

# Principal Component Regression as a Solution to Measurement Error Bias

Isaac Liu, Nicolás Martorell & Paul Opheim

June 4, 2021

## Abstract

We examine the use of Principal Component Regression (PCR) as a solution to bias introduced by measurement error in auxiliary covariates included in a regression. We show its usefulness through econometric theory and then use Monte Carlo simulations to show how it provides benefits for different parameters and correlations between the true covariate and the variable of interest. We also find that in a classical setting, PCR is outperformed by an instrumental variables approach, though it performs comparably in a situation with correlated errors. We then apply these methods to study the relationship between life expectancy and the level of government involvement in a country's healthcare system.

## Introduction

Many variables of interest in economics are not directly available as empirical data. Instead, economists often use other variables that are imperfect measurements of the true focus of their analysis. These available variables are known as *proxies* or “variables measured with error”, and, if they suffer from classical measurement error, are relevant for the model specification, and are correlated with the variable of interest, their use biases the coefficient of the variable of interest even if it does not suffer of measurement error. Traditionally, instrumental variables are used to get rid of measurement error induced bias.

As an alternative method of dealing with this problem, we propose the use of Principal Component Analysis (PCA) over several variables measured with error. When there are multiple observed variables driven by a single “true” one, we propose to use PCA over these variables to extract the “true” variable. One may then use this extracted value in a standard OLS regression (often referred to as PCR or Principal Components Regression), thus providing a way to identify the parameter of interest that does not require the assumptions of instrumental variable analysis. The method also allows for more complex and possibly more optimal weightings of mismeasurements relative to simple averaging, and is less vulnerable to the curse of dimensionality relative to the inclusion of many covariates.

This estimator ties into earlier literature considering the intersection of factor models and principal components analysis and measurement error and latent variables problems. Somewhat similarly to our methods, Nagasawa (2020) develops the use of a proxy variable to deal with unobserved heterogeneity in nuisance parameters and uses a partial effects method. Differing from most of scenarios considered, Schennach (2016) focuses on nonclassical measurement error and nonlinear cases and notes the usefulness of factor methods and some cases where they are of more use than instrumental variables. Wegge (1996) considers a setting in which measurement error regression models are factor analysis models, with the correct regressors being the factors. Latent factors are uncorrelated with the errors. Focusing on measurement error in the main regressor, Schofield (2015) combines solutions from structural equations modelling and item response theory to deal with misestimation. Finally, Heckman, Schennach and Williams (2010) considers a situation similar to ours, except involving matching estimators. Without correction, matching estimators can be harmed by mismeasured conditioning variables. However, average treatment effects can be identified using factor proxies, and without need for normalization.

In this paper, we present a theoretical framework and a Monte-Carlo analysis in order to show the properties and behavior of our estimator on large samples under standard assumptions. Additionally, we explore a basic empirical application of our method, by estimating the relationship between economic development on life expectancy at birth. Since there is no consensus on how to measure economic development, we take a sample of different variables that may measure economic development with error (GDP per capita, GNI per capita, Income per Employed Person, among others) over which we apply PCA to estimate coefficients. Our estimator generally behaves as expected in this empirical setting, though it is unclear whether it performs any better or worse than the direct inclusion of covariates, their averaging, or the instrumentation of mismeasured variables with each other.

## Theoretical framework

Consider a model where the outcome is denoted by  $y_i$ . This outcome depends on a variable of interest denoted by  $t_i$  and a vector of covariates denoted by  $X_i = (x_{i,1}, x_{i,2}, \dots, x_{i,p})'$ . Additionally, consider a vector of variables  $X_i^* = (x_{i,1}^*, x_{i,2}^*, \dots, x_{i,p}^*)'$  that correspond to the covariates  $X_i$  but observed with measurement error, where  $x_{i,k}^* = x_{i,k} + \eta_{i,k}$  with  $\eta_{i,k} \sim iid(0, \sigma_{\eta_k}^2)$ ,  $E(x_{i,k}'\eta_{i,k}) = 0, \forall i$ ,  $E(x_{i,k}'\eta_{j,l}) = 0, \forall i \neq j$  and  $k \neq l$ , and  $E(\eta_{i,k}'\eta_{j,l}) = 0, \forall i \neq j$  and  $k \neq l$ .

Therefore, each  $x_{i,k}^*$  suffers classical measurement error. Note that  $E(x_{i,k}) = E(x_{i,k}^*) = \mu_{x_k}$  and that  $V(x_{i,k}) = \sigma_{x_k}^2$  while  $V(x_{i,k}^*) = \sigma_{x_k}^2 + \sigma_{\eta_k}^2 \geq \sigma_{x_k}^2$ .

## Data Generation Process

Assume that the outcome  $y_i$  is determined by the following Data Generation Process (DGP):

$$y_i = \gamma t_i + X_i' \beta + \epsilon_i \quad (1)$$

where  $\gamma$  is the parameter of the variable of interest  $t_i$ ,  $\beta = (\beta_1, \beta_2, \dots, \beta_p)'$  is the vector of the parameters of the covariates  $X_i$  including a constant and  $\epsilon_i \sim \text{iid}(0, \sigma_\epsilon^2)$ . Under this specification, the coefficients are such that:

$$\begin{pmatrix} \gamma \\ \beta \end{pmatrix} = \begin{pmatrix} \sigma_t^2 & \Sigma_{tX} \\ \Sigma_{Xt} & \Sigma_X \end{pmatrix}^{-1} \begin{pmatrix} \Sigma_{yt} \\ \Sigma_{yX} \end{pmatrix} \quad (2)$$

Suppose that the econometrician has access to  $t_i$  but, instead of  $X_i$  she observes  $X_i^*$ . Then, she specifies the following linear model:

$$y_i = \gamma^* t_i + X_i^{*'} \beta^* + \zeta_i \quad (3)$$

The coefficients would be such that:

$$\begin{pmatrix} \gamma^* \\ \beta^* \end{pmatrix} = \begin{pmatrix} \sigma_t^2 & \Sigma_{tX} \\ \Sigma_{Xt} & \Sigma_X + \Sigma_\eta \end{pmatrix}^{-1} \begin{pmatrix} \sigma_t^2 & \Sigma_{tX} \\ \Sigma_{Xt} & \Sigma_X \end{pmatrix} \begin{pmatrix} \gamma \\ \beta \end{pmatrix} \quad (4)$$

Without loss of generality, assume that the abovementioned DGP consists only of two variables as follows:

$$y_i = \gamma^* t_i + \beta^* x_i^* + \zeta_i \quad (5)$$

Then the coefficient of our variable of interest will be biased.

**Claim 1.** *Under measurement error in the covariates, the coefficient is such that:*

$$\gamma^* = \gamma + \beta \frac{\text{cov}(t, x)(\sigma_{x^*}^2 - \sigma_x^2)}{\sigma_t^2 \sigma_{x^*}^2 - \text{cov}(t, x)^2} \quad (6)$$

**Proof.** See Theory Appendix. ■

From Claim 1 it is clear that when  $\text{cov}(t, x) \neq 0$  and that  $x$  is measured with error (i.e.  $\sigma_{x^*}^2 > \sigma_x^2$ ), the coefficient of our variable of interest is biased. If  $t$  and  $x$  are independent, then measurement error in  $x^*$  does not cause any bias. If there is no measurement error in  $x^*$ , then  $\sigma_{x^*}^2 = \sigma_x^2$  and so we would not face any kind of bias, as one would expect.

Claim 1 also allows us to know the direction of the bias. Given that we are facing measurement error in the covariate,  $\sigma_{x^*}^2 > \sigma_x^2$  which implies  $\sigma_{x^*}^2 - \sigma_x^2 > 0$ . Also, it follows from the *Cauchy-Schwarz* inequality that the denominator is also positive. Then, the direction of the bias will depend on the sign of  $\beta$  and the covariance of  $t$  and  $x$ , as Table 1 illustrates.

Table 1: Direction of the Bias due to Measurement Error in the Covariate

	$\beta > 0$	$\beta < 0$
$\text{cov}(t, x) > 0$	upward-biased	downward-biased
$\text{cov}(t, x) < 0$	downward-biased	upward-biased

## Principal Component Regression as a Bias Correction Method

The classical solution for the measurement-error induced bias in econometrics has been the usage of instrumental variables. Suppose we use as an instrument  $Z_i$  another measure of  $X_i$  so that

$$Z_i = X_i + \omega_i \quad (7)$$

where  $E(\omega_i) = 0$ ,  $\text{Cov}(\epsilon_i, \omega_i) = 0$  and that  $\omega_i$  brings new information so that  $\text{Cov}(\eta_i, \omega_i) = 0$ . Under these conditions,  $Z_i$  is a valid instrument.

**Claim 2.** Suppose  $Z_i = X_i + \omega_i$ . If  $E(\omega_i) = 0$ ,  $\text{Cov}(\epsilon_i, \omega_i) = 0$  and  $\text{Cov}(\eta_i, \omega_i) = 0$ . Then

$$E(Z_i \zeta_i) = 0 \text{ and } E(Z_i X_i) \neq 0 \quad (8)$$

And so  $\gamma$  can be identified through IV regression.

**Proof.** See Theory Appendix. ■

Alternatively, we propose an alternative bias-correction method when there are several mismeasured variables for each covariate; that is when we have more than one  $x_{i,k}^*$  for every  $x_{i,k}$ . Given that in all the mismeasured variables the underlying value is the real value, one could think of extracting the underlying true  $x_{i,k}$  through a linear combination of the different  $x_{i,k}^*$ . Then, we could treat all the  $x_{i,k}^*$  as variables that share components as follows:

$$h_j = \underset{h'h=1, h'h_1=0, \dots, h'h_{j-1}=0}{\text{argmax}} \quad \text{var} [h'X_k^*] \quad (9)$$

where  $h_j$  is the eigenvector of  $\Sigma$  associated with the  $j^{\text{th}}$  ordered eigenvalue  $\lambda_j$  of  $\Sigma_{X_k^*}$ , and the principal components of  $X_k^*$  are  $U_j = h_j'X_k^*$ , where  $h_j$  is the eigenvector of  $\Sigma$  associated with the  $j^{\text{th}}$  ordered eigenvalue  $\lambda_j$  of  $\Sigma$ . Then, we could then retrieve the vector of true variables  $X_i$ .

**Claim 3.**

$$X_i = HX_i^* \quad (10)$$

where  $H$  is a matrix composed of the  $h_k$  vectors of eigenvalues of  $x_{i,k}$ ,  $\forall i, k$ .

**Proof.** See Theory Appendix. ■

Our new linear model then would be:

$$y_i = \gamma^{\text{PCR}} t_i + HX_i^{*'} \beta^{\text{PCR}} + \epsilon_i \quad (11)$$

where  $\gamma$  is identified.

**Claim 4.** Consider equation (11). Then

$$\gamma^{PCR} = \gamma \quad (12)$$

**Proof.** See Theory Appendix. ■

Note that according to Claim 3, the true variable  $x_{i,k}$  is a linear combination of the mismeasured variables that the researcher may have, where the weights are such that equation (9) is satisfied. This allows us to think about other linear combinations that could be used as a bias-correction method which need not satisfy (9).

In particular, take the case in which  $h_k = (\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n})$ , where  $h_k$  is a row vector of dimension  $(1 \times J)$ , and  $J$  is the number of mismeasured variables for  $x_{i,k}$ . Then, Claim 3 will be

$$\tilde{x}_{i,k} = \begin{pmatrix} \frac{1}{n} & \frac{1}{n} & \dots & \frac{1}{n} \end{pmatrix} \begin{pmatrix} x_{i,1}^* \\ x_{i,2}^* \\ \vdots \\ x_{i,J}^* \end{pmatrix} \quad (13)$$

$$= \frac{1}{n} \sum_{j=1}^J x_{i,j}^* \quad (14)$$

That is, the average of the mismeasured variables for  $x_{i,k}$  is a feasible linear combination that may correct for the mismeasurement bias problem.

## Properties of the Estimator: Monte Carlo Simulations

We complement our theoretical analysis by using Monte Carlo Simulation to analyze the effects of using Principal Components Regression as a method of bias correction. For these simulations, we assume that the true DGP for the data is:

$$y_i = \beta_1 x_i + \beta_2 z_i + u_i$$

where  $x_i$  and  $z_i$  are single variables drawn from  $\mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}\right)$  and where  $\rho$  is the covariance between our main variable of interest  $x_i$  and the covariate  $z_i$ . The  $u_i$  is drawn from a white noise distribution ( $\mathcal{N}(0,1)$ ) that is uncorrelated with both  $x_i$  and  $z_i$ . We then assume (as with the theoretical analysis) that  $z_i$  is not directly observable and instead the researchers only have access to  $p$  measurements  $z_{i,j}^*$  where  $z_{i,j}^* = z_i + \eta_j$ , where  $\eta_j$  is drawn from a white noise distribution  $\mathcal{N}(\mathbf{0}, \Sigma)$ , where  $\mathbf{0}$  is a  $p$ -vector, and where  $\Sigma$  is the  $p$  by  $p$  identity matrix.

In our simulations, we assume default values of  $\rho = 0.5$ ,  $\beta_1 = \beta_2 = 1$ , and  $p = 5$ . We then vary each factor while holding the others fixed, and perform 1,000 simulations of the DGP followed by an OLS regression on either the PCA value from the  $p$  measurements of the true  $z_i$ , or on a single one of the measurements of  $z_i$ . For each simulation, we generate 3,000 observations of  $y_i, x_i$ , etc. These values were chosen to best match the empirical application later in this paper. The first two rows of each panel in Table 2 show the results for different values of  $\rho$ .

Table 2: Average Coefficients for Values of  $\rho$  ( $N = 3,000$ , No Exponential Transformation)

	$\rho$ Value				
	-0.9	-0.5	0	0.5	0.9
	<i>Coefficient on Main Variable</i>				
PCA	0.538 (0.034)	0.895 (0.022)	1.0 (0.019)	1.106 (0.022)	1.462 (0.034)
Single Measurement	0.244 (0.025)	0.714 (0.023)	0.999 (0.022)	1.286 (0.022)	1.756 (0.026)
All Measurements	0.538 (0.033)	0.895 (0.023)	1.0 (0.019)	1.106 (0.022)	1.462 (0.034)
Average of Measurements	0.538 (0.034)	0.895 (0.022)	1.0 (0.019)	1.106 (0.022)	1.462 (0.034)
Instrumental Variable	0.994 (0.116)	0.999 (0.048)	1.0 (0.037)	1.0 (0.039)	1.002 (0.115)
	<i>Mean Absolute Percentage Error</i>				
PCA	46.2%	10.5%	1.5%	10.6%	46.2%
Single Measurement	75.6%	28.6%	1.7%	28.6%	75.6%
All Measurements	46.2%	10.5%	1.5%	10.6%	46.2%
Average of Measurements	46.2%	10.5%	1.5%	10.6%	46.2%
Instrumental Variable	8.6%	3.3%	2.7%	3.0%	8.7%
Simulations	1,000	1,000	1,000	1,000	1,000

We first note that when  $\rho = 0.5$  or  $0.9$  then the coefficient on the variable of interest is artificially inflated when we use a single mismeasurement as a covariate (on average, for  $\rho = 0.5$ , it is roughly 1.29 instead of the true value of 1). Conversely, when  $\rho = -0.5$  or  $-0.9$  then the coefficient is artificially deflated. Using the PCA value as the covariate reduces this bias for both directions, and brings the main coefficient closer to its true value of 1.0. These results are consistent with our theoretical section, where we argued that a positive covariance between the main variable and the true covariate will lead to an inflation on the main coefficient, while a negative covariance will lead to a deflation of the coefficient. Separately, there is no bias when  $\rho = 0$ , as predicted. Since there no bias to correct, we do not see gains from using the PCA covariate method for that particular  $\rho$  value. These simulation results suggest that using PCR is more effective than a single mismeasured covariate, although there are no gains to using it when the covariance between the covariate and the main variable of interest is close to 0.

However, the performance advantages that we see from using PCR could be driven by the benefit of having multiple measurements of our true covariate of interest, as opposed to any special advantages from PCR specifically. We test this question by comparing the estimated  $\beta_1$  in our PCR regressions with the estimated  $\beta_1$  when we include all  $p$  measurements as separate covariates in the regression, and the  $\beta_1$  obtained when the covariate is the mean of all  $p$  measurements of the true covariate. We also show the results for an instrumental variables regression where we use other measurements of the true covariate as an instrument for a single measurement of the covariate. The results from these regressions for different values of  $\rho$  are shown in the bottom three rows of each panel of Table 2.

As one can see from these results, there does not seem to be a noticeable difference between using PCR, all measurements, or the average of measurements. However, the instrumental variable regression performs far better than these other methods, correcting for almost all bias and moving the estimated coefficient value close to 1 (although with noticeably larger standard errors than the other methods). Thus, our simulations suggest that, in this framework, there are major benefits to having multiple measurements of a latent covariate of interest. However, PCR does not noticeably improve on two other ways of incorporating these other measurements (taking their average or including all measurements as separate covariates) and performs much worse than using instrumental variables regression with these additional measurements.

We then varied different aspects of this framework in order to better understand the conditions under which

PCR performed especially well relative to the other possible methods for including multiple measurements of the same covariate. In this exploration, we discovered that when the ratio of  $N$  over  $p$  is especially low (there are many measurements relative to the number of observations) then PCA starts to perform better than the instrumental variables method and the method where all measurements are included as separate covariates. We can thus see that PCR does reasonably well under the conditions known as “The Curse of Dimensionality”.

Additionally, we wanted to see how the different methods performed when we used a different framework of measurement error. To do this, we transform half of the covariate measurements by taking the  $z_{i,j}^*$  value from before and transforming it, creating a new  $z'_{i,j}$  value such that  $z'_{i,j} = e^{z_{i,j}^*}$ . This allows us to analyze the performance of our different techniques in a situation where the measurements of the true covariate are on different scales from one another.<sup>1</sup> We found that PCR tends to outperform taking the average of all covariate measurements under this multiple-scale measurement error framework. We can see the effect of both a low  $N/p$  ratio and transformed measurement errors in Table 3:

Table 3: Average Coefficients for Values of  $p$  ( $N = 100$ , Exponential Transformation)

	$p$ Value				
	5	20	50	100	500
	<i>Coefficient on Main Variable</i>				
PCA	1.187 (0.137)	1.081 (0.127)	1.05 (0.122)	1.049 (0.123)	1.038 (0.118)
Single Measurement	1.286 (0.136)	1.289 (0.134)	1.284 (0.128)	1.291 (0.132)	1.282 (0.131)
All Measurements	1.162 (0.135)	1.056 (0.136)	1.019 (0.168)	0.971 (1.055)	0.098 (0.025)
Average of Measurements	1.279 (0.147)	1.194 (0.141)	1.157 (0.137)	1.15 (0.138)	1.139 (0.132)
Instrumental Variable	1.011 (0.213)	1.103 (0.156)	1.2 (0.133)	1.291 (0.132)	1.282 (0.131)
	<i>Mean Absolute Percentage Error</i>				
PCA	19.9%	12.0%	10.6%	10.5%	9.9%
Single Measurement	28.8%	29.1%	28.5%	29.3%	28.4%
All Measurements	17.8%	11.8%	13.4%	68.8%	90.2%
Average of Measurements	28.3%	20.4%	17.4%	16.8%	15.8%
Instrumental Variable	16.4%	15.0%	20.8%	29.3%	28.4%
Simulations	1,000	1,000	1,000	1,000	1,000

We can see here that PCR continues to outperform using a single measurement across all values of  $p$ . Additionally, we can note that when  $N = 100$  and  $p = 5$ , PCR performs better than taking the average of measurements, but worse than the other two methods of using multiple covariates. However, as  $p$  increases, IV and all measurements start to perform worse (as measured by the mean absolute percentage error of their coefficients) while PCR continues to improve (PCR continues to outperform the measurement average for all values of  $p$  on this chart). For  $p = 500$ , PCR’s MAPE is 10%, while all measurements’ is 90%, the average is 16%, and IV is 28%. It is worth noting that in a real-life situation where there were 500 measurements of a single covariate, a researcher could choose to disregard most measurements if adding them to their model reduced performance (as seems to be the case with IV here). This suggests that IV’s performance for  $p = 50, 100$ , and 500 is probably understated in these charts, as the researcher would instead ignore measurements that make the performance of the IV estimator worse. However, it is worth noting that, as measured by MAPE, PCR’s performance when  $p = 500$  (10%) notably outperforms IV’s performance when  $p = 5$  or 10 (16% and 15% respectively). While IV’s average coefficient across simulations is closer to the true value of 1 when  $p = 5$ , its larger standard deviation mean that its MAPE is higher

<sup>1</sup>This is often the case; for example, one could measure income inequality through a Gini coefficient and the percent of income that goes to the top 1% of income-earners, but these variables are on totally different scales from each other.

than PCR under  $p = 20, 50, 100$  or  $500$ . To better understand the relative performance of PCR compared to these other methods, the simulation appendix contains other useful charts.

Overall, these simulations tell us that PCR helps to reduce measurement error-induced bias in OLS regression with a latent covariate relative to including only a single measurement of that covariate. Additionally, within the framework most similar to our empirical application, it seems to perform in line with taking the average of those covariate measurements and using each measurement as a separate covariate in the regression, but worse than using an IV technique. However, there are circumstances where PCR is the best-performing technique of all the ones studied in this simulation section. This suggests that using PCR as a method for reducing measurement error-induced bias in the coefficient on the variable of interest can be helpful, especially under the specific circumstances mentioned earlier in this section.

## **Application: Government Share of Healthcare Spending and Life Expectancy**

We now examine the usage of the principal components estimator in an empirical setting with measurement error. One interesting question in public economics and public health is the study of the relationship between publicly and privately funded healthcare systems and outcomes such as life expectancy. To measure the public or private nature of a healthcare system we use the continuous variable of the government's share of total health expenditure in a given country and year.

Some previous work has covered this relationship. Considering that this topic has been studied in "relatively few papers," Linden and Ray (2017) focus on the relationship between life expectancy at birth and public and private health expenditures for 34 OECD countries from 1970-2012 and find that both public and private health spending are important to life expectancy and are associated with each other. In work similar to ours, Or (2000) predicts premature death in 21 OECD countries from 1970-1992, considering the public share of health expenditure, environmental factors, and GDP. He finds that a larger share of public spending is associated with lower rates of premature mortality for both males and females, and that controlling for GDP is important; it is also associated with less premature mortality. This work also demonstrates the importance of our methods of reducing the number of covariates considered, as it includes many economic variables and fixed effects but examines only several hundred observations; the estimators used may be subject to an significant amount of variance.

In this regression it is important to account for the role of a country's level of economic development. There is an extensive literature documenting the relationship between economic development and life expectancy. Ling et al. (2017) finds that economic growth is associated with increased life expectancy in Malaysia, while considering the reverse causal direction Acemoglu and Johnson (2007) finds improvements in life expectancy lead to little or no growth. Somewhat less obvious is the linkage between government provision of healthcare and development. In general, public goods provision and government spending, including in fields such as healthcare, has been linked to prosperity; low income countries may remain in such a state due to inefficient governments and inferior institutions (Wu, Tang and Lin, 2010).

However, economic development is liable to be measured with error. GDP measurements usually rely on company surveys, and the methodology within a country and for comparisons between countries through exchange rates or PPP adjustments may vary (Grishin, Ustyuzhanina and Pavlovna, 2019). Other sources of error include the presence of the informal economy and non-monetary but productive work, the challenge of accurately measuring the value of digital services which often do not have visible prices, and government incentives to manipulate official statistics about growth (Charmes, 2012; Ahmad, Ribarsky and Reinsdorf, 2017; Nakamura, Steinsson and



Liu, 2016).<sup>2</sup> Hence, this setup, with a covariate in regression subject to measurement error, fits the situation described in the theory and simulations in the previous sections.<sup>3</sup> In this case, we aim to reduce possible bias in the coefficient of the government's share of health spending by making appropriate use of multiple measures of economic development.

Our data on all measures comes from the World Bank (The World Bank, 2021). We remove country-years with missing values for any of the variables summarized in Table 4 and standardize the economic covariates in the 5 uppermost rows by subtracting the mean and dividing by the standard deviation (in order to enable interpretable principal component analysis).

Table 4: Summary Statistics

Variable	Obs	Mean	SD	Min	Med	Max
GDP Per Capita PPP (Current International \$)	3,143	16,443.76	19,173.04	435.08	9,331.99	141,634.96
GDP Per Capita (Current USD)	3,143	11,764	17,582.7	111.93	4,018.95	118,823.65
GNP Per Capita PPP (Current International \$)	3,143	15,980.12	18,478.14	410	9,080	132,440
GNP Per Capita (Current USD)	3,143	11,247.65	16,583.07	110	3,800	104,560
ILO GDP Per Person Employed	3,143	41,782.32	39,743.08	1,371.24	29,220.02	266,103.71
Life Expectancy at Birth (All Population)	3,143	69.55	9.26	39.44	71.78	84.21
Government Share of Health Expenditure	3,143	49.62	21.73	4.06	50.27	95.14

In Table 5 we apply our main set of estimators. The PCA estimator is found in column (1). In column (2), we instead control for just a single measurement of GDP per capita (PPP). In column (3) we directly include the full set of economic covariates (measurements) listed in the 5 uppermost rows of Table 4. In column (4), we use the mean of these mismeasured covariates, and in column (5) we perform instrumental variables regression using all the other economic indicators as an instrument for GDP per capita (PPP).

Table 5: Regressions of Life Expectancy on Government Share of Health Spending

	<i>Life Expectancy at Birth (Years)</i>				
	(1)	(2)	(3)	(4)	(5)
Govt. Share of Health Exp.	0.180*** (0.007)	0.194*** (0.007)	0.166*** (0.007)	0.180*** (0.007)	0.193*** (0.007)
Covariates	PCA	Single Measurement (GDP Per Capita PPP)	All Measurements	Average of Measurements	Instrumental Variable (GDP Per Capita PPP)
Observations	3,143	3,143	3,143	3,143	3,143

Note:

\*p<0.1; \*\*p<0.05; \*\*\*p<0.01

The PCA estimator indicates that a one percentage point increase in the government share of health expenditure is linked to a increase in life expectancy of 0.18 years. Notably, the results generally demonstrate that the methods using multiple measures of the covariates produce different coefficients for government health share, relative to the use of a single mismeasurement. The principal components estimator, the usage of many economic indicators directly, and their mean each produce far smaller coefficients. The IV estimator produces a slightly smaller coefficient.

<sup>2</sup>Due to differences in statistical capacity and the larger relative size of the informal economy, it is possible that mismeasurement of economic development is particularly severe in developing countries. On the other hand, the presence of the digital economy may mean mismeasurement is larger in developed nations. This would constitute the presence of non-classical measurement error, but we only consider classical measurement error in this paper. It is also possible that the interaction of many forms of measurement produces error which is closer to classical assumptions.

<sup>3</sup>It seems less likely that our variable of interest, the government's share of health spending, is measured with error. One would think most governments capable of monitoring their own spending better than economic activity in general. Furthermore, for this variable there would seem to be less governmental incentive to manipulate the statistic relative to GDP or other items.

Moreover, these different coefficients behave in a manner similar to that predicted by our theoretical development and simulations. In Table 2, we saw the impact of variation in  $\rho$ , the correlation between measurements for a values of  $p = 5$  and  $\beta_1, \beta_2 = 1$  for 3,000 observations. In the empirical setting it is difficult to tell what is a reasonable value of  $\beta$ . Nevertheless, we see that for a positive  $\rho$  value between 0 and 1 (as is likely to be the case in light the correlation between GDP and the government share of health spending and overall public goods), the coefficient obtained from using a single measurement is inflated relative to that from PCA, and presumably other methods combining multiple measures as in columns (3), (4), and (5) of Table 5.

Results using univariate OLS (with no covariates), country and year fixed effects models (with country clustered standard errors), and more principal components are in Table 6. Univariate OLS produces a large and inflated coefficient. Fixed effects coefficients greatly reduce the magnitude of any potential causal effects and are insignificant. The results in column (4) also show a reduction in the inflation of coefficients relative to column (2) of Table 5, as the inclusion of more principal components produces a small coefficient very similar to that obtained with just a single component.<sup>4</sup>

Table 6: Additional Regressions

	<i>Life Expectancy at Birth (Years)</i>			
	(1)	(2)	(3)	(4)
Govt. Share of Health Exp.	0.276*** (0.006)	-0.003 (0.011)	-0.003 (0.011)	0.181*** (0.007)
Covariates	None	None	PCA	PC 1-2
Fixed Effects	No	Yes	Yes	No
Observations	3,143	3,143	3,143	3,143
<i>Note:</i>		*p<0.1; **p<0.05; ***p<0.01		

## Conclusion

In this paper, we have explored the usefulness of using Principal Component Regression (PCR) as a method of mitigating the bias on the coefficient for the main variable of interest that is induced by measurement error in an auxiliary covariate. Using theory and Monte Carlo simulations we have shown how this technique improves upon using a single measurement of the true covariate, and we have shown how its performance compares to that of some other methods for mitigating this bias, while mostly falling short of the performance of instrumental variables. We then applied this technique to an empirical question (the relationship between the government share of health expenditure and life expectancy).

Other solutions to measurement error meriting further exploration include the usage of general factor models, which could also involve the exploitation of the panel structure of relevant data. Measurements within time periods or units may provide information about the true value of variables. There are also likely other ways to summarize and control the dimension of the information presented by multiple measurements. Another option could be the usage of dimensionality-reduction techniques within an instrumental variables regression.

<sup>4</sup>This is likely due to the high explanatory power of just the first principal component, as is clear in Appendix Figure 3.

## References

- Acemoglu, Daron, and Simon Johnson.** 2007. "Disease and Development: The Effect of Life Expectancy on Economic Growth." *Journal of Political Economy*, 115(6): 925–985. Publisher: The University of Chicago Press.
- Ahmad, Nadim, Jennifer Ribarsky, and Marshall Reinsdorf.** 2017. "Can potential mismeasurement of the digital economy explain the post-crisis slowdown in GDP and productivity growth?" Publisher: OECD.
- Charmes, Jacques.** 2012. "The Informal Economy Worldwide: Trends and Characteristics." *Margin: The Journal of Applied Economic Research*, 6(2): 103–132. Publisher: SAGE Publications India.
- Grishin, Victor Ivanovich, Elena Vladimirovna Ustyuzhanina, and Irina Pavlovna.** 2019. "Main Problems with Calculating GDP as and Indicator of Economic Health of the County." 9.
- Heckman, James, Susanne Schennach, and Benjamin D Williams.** 2010. "Matching on Proxy Variables." 15.
- Linden, Mikael, and Deb Ray.** 2017. "Life expectancy effects of public and private health expenditures in OECD countries 1970–2012: Panel time series approach." *Economic Analysis and Policy*, 56: 101–113.
- Ling, Chong Hui, Khalid Ahmed, Rusnah Muhamad, Muhammad Shahbaz, and Nanthakumar Loganathan.** 2017. "Testing the Social Cost of Rapid Economic Development in Malaysia: The Effect of Trade on Life Expectancy." *Social Indicators Research*, 130(3): 1005–1023.
- Nagasawa, Kenichi.** 2020. "Identification and Estimation of Partial Effects with Proxy Variables." 23.
- Nakamura, Emi, Jón Steinsson, and Miao Liu.** 2016. "Are Chinese Growth and Inflation Too Smooth? Evidence from Engel Curves." *American Economic Journal: Macroeconomics*, 8(3): 113–144.
- Or, Zeynep.** 2000. "Determinants of Health Outcomes in Industrialised Countries: a Pooled, Cross-country, Time-series Analysis." *OECD Economic Studies*, 25.
- Schennach, Susanne M.** 2016. "Recent Advances in the Measurement Error Literature." *Annual Review of Economics*, 8(1): 341–377.
- Schofield, Lynne Steuerle.** 2015. "Correcting for Measurement Error in Latent Variables Used as Predictors." *The annals of applied statistics*, 9(4): 2133–2152.
- The World Bank.** 2021. "Indicators | Data."
- Wegge, Leon L.** 1996. "Local identifiability of the factor analysis and measurement error model parameter." *Journal of Econometrics*, 70(2): 351–382.
- Wu, Shih-Ying, Jenn-Hong Tang, and Eric S. Lin.** 2010. "The impact of government expenditure on economic growth: How sensitive to the level of development?" *Journal of Policy Modeling*, 32(6): 804–817.

# Theory Appendix

## Short Proofs

**Claim 1.** Under measurement error in the covariates, the coefficient is such that:

$$\gamma^* = \gamma + \beta \frac{\text{cov}(t, x)(\sigma_{x^*}^2 - \sigma_x^2)}{\sigma_t^2 \sigma_{x^*}^2 - \text{cov}(t, x)^2}$$

**Proof.** Equations (2) and (4) will be such that

$$\begin{pmatrix} \gamma^* \\ \beta^* \end{pmatrix} = \begin{pmatrix} \sigma_t^2 & \text{cov}(t, x^*) \\ \text{cov}(x^*, t) & \sigma_{x^*}^2 \end{pmatrix}^{-1} \begin{pmatrix} \text{cov}(y, t) \\ \text{cov}(y, x^*) \end{pmatrix}$$

■

**Claim 2.** Suppose  $Z_i = X_i + \omega_i$ . If  $E(\omega_i) = 0$ ,  $\text{Cov}(\epsilon_i, \omega_i) = 0$  and  $\text{Cov}(\eta_i, \omega_i) = 0$ . Then

$$E(Z_i \zeta_i) = 0 \text{ and } E(Z_i X_i) \neq 0$$

And so  $\gamma$  can be identified through IV regression.

**Proof.** Suppose an instrument  $Z_i$  that satisfies the relevance condition  $E(Z_i' X_i) \neq 0$  and  $E(Z_i' t_i) \neq 0$ , and also the exclusion restriction  $E(Z_i' \epsilon_i) = E(Z_i' \zeta_i) = E(Z_i' \eta_{i,k}) = 0$ , for all  $i$  and  $k$ . Then premultiplying by  $Z_i$  we have

$$Z_i' y_i = Z_i' \gamma^* t_i + Z_i' X_i^* \beta^* + Z_i' \zeta_i$$

and so

$$\begin{aligned} \begin{pmatrix} \gamma^{IV} \\ \beta^{IV} \end{pmatrix} &= \begin{pmatrix} \Sigma_{Zt} & \Sigma_{ZX, Zt} \\ \Sigma_{Zt, ZX} & \Sigma_{ZX} + \Sigma_{Z\eta} \end{pmatrix}^{-1} \begin{pmatrix} \Sigma_{Zt} & \Sigma_{ZX, Zt} \\ \Sigma_{Zt, ZX} & \Sigma_{ZX} \end{pmatrix} \begin{pmatrix} \gamma \\ \beta \end{pmatrix} \\ &= \begin{pmatrix} \Sigma_{Zt} & \Sigma_{ZX, Zt} \\ \Sigma_{Zt, ZX} & \Sigma_{ZX} \end{pmatrix}^{-1} \begin{pmatrix} \Sigma_{Zt} & \Sigma_{ZX, Zt} \\ \Sigma_{Zt, ZX} & \Sigma_{ZX} \end{pmatrix} \begin{pmatrix} \gamma \\ \beta \end{pmatrix} \\ \begin{pmatrix} \gamma^{IV} \\ \beta^{IV} \end{pmatrix} &= \begin{pmatrix} \gamma \\ \beta \end{pmatrix} \end{aligned}$$

However, finding a reliable source of exogeneity is sometimes difficult, as is demonstrating the exclusion restriction.

Suppose now that as instrument we have another measure of  $X_i$  so that

$$Z_i = X_i + \omega_i$$

where  $E(\omega_i) = 0$ ,  $\text{Cov}(\epsilon_i, \omega_i) = 0$  and that  $\omega_i$  brings new information so that  $\text{Cov}(\eta_i, \omega_i) = 0$ . Then, if  $Z_i$

satisfies exogeneity and relevance we will be able to identify the parameters without any bias. In fact:

$$\begin{aligned}
E(Z_i \zeta_i) &= E(Z_i(\epsilon_i - \eta_i \beta)) \\
&= E(Z_i \epsilon_i) - E(Z_i \eta_i) \beta \\
&= E((X_i + \omega_i) \epsilon_i) - E((X_i + \omega_i) \eta_i) \beta \\
&= E(X_i \epsilon_i) + E(\omega_i \epsilon_i) - (E(X_i \eta_i) + E(\omega_i \eta_i)) \beta
\end{aligned}$$

And so  $Z_i$  is exogenous. Given  $Z_i = X_i + \omega_i$  it is clear that  $E(Z_i X_i) \neq 0$  and so relevance is also satisfied. Thus,  $\gamma$  and  $\beta$  may be identified using this kind of instrument. ■

**Claim 3.**

$$X_i = H X_i^*$$

where  $H$  is a matrix compound of the  $h_k$  vectors of eigenvalues of  $x_{i,k}$ ,  $\forall i, k$

**Proof.** Under our assumptions, the vector of mismeasured values  $X_k^*$  of  $x_{i,k}$ , share only one principal component which is precisely  $x_{i,k}$ . Then, we only have one principal component,  $x_{i,k}$ , and so the  $x_{i,k}$  is such that

$$x_{i,k} = h_k' X_k^*$$

Finally, we could then retrieve the vector of true variables  $X_i$

$$X_i = H X_i^*$$

where  $H$  is a matrix such that

$$H = \begin{pmatrix} h_1 & 0 & 0 & \dots & 0 \\ 0 & h_2 & 0 & \dots & 0 \\ \vdots & \ddots & h_3 & \ddots & \vdots \\ 0 & \dots & \dots & \dots & h_p \end{pmatrix}$$

and  $h_k$  is the vector of eigenvalues for the variable  $x_{i,k}$ . ■

**Claim 4.** Consider equation (11). Then

$$\gamma^{PCR} = \gamma$$

**Proof.** The coefficients are as follows

$$\begin{aligned}
\begin{pmatrix} \gamma^{PCR} \\ \beta^{PCR} \end{pmatrix} &= \begin{pmatrix} \sigma_t^2 & \Sigma_{t,HX^*} \\ \Sigma_{HX^*,t} & \Sigma_{HX^*} \end{pmatrix}^{-1} \begin{pmatrix} \Sigma_{yt} \\ \Sigma_{y,HX^*} \end{pmatrix} \\
&= \begin{pmatrix} \sigma_t^2 & \Sigma_{t,HX^*} \\ \Sigma_{HX^*,t} & \Sigma_{HX^*} \end{pmatrix}^{-1} \begin{pmatrix} \sigma_t^2 & \Sigma_{tX} \\ \Sigma_{Xt} & \Sigma_X \end{pmatrix} \begin{pmatrix} \gamma \\ \beta \end{pmatrix} \\
&= \begin{pmatrix} \gamma \\ \beta \end{pmatrix}
\end{aligned}$$

where the second equality comes from Equation (4) and Claim 3 and the last equality comes from Claim 3. ■

## Simulation Appendix

Table 7: Average Coefficients for Values of  $p$  ( $N = 100$ , No Exponential Transformation)

	$p$ Value				
	5	20	50	100	500
	<i>Coefficient on Main Variable</i>				
PCA	1.105 (0.124)	1.028 (0.121)	1.012 (0.117)	1.012 (0.114)	0.996 (0.118)
Single Measurement	1.283 (0.134)	1.289 (0.127)	1.286 (0.126)	1.288 (0.129)	1.281 (0.13)
All Measurements	1.104 (0.127)	1.028 (0.135)	1.012 (0.175)	1.062 (1.168)	0.156 (0.029)
Average of Measurements	1.104 (0.124)	1.028 (0.121)	1.012 (0.117)	1.012 (0.114)	0.996 (0.118)
Instrumental Variable	1.027 (0.243)	1.095 (0.157)	1.2 (0.134)	1.288 (0.129)	1.281 (0.13)
	<i>Mean Absolute Percentage Error</i>				
PCA	13.2%	9.9%	9.4%	9.3%	9.4%
Single Measurement	28.4%	28.9%	28.8%	28.9%	28.3%
All Measurements	13.3%	10.8%	13.8%	72.2%	84.4%
Average of Measurements	13.2%	9.8%	9.4%	9.3%	9.4%
Instrumental Variable	18.2%	14.9%	20.9%	28.9%	28.3%
Simulations	1,000	1,000	1,000	1,000	1,000

Table 8: Average Coefficients for Values of  $p$  ( $N = 1,000$ , No Exponential Transformation)

	$p$ Value				
	5	20	50	100	500
	<i>Coefficient on Main Variable</i>				
PCA	1.104 (0.038)	1.032 (0.038)	1.012 (0.037)	1.007 (0.036)	1.002 (0.037)
Single Measurement	1.284 (0.04)	1.287 (0.041)	1.285 (0.04)	1.286 (0.04)	1.287 (0.04)
All Measurements	1.104 (0.038)	1.032 (0.039)	1.011 (0.038)	1.007 (0.039)	1.001 (0.052)
Average of Measurements	1.104 (0.038)	1.032 (0.038)	1.012 (0.037)	1.007 (0.036)	1.002 (0.037)
Instrumental Variable	0.999 (0.076)	1.011 (0.054)	1.029 (0.052)	1.059 (0.047)	1.201 (0.042)
	<i>Mean Absolute Percentage Error</i>				
PCA	10.4%	4.1%	3.1%	2.9%	2.9%
Single Measurement	28.4%	28.7%	28.5%	28.6%	28.7%
All Measurements	10.4%	4.2%	3.1%	3.1%	4.2%
Average of Measurements	10.4%	4.1%	3.1%	2.9%	2.9%
Instrumental Variable	5.4%	4.4%	4.8%	6.4%	20.1%
Simulations	1,000	1,000	1,000	1,000	1,000

Table 9: Average Coefficients for Values of  $p$  ( $N = 3,000$ , No Exponential Transformation)

	$p$ Value				
	5	20	50	100	500
	<i>Coefficient on Main Variable</i>				
PCA	1.106 (0.022)	1.03 (0.021)	1.013 (0.021)	1.007 (0.021)	1.0 (0.02)
Single Measurement	1.286 (0.022)	1.285 (0.023)	1.285 (0.022)	1.286 (0.024)	1.285 (0.023)
All Measurements	1.106 (0.022)	1.03 (0.021)	1.013 (0.021)	1.007 (0.022)	1.0 (0.022)
Average of Measurements	1.106 (0.022)	1.03 (0.021)	1.013 (0.021)	1.007 (0.021)	1.0 (0.02)
Instrumental Variable	1.0 (0.039)	1.004 (0.031)	1.011 (0.029)	1.022 (0.029)	1.089 (0.026)
	<i>Mean Absolute Percentage Error</i>				
PCA	10.6%	3.2%	1.9%	1.8%	1.6%
Single Measurement	28.6%	28.5%	28.5%	28.6%	28.5%
All Measurements	10.6%	3.2%	1.9%	1.8%	1.8%
Average of Measurements	10.6%	3.2%	1.9%	1.8%	1.6%
Instrumental Variable	3.0%	2.5%	2.5%	2.9%	8.9%
Simulations	1,000	1,000	1,000	1,000	1,000

Table 10: Average Coefficients for Values of  $p$  ( $N = 1,000$ , Exponential Transformation)

	$p$ Value				
	5	20	50	100	500
	<i>Coefficient on Main Variable</i>				
PCA	1.204 (0.042)	1.083 (0.039)	1.054 (0.037)	1.048 (0.037)	1.039 (0.039)
Single Measurement	1.287 (0.039)	1.287 (0.041)	1.285 (0.039)	1.287 (0.04)	1.287 (0.042)
All Measurements	1.171 (0.038)	1.056 (0.038)	1.023 (0.037)	1.015 (0.037)	1.002 (0.052)
Average of Measurements	1.315 (0.052)	1.22 (0.053)	1.186 (0.049)	1.176 (0.05)	1.164 (0.053)
Instrumental Variable	1.003 (0.059)	1.014 (0.052)	1.029 (0.049)	1.061 (0.047)	1.201 (0.044)
	<i>Mean Absolute Percentage Error</i>				
PCA	20.4%	8.3%	5.6%	5.2%	4.5%
Single Measurement	28.7%	28.7%	28.5%	28.7%	28.7%
All Measurements	17.1%	5.9%	3.5%	3.2%	4.2%
Average of Measurements	31.5%	22.0%	18.6%	17.6%	16.4%
Instrumental Variable	4.7%	4.3%	4.5%	6.5%	20.1%
Simulations	1,000	1,000	1,000	1,000	1,000



Table 11: Average Coefficients for Values of  $p$  ( $N = 3,000$ , Exponential Transformation)

	<i>p Value</i>				
	5	20	50	100	500
	<i>Coefficient on Main Variable</i>				
PCA	1.208 (0.028)	1.08 (0.023)	1.054 (0.023)	1.045 (0.022)	1.037 (0.021)
Single Measurement	1.285 (0.024)	1.285 (0.024)	1.285 (0.023)	1.286 (0.023)	1.286 (0.023)
All Measurements	1.173 (0.024)	1.053 (0.022)	1.023 (0.022)	1.013 (0.022)	1.002 (0.023)
Average of Measurements	1.325 (0.036)	1.223 (0.035)	1.191 (0.034)	1.18 (0.033)	1.169 (0.032)
Instrumental Variable	1.001 (0.034)	1.003 (0.031)	1.01 (0.029)	1.02 (0.028)	1.091 (0.026)
	<i>Mean Absolute Percentage Error</i>				
PCA	20.8%	8.0%	5.4%	4.5%	3.8%
Single Measurement	28.5%	28.5%	28.5%	28.6%	28.6%
All Measurements	17.3%	5.3%	2.7%	2.0%	1.8%
Average of Measurements	32.5%	22.3%	19.1%	18.0%	16.9%
Instrumental Variable	2.7%	2.5%	2.5%	2.8%	9.1%
Simulations	1,000	1,000	1,000	1,000	1,000

# Application Appendix

Figure 1: Correlations Between Variables



Figure 2: Economic Measures PCA Loadings



Figure 3: Economic Measures PCA Share of Variance Explained

