worse predictions of the “true” wage gains. Figure 2 (and the bottom two rows of Table 1) show the results of the horse race, where the RMSE and AAD of $SR(\overline{TE})$ are normed to 1. OLS country by country produces an RMSE about half that of $SR(\overline{TE})$, and performs even worse for AAD. $SR(\overline{SBU})$ outperforms either $SR(\overline{TE})$ (by a wide margin) or OLS (by a modest margin).

The intuition of this result is clear. In this data standard deviation of the $TE(O)^c$ of 3.3 (Table 1) is much larger than the typical OLS SBU of 0.84. Hence, ignoring the cross-national variance in the true TE in order to use an average of the “rigorous” estimates leads to worse predictions on average.

3.4 | $SR(\overline{TE})$ is not based on any plausible theory

A prediction of the wage gain from actions that allow a worker i to move from country h to country d should be based on our best available model predicting the wages of the mover in those two places. Implicitly, assumptions of external validity of estimated policy-relevant quantities and the proposed method of $SR(\overline{TE})$ assume that most (or all) of the variation in the biased observational estimates of the TE are due to flawed methods and not due to true cross-contextual variation in the TE. Therefore if substantial variation in the $TE(O)^c$ estimates are associated with *any* model of cross-national wage differences for workers with equal intrinsic (personal) productivity this assumption is false and hence $SR(\overline{TE})$ is scientifically dubious.

In Solow–Swan models there are cross-national differences in A or TFP and these differences imply different marginal products of factors, capital, or human capital. If factors are paid their marginal product then wages in countries h and d for a worker with the same human capital will be higher where A is higher. Figure 3 shows the scatter plot of the 29 countries that have both a Penn World Table 10.0 estimate of TFP relative to the USA at current PPPs and also a CMP (2019) estimate of wage differences. The association between country-level TFP relative to the USA and the estimated $TE(O)^c$ wage ratios is strong and negative (there are two large, Nigeria [NGA] and Egypt [EGY], and the regression includes a dummy for each of those countries). The regression is strongly consistent with the idea that equal intrinsic productivity workers gain more by moving to the USA from moving from countries with lower TFP relative to the USA.

This is not to say this simply cross-national association based on a simple model of aggregate output is the “best” model of gains from migration, the point that even this “quick and dirty” economics suggests that “external validity” cannot be assumed as it is not based on *any* economics at all and contradicts even simple, but empirically validated, models.

3.5 | $SR(\overline{SBU})$ is also not based on any valid theory

The primary reason why it is so difficult to recover reliable causal estimates about economic phenomena from observational data is that the data reflect the results of agents making purposive

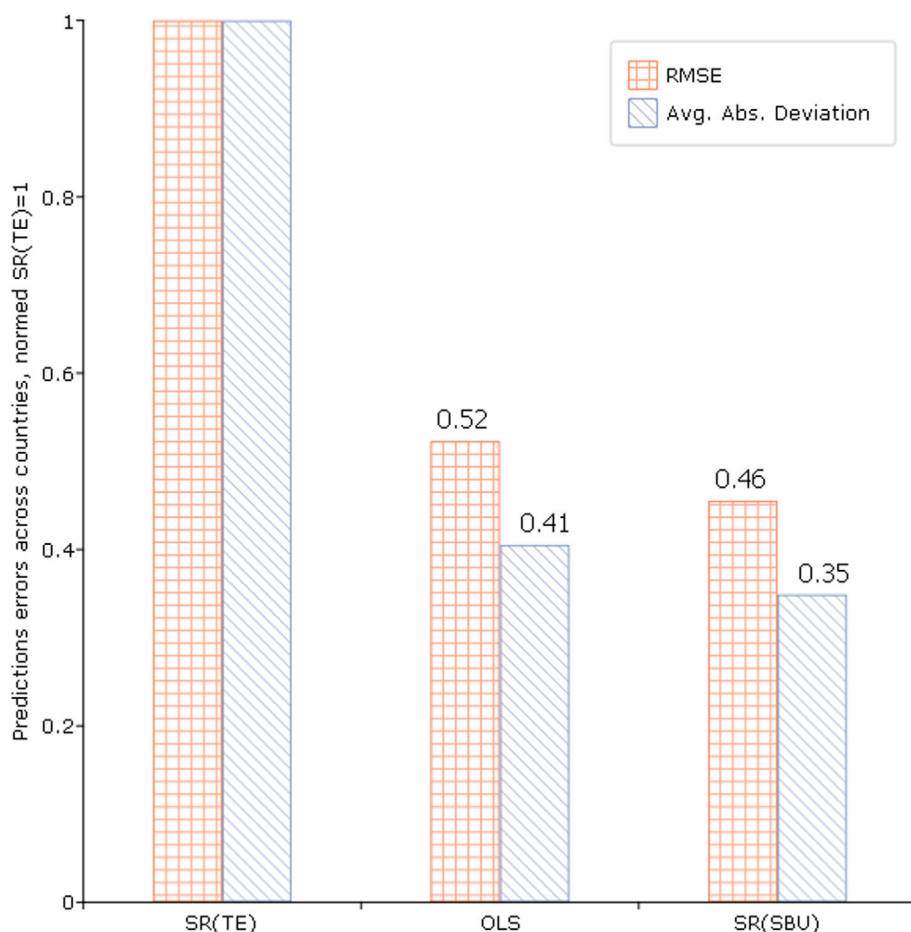


FIGURE 2 Prediction errors from $SR(\overline{TE})$, using the average estimated treatment effect from a systematic review, the standard interpretation of “rely on the rigorous evidence,” are much worse than OLS or $SR(\overline{SBU})$. *Source:* Author’s calculations with data in Figure 1 and Table 1. [Colour figure can be viewed at wileyonlinelibrary.com]

decisions. But this implies that the *magnitude* of selectivity bias depends on the underlying economics of the choices agents make, subject to the constraints they face. This then implies that the selectivity bias, on both observables and unobservables, may vary from context to context. Assuming external validity for $SR(\overline{SBU})$ faces the same problems as assuming external validity for TEs, the assumption is not based on any (much less the best available) understanding/model/theory.

For instance, it is plausible that the higher the fixed costs of a given move the larger the selectivity bias, as only those who anticipate larger gains are willing to make the move. Figure 4 shows a simple scatter plot between an estimate of selectivity of migrants on both observables and unobservables (the gap between the country raw wage ratio and $TE(O)^c$) and the distance from the country to the USA (Mayer & Zignago, 2011). As can be seen there are some massive outliers (Haiti, Yemen, Cambodia [KHM]) but if one allows for dummy variables there is a strong positive association between distance and estimates of selectivity bias consistent with a simple economic model. Again, the point is not that Figure 4 illustrates a complete and correct model of cross-national differences in the selectivity bias of wage gains for (low schooling level)

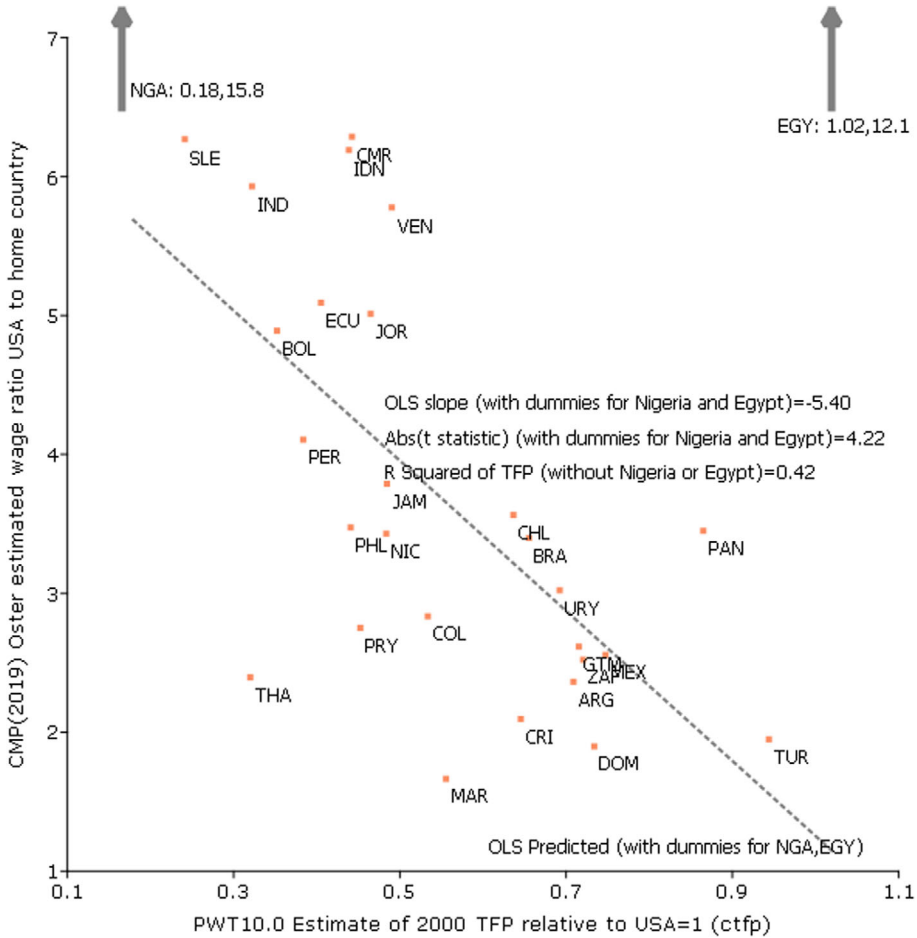


FIGURE 3 Wage gains of movers to the USA and TFP relative to the USA are strongly (negatively) correlated. *Source:* Author’s calculations with CMP (2019) estimates. OLS regression includes dummy variables for the extreme observations for Nigeria (NGA), Egypt (EGY) (whose data are indicated in the graph). [Colour figure can be viewed at wileyonlinelibrary.com]

migrants, but rather only that assuming “external validity” of bias estimates has no justification as the best available understanding of the migrant selection process. “External validity” of selectivity estimates would imply that the variance in selectivity across countries should not be predictable on the basis of *any* economic model.

4 | SECOND EMPIRICAL EXAMPLE: PRIVATE SCHOOL LEARNING PREMIUM

The sections above are a complete paper and lay out the insuperable conceptual and empirical objections to adopting the implicit assumptions embedded in the apparently straightforward and seemingly good advice to “rely on the rigorous evidence.” This section adds to that by showing that all of the conceptual and empirical points made above using the cross-national

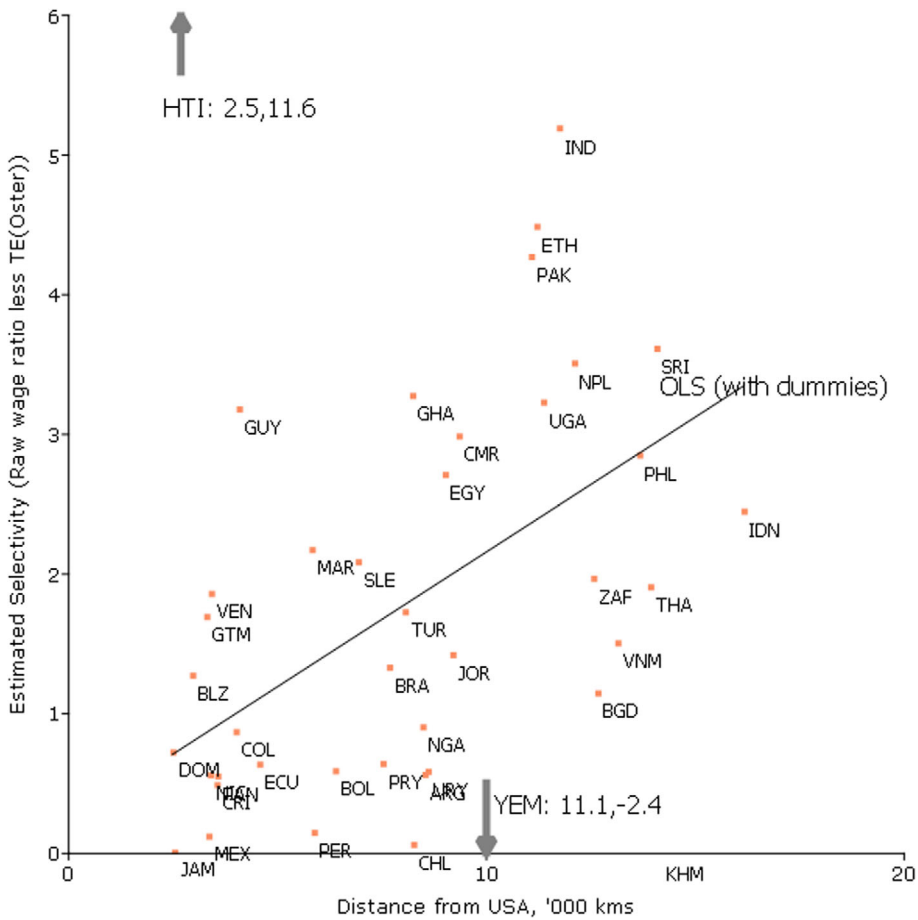


FIGURE 4 Estimated association of impact of selectivity in migration to the USA on estimated wage gains and distance to the USA. *Source:* Author's calculations from CMP (2019) estimates and data on distance from CPEII. OLS regression includes a binary variable for Haiti (HTI), Yemen (YEM), and Cambodia (KHM). [Colour figure can be viewed at wileyonlinelibrary.com]

data on wage ratios hold for a completely different phenomenon, the private sector learning premia. This is important in reassuring the reader the wage ratio example was not a “special case” that is uniquely unfavorable for the case for “rely on the rigorous evidence” but that the problems raised by large variability across countries in the true TE are generic. Given the space constraints of this paper, the hopeful explication of the points above, and the fuller discussion in the working paper version (Pritchett, 2021), this section will be telegraphic.

A recent study by Patel and Sandefur (2020) uses a “Rosetta Stone” approach of giving students in a single setting an assessment with items from different assessments to create comparable estimates of mathematics capabilities for large samples of individual students in 29 developing countries. They estimated: (i) the “raw” private sector learning premium (PSLP), (ii) an OLS estimate of the PSLP controlling for a set of student and household covariates, and (iii) TE(O)^c estimates of the PSLP, using $\delta = 1$ and $\Pi = 1.3$.

Figure 5 shows the same four key points about the distribution of the empirical estimates of the PSLP as shown for the wage ratios. One, on average the TE(O) is substantial, with an

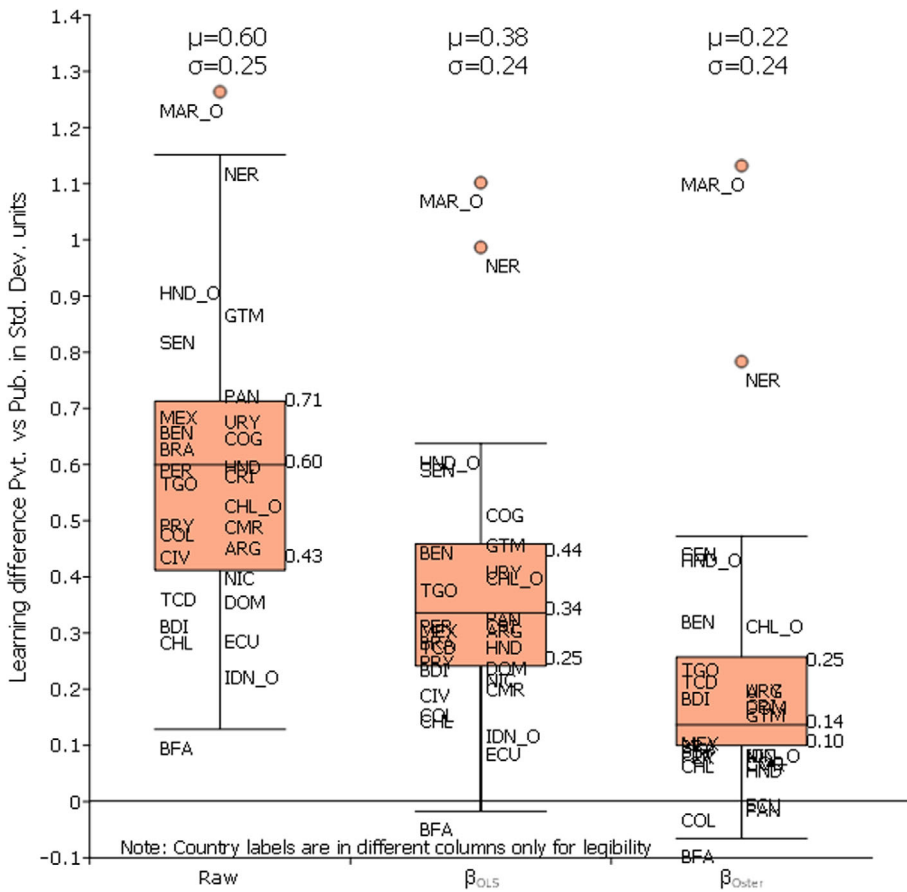


FIGURE 5 Box plots of Raw, OLS, and TE(O) estimates of the private sector learning premium (PSLP). Source: Patel and Sandefur (2020). The notation ABC_O in the country labels (e.g., CHL_O) for Chile indicates estimates of “original” data, not based on the “Rosetta Stone” estimates. [Colour figure can be viewed at wileyonlinelibrary.com]

average of 0.22 standard deviations (a median of 0.14, as the estimates are skewed). Two, the heterogeneity of the TE(O) is substantial, the 25th–75th spread is 0.15 and the standard deviation is 0.24. Three, there is quite a strong selectivity on observables as the average OLS is 0.38 versus the average raw PSLP is 0.60. Four, (and again this has to be inferred from comparison across the box plots) there are substantial differences in the degree of selectivity on observables and unobservables, from quite small for Morocco (MAR_O) to very large for Mexico (MEX), which falls from a raw of 0.72 to a TE(O) of only 0.14.

Table 2—which are the same calculations as in Table 1—illustrates with PSLP data that the assumption of external validity of the TE estimates and the identity linking OLS, TE, and bias necessarily leads to unreasonable implications. For instance, in Ecuador (ECU) the estimates imply strong positive selectivity bias on observables, the gap between the raw and the OLS is 0.20 (0.32 – 0.12), but the $SR(\overline{TE})$ estimate of 0.22 implies selection into private school on unobservables must be *negative*, as 0.22 is *higher* than Ecuador's OLS estimate of the PSLP of 0.12. The assumption of $SR(\overline{TE})$ implies there are five countries with *positive* selectivity on observables but *negative* selectivity on unobservables. Conversely, in

TABLE 2 Estimates of the private sector learning premium, $SR(\overline{TE})$, and $SR(\overline{SBU})$.

Country Column number	Patel and Sandefur (2020) estimates						SR(\overline{TE})			SR(\overline{SBU})		
	Raw OLS		Implied OLS SBU		Average TE(O)		Gap with TE(O) ^c		Implied OLS SBU		TE estimate	
	I	II	III	IV	V	VI (= V - III)	VII (= II - V)	VIII (= II - V)	IX (VIII - III)	X (= IV - IV)(avg)		
Burkina Faso	0.13	-0.02	-0.07	0.05	0.22	0.29	-0.24	-0.17	-0.10	0.10		
Columbia	0.51	0.19	0.00	0.19	0.22	0.23	-0.04	0.04	0.04	-0.04		
Panama	0.75	0.36	0.02	0.34	0.22	0.21	0.13	0.21	0.19	-0.19		
Ecuador	0.32	0.12	0.02	0.09	0.22	0.20	-0.11	-0.04	-0.06	0.06		
Honduras	0.63	0.31	0.09	0.22	0.22	0.14	0.08	0.16	0.07	-0.07		
Chile	0.31	0.18	0.10	0.08	0.22	0.13	-0.05	0.03	-0.07	0.07		
Cameroon	0.52	0.23	0.10	0.14	0.22	0.13	0.01	0.08	-0.02	0.02		
Rep. of Congo	0.68	0.54	0.10	0.44	0.22	0.12	0.32	0.39	0.29	-0.29		
Paraguay	0.52	0.28	0.11	0.17	0.22	0.11	0.05	0.13	0.02	-0.02		
Peru	0.62	0.34	0.11	0.23	0.22	0.11	0.12	0.19	0.08	-0.08		
Nicaragua	0.43	0.25	0.11	0.14	0.22	0.11	0.02	0.10	-0.02	0.02		
Indonesia	0.25	0.15	0.12	0.03	0.22	0.11	-0.08	0.00	-0.12	0.12		
Côte d'Ivoire	0.47	0.22	0.12	0.11	0.22	0.11	0.00	0.07	-0.05	0.05		
Brazil	0.66	0.32	0.13	0.18	0.22	0.09	0.09	0.16	0.03	-0.03		
Mexico	0.72	0.34	0.14	0.20	0.22	0.09	0.11	0.19	0.05	-0.05		
Guatemala	0.90	0.49	0.19	0.30	0.22	0.04	0.27	0.34	0.15	-0.15		
Dom. Rep.	0.39	0.27	0.20	0.07	0.22	0.02	0.05	0.12	-0.08	0.08		
Costa Rica	0.61	0.35	0.20	0.14	0.22	0.02	0.12	0.19	-0.01	0.01		
Burundi	0.34	0.27	0.22	0.05	0.22	0.01	0.04	0.11	-0.10	0.10		
Uruguay	0.71	0.44	0.23	0.21	0.22	0.00	0.22	0.29	0.06	-0.06		

TABLE 2 (Continued)

Country Column number	Patel and Sandefur (2020) estimates									
	SR(TE)					SR(SBU)				
	Raw OLS	TE(O)	Implied OLS SBU	Average TE(O)	V	Gap with TE(O) ^c VI (= V – III)	Implied OLS SBU	TE estimate VIII (= II – V)	Gap with TE(O) ^c IX (VIII – III)	SBU gap X (= IV – IV(avg))
Argentina	0.48	0.34	0.23	0.11	0.22	0.00	0.11	0.18	-0.04	0.04
Chad	0.39	0.31	0.25	0.06	0.22	-0.02	0.08	0.15	-0.09	0.09
Togo	0.60	0.41	0.27	0.14	0.22	-0.04	0.18	0.26	-0.01	0.01
Chile_O	0.56	0.43	0.34	0.09	0.22	-0.12	0.20	0.28	-0.07	0.07
Benin	0.69	0.48	0.35	0.13	0.22	-0.13	0.25	0.33	-0.03	0.03
Honduras_O	0.94	0.64	0.46	0.17	0.22	-0.24	0.41	0.49	0.02	-0.02
Senegal	0.85	0.62	0.47	0.15	0.22	-0.25	0.40	0.47	0.00	0.00
Niger	1.15	0.99	0.78	0.20	0.22	-0.56	0.76	0.83	0.05	-0.05
Morocco_O	1.26	1.10	1.13	-0.03	0.22	-0.91	0.88	0.95	-0.18	0.18
Mean	0.60	0.38	0.22	0.15	0.22	0.00	0.15	0.22	0.00	0.00
Median	0.60	0.34	0.14	0.14	0.22	0.09	0.11	0.18	-0.01	0.01
Root-mean-squared error				0.18		0.24			0.10	
Average absolute deviation				0.15		0.16			0.07	

Source: Author's calculations with estimates from Patel and Sandefur (2020).

Niger (NER) the selection of observables is relatively strong and the $TE(O)$ is 0.20 units lower than the OLS. But the $SR(\overline{TE})$ of 0.22 implies the OLS SBU in Niger was not 0.20, but three times larger, 0.76. Again, these dubious empirical implications are *necessarily* implied by the $SR(\overline{TE})$ assumption of the external validity of TE estimates.

Again the SBU estimates implied by assuming the external validity of TEs must fall onto exactly (the arbitrary set of measure zero) results in Column VII. As argued above, there is no reason to prefer $SR(\overline{TE})$ over $SR(\overline{SBU})$ as it is a priori at least as plausible there is cross-country external validity in the estimates of the OLS SBU and hence that the country-specific estimates of the TE should be the result of adjusting the country-specific OLS estimates for the average cross-national estimated bias (as in Column VIII).

Table 2 shows the RMSE and AAD from using either $SR(\overline{TE})$, country-specific OLS, or $SR(\overline{SBU})$, assuming $TE(O)^c$ estimate is the “true” TE for each country. As with wage ratios, the RMSE error for $SR(\overline{TE})$ is more than twice as large as that for $SR(\overline{SBU})$ (0.24 vs. 0.10) and larger than for OLS (0.24 vs. 0.18). In this case, there is a modest caveat as the AAD is only slightly lower for OLS than $SR(\overline{TE})$ and the RMSE without Morocco is slightly larger for OLS than for $SR(\overline{TE})$.

The PSLP results reinforce the intuition that $SR(\overline{TE})$ will produce worse prediction errors than OLS when the cross-national variation in the true TE is large relative to the selectivity bias (average OLS SBU). With wage ratios the variation in the true TE was large and OLS SBU modest, whereas with the PSLP these are of roughly similar magnitude and hence $SR(\overline{TE})$ and OLS RMSE prediction errors (excluding Morocco) are quite similar.

A model that the PSLP is constant across all countries is easily rejected. In the Patel and Sandefur (2020) data there is a strong negative, nonlinear association (R^2 of .305) between the $TE(O)^c$ PSLP estimates and the math assessment results in the public sector. One need not have a complete and fully articulated model of the cross-national variation in the PSLP to think it is plausible that some governments are reasonably effective and can produce learning outcomes near the efficiency frontier and hence in those cases the PSLP will be low. But when governments are not effective (and economists have no validated positive model suggesting all governments are equally effective) they may be very bad at producing learning. This low government efficacy creates the possibility of the private sector outperforming the public sector by a wide margin and the PSLP is high.

Similarly, there is no plausible case that there is external validity in the estimates of selectivity bias. Patel and Sandefur (2020) show (fig. 17 in their paper) that the measured total selectivity bias—the gap between the raw PSLP and the $TE(O)$ estimate—is associated with the country's income inequality. Countries with larger inequality (e.g., Guatemala and Honduras) tend to have a higher selectivity bias than do lower income inequality countries (e.g., Indonesia or Morocco). Assuming equal selectivity bias across countries is not consistent with the data.

5 | CONCLUSION

This paper, together with Pritchett and Sandefur (2015), makes the horse race score in predicting the actual country-specific TEs 3–0, with the most naïve use of OLS winning three and “rigorous evidence” never winning. Across three very different subjects—microcredit, wage gains from migration, private school learning premium—using naïve OLS from the specific country gives better RMSE than “rely on the rigorous evidence” interpreted as $SR(\overline{TE})$ —or, in the notation of Equations (7d) and (7e) $RORE(1,1)$. And worse, as emphasized in Pritchett and

Sandefur (2014), the standard approach to “systematic reviews” which focuses on summarizing TEs on the premise these summaries have evidentiary value relies on assumptions about external validity that are completely indefensible: indefensible conceptually—there is no reason why TE estimates have external validity and estimates of bias do not; indefensible theoretically—economic models predict heterogeneity of TEs and empirical associations of TEs with theoretical predictions are inconsistent with external validity; and indefensible empirically—assuming external validity of TE implies impossible beliefs about the cross-national variance in observational estimates.

To conclude on a more positive note, the alternative to the naïve slogans about an “evidence-based” approach to decision-making is “understanding-based” decision-making. This paper had its origins as an encomium to Edward Leamer, who, for me, emphasized that the goal of empirical economics was a correct understanding rather than methodological purity. A correct understanding of the relevant phenomena needed to be capable of *encompassing* all of the available evidence into an overall theory or model or narrative, which in turn means that our understanding needs to be dynamic, as even past reliable empirical associations can break down in new circumstances (Leamer, 2010). In a social science like economics, this is the sense understanding of the German word *verstehen* (a concept whose consequences for the method were elaborated (for me) by Gadamer 1975), an interpretive understanding, while, in a social science like economics, appreciating that this interpretive understanding needs to encompass and embed empirical findings. And the application of a correct understanding to concrete decisions needs to reflect something like the Greek word *phronesis*, or practical wisdom. And the effective implementation of policies, programs, and projects that actually leads to improvements often relies on knowledge as both *metis* and *techné*, as articulated by Scott (1998).

Many RCT papers are important. But that is not because their specific numerical findings have immediate applicability to “policy” or, much less, have external validity. Rather RCTs provide well-documented and empirical anecdotes about directed perturbations to the existing equilibria of complex systems. Hence, they force us to adapt our interpretive understanding and expand our collective practical wisdom and hence, at times, point to potential new pathways of betterment. I do not “reject” RCTs as a research method: there are dozens of RCT studies that have changed my understanding and which I routinely use in my writing and teaching (e.g., Andrabi et al. 2017; Andrabi et al., 2020; Banerjee et al., 2008; Bertrand et al., 2007; Dhaliwal & Hanna, 2017; Glewwe et al., 2009; Kerwin & Thornton, 2021; Muralidharan & Singh, 2020; Muralidharan & Sundararaman, 2015; Olken, 2007). However, none of these excellent papers provide evidence that is “rigorous” about anything other than exactly what is reported.

ACKNOWLEDGMENTS

Would like to thank the editors of the special issue, Anke Hoeffler and Reetika Khera, the anonymous referees, Justin Sandefur, and Gulzar Natarajan for helping me produce this work.

CONFLICT OF INTEREST STATEMENT

None.

DATA AVAILABILITY STATEMENT

All of the data used in this paper is in the public domain.

ORCID

Lant Pritchett  <https://orcid.org/0000-0002-3736-3118>

REFERENCES

- Altonji, J., Elder, T., & Taber, C. (2005). Selection on observed and unobserved variables: Assessing the effectiveness of Catholic schools. *Journal of Political Economy*, 113(1), 151–184.
- Andrabi, T., Das, J., Khwaja, A., Ozyurt, S., & Singh, N. (2020). Upping the ante: The equilibrium effects of unconditional grants to private schools. *American Economic Review*, 110(10), 3315–3349.
- Andrabi, T., Das, J., & Khwaja, A. (2017). Report cards: The impact of providing school and child test scores on educational markets. *American Economic Review*, 107(6), 1535–63.
- Angrist, N., & Meager, R. (2022). *The role of implementation in generalizability: A synthesis of evidence on targeted educational instruction and a new randomized trial* (CEDIL syntheses working paper, 4). CEDIL.
- Banerjee, A., Duflo, E., & Glennerster, R. (2008). Putting a band-aid on a corpse: Incentives for nurses in the Indian public health care system. *Journal of European Economic Association*, 6(2–3), 487–500.
- Banerjee, A., Karlan, D., & Zinman, J. (2015). Six randomized evaluations of microcredit: Introduction and further steps. *American Economic Journal: Applied Economics*, 7(1), 1–21.
- Bertrand, M., Djankov, S., Hanna, R., & Mullainathan, S. (2007). Obtaining a driver's license in India: An experimental approach to studying corruption. *The Quarterly Journal of Economics*, 122(4), 1639–1676.
- Burde, D., & Linden, L. L. (2013). Bringing education to Afghan girls: A randomized controlled trial of village-based schools. *American Economic Journal: Applied Economics*, 5(3), 27–40.
- Clemens, M. A., & Mendola, M. (2020). *Migration from developing countries: Selection, income elasticity, and Simpson's paradox* (IZA Discussion Papers 13612). Institute of Labor Economics (IZA).
- Clemens, M. A., Montenegro, C. E., & Pritchett, L. (2019). The place premium: Bounding the Price equivalent of migration barriers. *The Review of Economics and Statistics*, 101(2), 201–213.
- Dhaliwal, I., & Hanna, R. (2017). The devil is in the details: The successes and limitations of bureaucratic reform in India. *Journal of Development Economics*, 124, 1–21.
- Filmer, D. (2007). If you build it, will they come? School Availability and school enrolment in 21 poor countries. *The Journal of Development Studies*, 43(5), 901–928.
- Gadamer, H.-G. (1975). *Truth and method*. Continuum.
- Glewwe, P., Kremer, M., & Moulin, S. (2009). Many children left behind? Textbooks and test scores in Kenya. *American Economic Journal: Applied Economics*, 1(1), 112–135.
- Kerwin, J. T., & Thornton, R. L. (2021). Making the grade: The sensitivity of education program effectiveness to input choices and outcome measures. *The Review of Economics and Statistics*, 103(2), 251–264.
- Leamer, E. E. (1983). Let's take the con out of econometrics. *The American Economic Review*, 73(1), 31–43.
- Leamer, E. E. (2010). Tantalus on the road to Asymptopia. *Journal of Economic Perspectives*, 24(2), 31–46.
- Mayer, T., & Zignago, S. (2011). Notes on Cepii's distances measures: The geodist database. *CEPII Working Papers*, 2011–25.
- McKenzie, D., Stillman, S., & Gibson, J. (2010). How important is selection? Experimental vs. non-experimental measures of the income gains from migration. *Journal of the European Economic Association*, 8(4), 913–945.
- Muralidharan, K., & Singh, A. (2020). *Improving public sector management at scale? Experimental evidence on school governance in India* (RISE Working Paper, 20/056). National Bureau of Economic Research.
- Muralidharan, K., & Sundaraman, V. (2015). The aggregate effect of school choice: Evidence from a two-stage experiment in India. *The Quarterly Journal of Economics*, 130(3), 1011–1066.
- Olken, B. A. (2007). Monitoring corruption: Evidence from a field experiment in Indonesia. *Journal of Political Economy*, 115(2), 200–249.
- Oster, E. (2019). Unobservable selection and coefficient stability: Theory and evidence. *Journal of Business and Economic Statistics*, 37(2), 187–204.
- Patel, D., & Sandefur, J. (2020). A Rosetta stone for human capital. Center for Global Development Working Paper 550.
- Pritchett, L. (2020). Why “feigned ignorance” is not good economics (or science generally). *LantRant Blog*.
- Pritchett, L. (2021). *Let's take the con out of randomized control trials in development: The puzzles and paradoxes of external validity, empirically illustrated* (CID Faculty Working Paper Series, 399). Center for International Development.
- Pritchett, L., & Hani, F. (2020). The economics of international wage differentials and migration. In *Oxford research encyclopedias*. Oxford University Press.

- Pritchett, L., & Sandefur, J. (2014). Context matters for size: Why external validity claims and development practice do not mix. *Journal of Globalization and Development*, 42, 161–197.
- Pritchett, L., & Sandefur, J. (2015). Learning from experiments when context matters. *American Economic Review*, 105(5), 471–475.
- Scott, J. (1998). *Seeing like a state: How certain schemes to improve the human condition have failed*. Yale University Press.
- Vivalt, E. (2020). How much can we generalize from impact evaluations? *Journal of the European Economic Association*, 18(6), 3045–3089.
- World Bank. (2020). Cost-effective approaches to improve global learning: What does recent evidence tell us are “smart buys” for improving learning in low and middle income countries. World Bank Group.

How to cite this article: Pritchett, L. (2023). “Rely (only) on the rigorous evidence” is bad advice. *Review of Development Economics*, 1–25. <https://doi.org/10.1111/rode.13037>