

# Attenuation Bias, Measurement Error & Principal Component Analysis

Isaac Liu, Nicolás Martorell & Paul Opheim

May 19, 2021

## **Abstract**

Shorter version of the abstract (I would say 4-5 sentences in a single paragraph max) goes here

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, their use causes *attenuation bias* when they are used as independent variables in econometric estimation. Traditionally, instrumental variables are used as a shock of exogeneity to get rid of this bias, but finding truly exogenous variables that satisfy the exclusion restriction is difficult, and so this method can often not be feasibly applied.

As an alternative to dealing with attenuation bias, 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. We then use this extracted value and use it in a standard OLS regression, thus providing a solution to attenuation bias that does not require the strong assumptions of instrumental variable analysis.

To show the properties and behaviour of our estimator on large samples under standard assumptions, we present a theoretical framework and a Monte-Carlo analysis. Additionally, we explore a basic empirical application to our method, by estimating the effect of 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, Household Income Per Capita, among others) over which we apply PCA to apply our identification strategy.

## Literature

Brief discussion of <https://warwick.ac.uk/fac/soc/economics/staff/knagasawa/PartialEffects.pdf>, as well as anything else important that comes up on Google Scholar

## 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 from 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 Generating 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  $\epsilon_i \sim iid(0, \sigma_\epsilon^2)$ . 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

$$\begin{pmatrix} \gamma^* \\ \beta^* \end{pmatrix} = \begin{pmatrix} \sigma_t^2 & \Sigma_{tX^*} \\ \Sigma_{X^*t} & \Sigma_{X^*} \end{pmatrix}^{-1} \begin{pmatrix} \Sigma_{yt} \\ \Sigma_{yX^*} \end{pmatrix} \quad (4)$$

$$= \begin{pmatrix} \sigma_t^2 & \Sigma_{tX} \\ \Sigma_{Xt} & \Sigma_X + \Sigma_\eta \end{pmatrix}^{-1} \begin{pmatrix} \sigma_t^2 & \Sigma_{tX} \\ \Sigma_{Xt} & \Sigma_X \end{pmatrix} \begin{pmatrix} \gamma \\ \beta \end{pmatrix} \quad (5)$$

To see the implications of the of this measurement error in the covariates, consider a simple case where the DGP depends only of the variable of interest and a covariate such that:

$$\begin{pmatrix} \gamma \\ \beta \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad (6)$$

and with  $\sigma_t^2 = \Sigma_X = \Sigma_\eta = 1$  while  $\Sigma_{Xt} = 0.6$ . Then

$$\begin{aligned} \begin{pmatrix} \gamma^* \\ \beta^* \end{pmatrix} &= \begin{pmatrix} 1 & 0.6 \\ 0.6 & 2 \end{pmatrix}^{-1} \begin{pmatrix} 1 & 0.6 \\ 0.6 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \\ \begin{pmatrix} \gamma^* \\ \beta^* \end{pmatrix} &= \begin{pmatrix} 1.37 \\ 0.39 \end{pmatrix} \end{aligned}$$

Clearly, both coefficients shows bias when the econometrician assumes a DGP with  $X_i^*$ : while there is attenuation bias on the coefficient of the covariate, the coefficient of the variable of interest is biased upward given that some of the effect of the covariates is “omitted” given this attenuation.

## Instrumental Variables 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 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

$$Z_i' y_i = Z_i' \gamma^* t_i + Z_i' X_i^{*'} \beta^* + Z_i' \zeta_i \quad (7)$$

and so

$$\begin{pmatrix} \gamma^{IV} \\ \beta^{IV} \end{pmatrix} = \begin{pmatrix} \Sigma_{Zt} & \Sigma_{ZX,Zt} \\ \Sigma_{Zt,ZX} & \Sigma_{ZX} + \Sigma_{Z\eta} \end{pmatrix}^{-1} \begin{pmatrix} \Sigma_{Zt} & \Sigma_{ZX,Zt} \\ \Sigma_{Zt,ZX} & \Sigma_{ZX} \end{pmatrix} \begin{pmatrix} \gamma \\ \beta \end{pmatrix} \quad (8)$$

$$= \begin{pmatrix} \Sigma_{Zt} & \Sigma_{ZX,Zt} \\ \Sigma_{Zt,ZX} & \Sigma_{ZX} \end{pmatrix}^{-1} \begin{pmatrix} \Sigma_{Zt} & \Sigma_{ZX,Zt} \\ \Sigma_{Zt,ZX} & \Sigma_{ZX} \end{pmatrix} \begin{pmatrix} \gamma \\ \beta \end{pmatrix} \quad (9)$$

$$\begin{pmatrix} \gamma^{IV} \\ \beta^{IV} \end{pmatrix} = \begin{pmatrix} \gamma \\ \beta \end{pmatrix} \quad (10)$$

However, finding a reliable source of exogeneity is difficult, and it is impossible to conclusively prove a suitable exclusion restriction. The use of IV as a bias-correction method is thus often unfeasible.

## Principal Component Regression as Bias-Correction Method

Alternatively, we propose an alternative bias-correction method for 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:

$$h_j = \underset{h'h=1, h'h_1=0, \dots, h'h_{j-1}=0}{\operatorname{argmax}} \operatorname{var}[h'X_k^*] \quad (11)$$

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  $\Sigma$ .

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

$$x_{i,k} = h_k'X_k^* \quad (12)$$

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

$$X_i = HX_i^* \quad (13)$$

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  $x_{i,k}$ .

Our new linear model then becomes

$$y_i = \gamma^{PCR}t_i + HX_i^*\beta^{PCR} + \epsilon_i \quad (14)$$

where the coefficients are as follows

$$\begin{pmatrix} \gamma^{PCR} \\ \beta^{PCR} \end{pmatrix} = \begin{pmatrix} \sigma_t^2 & \Sigma_{t, HX^*} \\ \Sigma_{HX^*, t} & \Sigma_{HX^*} \end{pmatrix}^{-1} \begin{pmatrix} \Sigma_{yt} \\ \Sigma_{y, HX^*} \end{pmatrix} \quad (15)$$

$$= \begin{pmatrix} \sigma_t^2 & \Sigma_{t, HX^*} \\ \Sigma_{HX^*, t} & \Sigma_{HX^*} \end{pmatrix}^{-1} \begin{pmatrix} \sigma_t^2 & \Sigma_{tX} \\ \Sigma_{Xt} & \Sigma_X \end{pmatrix} \begin{pmatrix} \gamma \\ \beta \end{pmatrix} \quad (16)$$

$$= \begin{pmatrix} \gamma \\ \beta \end{pmatrix} \quad (17)$$

where the last equality comes from (13).

# Properties of the Estimator: Monte Carlo Simulations

We then 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  $\mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}\right)$ , where  $\rho$  is some 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  $p$  many 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  $\Sigma$  is a diagonal  $p$  by  $p$  matrix with only 1s on the diagonal.

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 100 observations of  $y_i, x_i$ , etc. Below are the results for different values of  $p$ :

	<i>Number of p</i>			
	5	10	20	50
<i>Coefficient on Main Variable</i>				
PCA	1.105 (0.121)	1.066 (0.122)	1.033 (0.119)	1.022 (0.117)
Single Measurement	1.280 (0.124)	1.283 (0.129)	1.282 (0.131)	1.292 (0.167)
<i>Absolute Percentage Error</i>				
PCA	13.1% (9.3 ppts)	11.1% (8.3 ppts)	10.0% (7.3 ppts)	9.4% (7.3 ppts)
Single Measurement	28.2% (12.6 ppts)	28.5% (12.6 ppts)	28.3% (12.6 ppts)	29.3% (12.7 ppts)
Observations	1,000	1,000	1,000	1,000

We can see that using PCA to extract the latent covariate driving the mismeasured covariates noticeably outperforms using a single mismeasured covariate across several values of  $p$ . Both the average coefficient on  $\beta_1$  obtained when including the PCA output in the regression, and the mean absolute percentage error obtained on the 1,000 simulations are both much closer to the target values with the PCA-based regression than with the single measurement regression. Additionally, we can see that as  $p$  increases the estimated  $\beta_1^*$  coefficient in the PCA regression gets steadily closer to the true  $\beta_1$  value of 1. Appendix 1 contains charts that show that this increase in performance

is also true for different values of  $\beta_1$  and  $\beta_2$ .

However, there are certain circumstances where the PCA method does not lead to more accurate estimates of  $\beta_1^*$ . Let's now look at the simulation results for different values of  $\rho$  (the covariance between the main variable of interest  $x_i$  and the true latent covariate  $z_i$ ):

	$\rho$ Value				
	-1	-0.5	0	0.5	1
<i>Coefficient on Main Variable</i>					
PCA	-0.006 (0.238)	0.900 (0.120)	0.996 (0.111)	1.105 (0.121)	2.009 (0.242)
Single Measurement	-0.002 (0.142)	0.720 (0.130)	0.998 (0.127)	1.280 (0.129)	2.003 (0.147)
<i>Absolute Percentage Error</i>					
PCA	100.6% (23.8 ppts)	12.7% (9.1 ppts)	8.9% (6.6 ppts)	13.1% (9.3 ppts)	100.9% (24.2 ppts)
Single Measurement	100.2% (14.2 ppts)	28.1% (12.7 ppts)	10.2% (7.6 ppts)	28.2% (12.6 ppts)	100.3% (14.7 ppts)
Observations	1,000	1,000	1,000	1,000	

When the covariance between  $x_i$  and  $z_i$  is equal to 0,  $-1$ , or  $1$  then there is no notable improvement from using the PCA-extracted latent variable (and notice that since the variances of  $x_i$  and  $z_i$  are 1, this means that the covariance is equal to the correlation in these simulations). These simulation results suggest that so long as the correlation between  $x_i$  and  $z_i$  is not close to  $-1, 0$ , or  $1$ , there are noticeable performance gains from using PCA to extract the true covariate from a collection of observed variables that try to measure that true covariate.

However, the performance advantages that we see from using PCA could be driven by the benefit of having multiple measurements of our true covariate of interest, as opposed to any special advantages from PCA specifically. We test this question by comparing the estimated  $\beta_1^*$  in our PCA 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. The results from these regressions for different values of  $p$  is shown below:

As one can see from these results (and results for different values of  $\beta_1$ ,  $\beta_2$ , and  $\rho$  in Appendix 2), there does not seem to be a noticeable difference between these three regression methods (across any values of  $p$ ,  $\beta_1$ ,  $\beta_2$ , and  $\rho$ ). Thus, our simulations suggest that there are major benefits to having multiple measurements of a latent covariate of interest, but that using PCA, taking the average of these measurements, and including all measurements as separate covariates seem to give similar benefits to the performance of the regression.

	<i>Number of p</i>			
	5	10	20	50
	<i>Coefficient on Main Variable</i>			
PCA	1.105 (0.121)	1.066 (0.122)	1.033 (0.119)	1.022 (0.117)
All Measurements	1.100 (0.124)	1.061 (0.129)	1.025 (0.131)	1.010 (0.167)
Average of Measurements	1.100 (0.121)	1.060 (0.122)	1.026 (0.119)	1.015 (0.117)
	<i>Absolute Percentage Error</i>			
PCA	13.1% (9.3 ppts)	11.1% (8.3 ppts)	10.0% (7.3 ppts)	9.4% (7.3 ppts)
All Measurements	12.9% (9.3 ppts)	11.4% (8.5 ppts)	10.7% (7.9 ppts)	13.2% (10.2 ppts)
Average of Measurements	12.8% (9.2 ppts)	10.9% (8.2 ppts)	9.8% (7.2 ppts)	9.3% (7.2 ppts)
Observations	1,000	1,000	1,000	1,000

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

Explain economic importance/interesting-ness of the chosen application

Explain how GDP/economic development is measured with error

It is very difficult to find an instrumental variable for economic development which satisfies a reasonable exclusion restriction.

In the left column in the table below I first regress the life expectancy at birth for all individuals in a given country and year on a measure of government spending as a share of total health expenditure. In the middle column I include the economic controls/covariates of GDP per capita (PPP), GNI per capita (PPP), Survey Mean Income/Consumption Per Capita, ILO GDP per person employed, and Net Foreign Assets Per Capita, all from the World Bank. In the rightmost column I instead use the first principal component combining these covariates.

I standardize all variables by subtracting the mean and dividing by the standard deviation, linearly interpolate data between known observations, and remove country-years with missing values for any of the economic indicators.

<i>Life Expectancy at Birth (Years)</i>					
	(1)	(2)	(3)	(4)	(5)
Govt. Share of Health Exp.	0.613*** (0.019)	0.308*** (0.020)	-0.016 (0.019)	0.343*** (0.019)	-0.012 (0.018)
Covariates	None	Econ Indicators	Econ Indicators	PCs	PCs
Fixed Effects	No	No	Yes	No	Yes
Observations	1,799	1,799	1,799	1,799	1,799
$R^2$	0.376	0.586	0.987	0.536	0.987
Adjusted $R^2$	0.375	0.582	0.985	0.535	0.985
Residual Std. Error	0.791	0.646	0.122	0.682	0.123
F Statistic	1081.530***	168.157***	922.656***	1036.417***	1114.735***

*Note:*

\*p<0.1; \*\*p<0.05; \*\*\*p<0.01  
All variables are standardized



## With Ginis instead of econ controls

	<i>Life Expectancy at Birth (Years)</i>				
	(1)	(2)	(3)	(4)	(5)
Govt. Share of Health Exp.	0.697*** (0.034)	0.652*** (0.040)	0.028 (0.024)	0.692*** (0.037)	0.026 (0.023)
Covariates	None	Ginis	Ginis	PCs	PCs
Fixed Effects	No	No	Yes	No	Yes
Observations	322	322	322	322	322
$R^2$	0.566	0.596	0.999	0.566	0.999
Adjusted $R^2$	0.565	0.583	0.998	0.564	0.998
Residual Std. Error	0.596	0.583	0.040	0.597	0.040
F Statistic	417.382***	45.905***	1421.063***	208.238***	5197.564***

*Note:*

\*p<0.1; \*\*p<0.05; \*\*\*p<0.01  
All variables are standardized.

## Health share of gdp

	<i>Life Expectancy at Birth (Years)</i>				
	(1)	(2)	(3)	(4)	(5)
GDP Share of Health Exp.	0.258*** (0.023)	0.184*** (0.019)	-0.024 (0.031)	0.144*** (0.017)	-0.020 (0.030)
Covariates	None	Econ Indicators	Econ Indicators	PCs	PCs
Fixed Effects	No	No	Yes	No	Yes
Observations	1,799	1,799	1,799	1,799	1,799
$R^2$	0.067	0.553	0.987	0.467	0.987
Adjusted $R^2$	0.066	0.549	0.985	0.467	0.985
Residual Std. Error	0.967	0.672	0.122	0.730	0.123
F Statistic	128.396***	146.946***	772.041***	787.819***	384.737***

*Note:*

\*p<0.1; \*\*p<0.05; \*\*\*p<0.01  
All variables are standardized.

## Conclusion

## Appendix 1

	<i>True <math>\beta_1</math></i>			
	0.1	1	10	100
<i>Coefficient on Main Variable</i>				
PCA	0.207 (0.121)	1.105 (0.121)	10.104 (0.123)	100.117 (0.124)
Single Measurement	0.383 (0.128)	1.280 (0.129)	10.278 (0.131)	100.289 (0.133)
<i>Absolute Percentage Error</i>				
PCA	131.1% (95.1 ppts)	13.1% (9.3 ppts)	1.3% (0.9 ppts)	0.1% (0.1 ppts)
Single Measurement	283.6% (126.6 ppts)	28.2% (12.6 ppts)	2.8% (1.3 ppts)	0.3% (0.1 ppts)
Observations	1,000	1,000	1,000	1,000

	<i>True <math>\beta_2</math></i>			
	0.1	1	10	100
<i>Coefficient on Main Variable</i>				
PCA	1.018 (0.115)	1.105 (0.121)	2.112 (0.477)	12.171 (4.555)
Single Measurement	1.034 (0.107)	1.280 (0.129)	3.865 (0.703)	29.751 (7.231)
<i>Absolute Percentage Error</i>				
PCA	9.4% (7.0 ppts)	13.1% (9.3 ppts)	111.6% (47.0 ppts)	1,119.6% (449.4 ppts)
Single Measurement	8.9% (6.8 ppts)	28.2% (12.6 ppts)	286.5% (70.3 ppts)	2,875.1% (723.1 ppts)
Observations	1,000	1,000	1,000	1,000

## Appendix 2

## Appendix 3

	<i>Number of p</i>			
	5	10	20	50
	<i>Coefficient on Main Variable</i>			
PCA	1.105 (0.121)	1.066 (0.122)	1.033 (0.119)	1.022 (0.117)
Single Measurement	1.280 (0.124)	1.283 (0.129)	1.282 (0.131)	1.292 (0.167)
	<i>Absolute Percentage Error</i>			
PCA	13.1% (9.3 ppts)	11.1% (8.3 ppts)	10.0% (7.3 ppts)	9.4% (7.3 ppts)
Single Measurement	28.2% (12.6 ppts)	28.5% (12.6 ppts)	28.3% (12.6 ppts)	29.3% (12.7 ppts)
Observations	1,000	1,000	1,000	1,000

	<i><math>\rho</math> Value</i>				
	-1	-0.5	0	0.5	1
	<i>Coefficient on Main Variable</i>				
PCA	-0.006 (0.238)	0.900 (0.120)	0.996 (0.111)	1.105 (0.121)	2.009 (0.242)
Single Measurement	-0.002 (0.142)	0.720 (0.130)	0.998 (0.127)	1.280 (0.129)	2.003 (0.147)
	<i>Absolute Percentage Error</i>				
PCA	100.6% (23.8 ppts)	12.7% (9.1 ppts)	8.9% (6.6 ppts)	13.1% (9.3 ppts)	100.9% (24.2 ppts)
Single Measurement	100.2% (14.2 ppts)	28.1% (12.7 ppts)	10.2% (7.6 ppts)	28.2% (12.6 ppts)	100.3% (14.7 ppts)
Observations	1,000	1,000	1,000	1,000	

	<i>True <math>\beta_1</math></i>			
	0.1	1	10	100
	<i>Coefficient on Main Variable</i>			
PCA	0.207 (0.121)	1.105 (0.121)	10.104 (0.123)	100.117 (0.124)
All Measurements	0.201 (0.123)	1.100 (0.124)	10.098 (0.126)	100.11 (0.127)
Average of Measurements	0.202 (0.121)	1.100 (0.121)	10.098 (0.123)	100.111 (0.124)
	<i>Absolute Percentage Error</i>			
PCA	131.1% (95.1 ppts)	13.1% (9.3 ppts)	1.3% (0.9 ppts)	0.1% (0.1 ppts)
All Measurements	128.9% (93.6 ppts)	12.9% (9.3 ppts)	1.3% (1.0 ppts)	0.1% (0.1 ppts)
Average of Measurements	127.9% (93.1 ppts)	12.8% (9.2 ppts)	1.3% (0.9 ppts)	0.1% (0.1 ppts)
Observations	1,000	1,000	1,000	1,000

	<i>True <math>\beta_2</math></i>			
	0.1	1	10	100
	<i>Coefficient on Main Variable</i>			
PCA	1.018 (0.115)	1.105 (0.121)	2.112 (0.477)	12.171 (4.555)
All Measurements	1.02 (0.118)	1.100 (0.124)	2.067 (0.477)	11.664 (4.519)
Average of Measurements	1.017 (0.115)	1.100 (0.121)	2.061 (0.470)	11.625 (4.415)
	<i>Absolute Percentage Error</i>			
PCA	9.4% (7.0 ppts)	13.1% (9.3 ppts)	111.6% (47.0 ppts)	1,119.6% (449.4 ppts)
All Measurements	9.7% (7.0 ppts)	12.9% (9.3 ppts)	107.1% (46.8 ppts)	1,069.3% (445.0 ppts)
Average of Measurements	9.4% (6.9 ppts)	12.8% (9.2 ppts)	106.5% (46.0 ppts)	1,065.2% (435.0 ppts)
Observations	1,000	1,000	1,000	1,000

	<i>Number of p</i>			
	5	10	20	50
<i>Coefficient on Main Variable</i>				
PCA	1.105 (0.121)	1.066 (0.122)	1.033 (0.119)	1.022 (0.117)
All Measurements	1.100 (0.124)	1.061 (0.129)	1.025 (0.131)	1.010 (0.167)
Average of Measurements	1.100 (0.121)	1.060 (0.122)	1.026 (0.119)	1.015 (0.117)
<i>Absolute Percentage Error</i>				
PCA	13.1% (9.3 ppts)	11.1% (8.3 ppts)	10.0% (7.3 ppts)	9.4% (7.3 ppts)
All Measurements	12.9% (9.3 ppts)	11.4% (8.5 ppts)	10.7% (7.9 ppts)	13.2% (10.2 ppts)
Average of Measurements	12.8% (9.2 ppts)	10.9% (8.2 ppts)	9.8% (7.2 ppts)	9.3% (7.2 ppts)
Observations	1,000	1,000	1,000	1,000

	<i><math>\rho</math> Value</i>				
	-1	-0.5	0	0.5	1
<i>Coefficient on Main Variable</i>					
PCA	-0.006 (0.238)	0.900 (0.120)	0.996 (0.111)	1.105 (0.121)	2.009 (0.242)
All Measurements	-0.007 (0.249)	0.904 (0.122)	0.996 (0.112)	1.100 (0.124)	2.011 (0.249)
Average of Measurements	-0.005 (0.243)	0.905 (0.120)	0.996 (0.110)	1.100 (0.121)	2.010 (0.246)
<i>Absolute Percentage Error</i>					
PCA	100.6% (23.8 ppts)	12.7% (9.1 ppts)	8.9% (6.6 ppts)	13.1% (9.3 ppts)	100.9% (24.2 ppts)
All Measurements	100.7% (24.9 ppts)	12.6% (9.0 ppts)	9.0% (6.7 ppts)	12.9% (9.3 ppts)	101.1% (24.9 ppts)
Average of Measurements	100.5% (24.3 ppts)	12.4% (8.9 ppts)	8.9% (6.6 ppts)	12.8% (9.2 ppts)	101.0% (24.6 ppts)
Observations	1,000	1,000	1,000	1,000	

Figure 1: Correlations Between Covariates and Life Expectancy

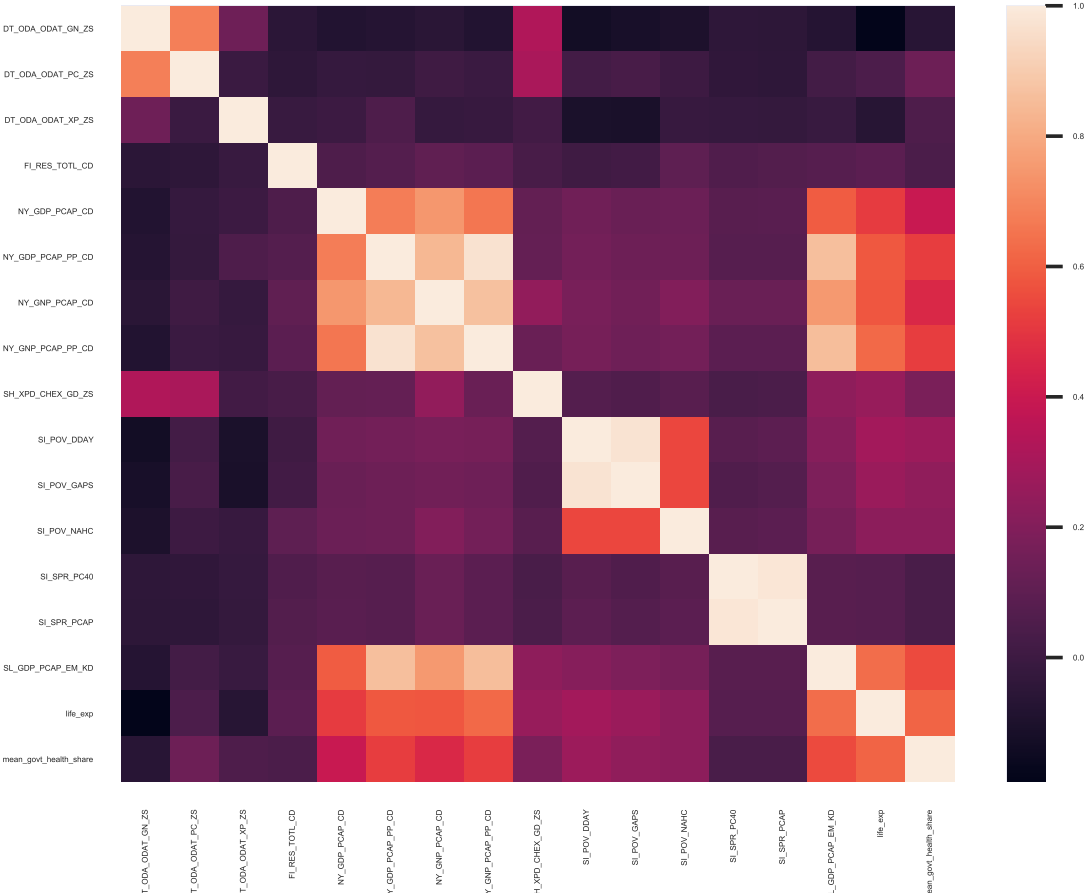


Figure 2: Economic Measures PCA Loadings

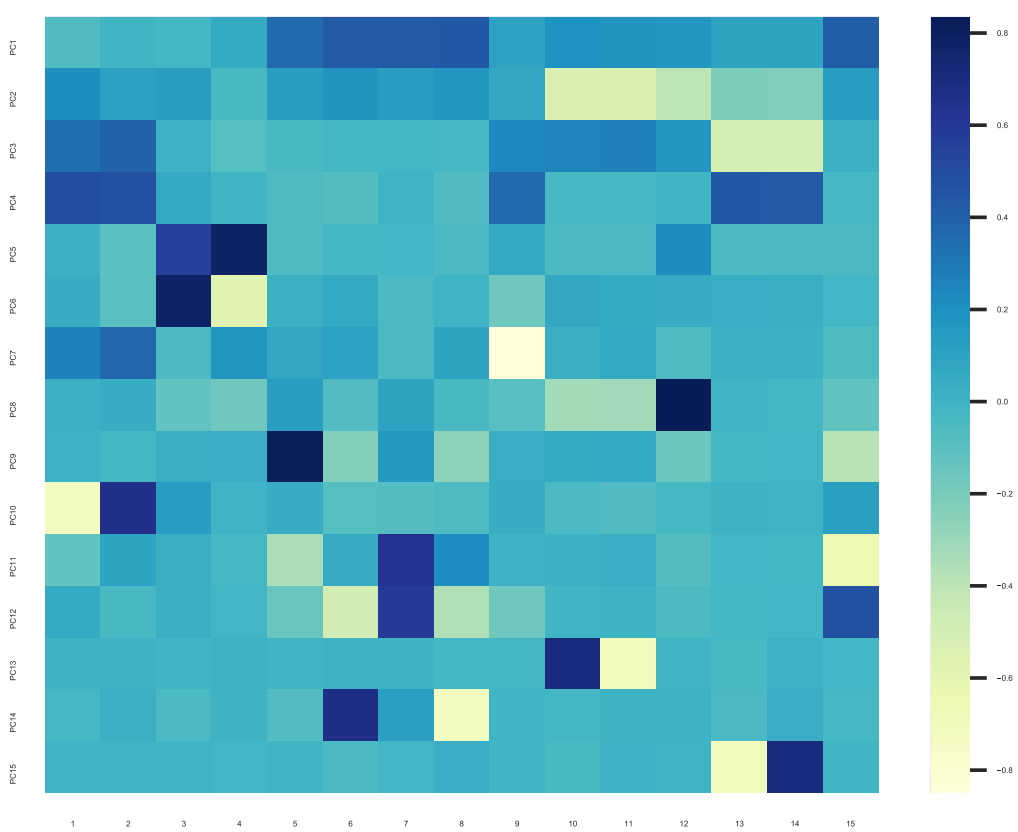




Figure 3: Economic Measures PCA Share of Variance Explained

