

Attenuation Bias, Measurement Error & Principal Component Analysis

Isaac Liu, Nicolás Martorell & Paul Opheim

May 25, 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.

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 our setting, (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) 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, (Schoefield) combines solutions from structural equations modelling and item response theory to deal with misestimation. Finally, (Heckman 2010) considers a situation similar to ours, except involving matching estimators. In this case, these 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, in order 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. Our estimator generally behaves as expected in this empirical setting, though it is unclear whether it performs any better or worse than other methods such as the direct inclusion of covariates or their averaging.

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 γ 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

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

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

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

$$y_i = \gamma^* t_i + \beta^* x_i + \zeta_i \quad (6)$$

Then, equations (4) and (5) will be such that

$$\begin{pmatrix} \gamma^* \\ \beta^* \end{pmatrix} = \begin{pmatrix} \sigma_t^2 & \text{cov}(t, x^*) \\ \text{cov}(x^*, t) & \sigma_{x^*}^2 \end{pmatrix}^{-1} \begin{pmatrix} \text{cov}(y, t) \\ \text{cov}(y, x^*) \end{pmatrix} \quad (7)$$

$$(8)$$

iiiiiii HEAD

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

$$y_i = \gamma^* t_i + \beta^* x_i + \zeta_i \quad (9)$$

Then, equations (4) and (5) will be such that

$$\begin{pmatrix} \gamma^* \\ \beta^* \end{pmatrix} = \begin{pmatrix} \sigma_t^2 & \text{cov}(t, x^*) \\ \text{cov}(x^*, t) & \sigma_{x^*}^2 \end{pmatrix}^{-1} \begin{pmatrix} \text{cov}(y, t) \\ \text{cov}(y, x^*) \end{pmatrix} \quad (10)$$

$$(11)$$

===== 5a34fa94c20f60fb8872147c40d789ba110ff14c

Thus, the coefficient of our variable of interest will be

$$\gamma^* = \frac{\sigma_{x^*}^2 \text{cov}(y, t) - \text{cov}(t, x^*) \text{cov}(y, x^*)}{\sigma_t^2 \sigma_{x^*}^2 - \text{cov}(t, x^*)^2} \quad (12)$$

And then

$$\gamma^* = \gamma + \beta \frac{\text{cov}(t, x)(\sigma_{x^*}^2 - \sigma_x^2)}{\sigma_t^2 \sigma_{x^*}^2 - \text{cov}(t, x)^2} \quad (13)$$

||||| HEAD From (10) it is clear that when $\text{cov}(t, x) \neq 0$ and that x is measure 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 be facing any kind of bias, as one would expect.

Equation (10) 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 *Cauchy-Schwarz* inequality that the denominator is also positive. Then, the direction of the bias will depend on the sign of β and the covariance of t and x , as Table one illustrates. ===== From (10) 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 be facing any kind of bias, as one would expect.

Equation (10) 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 β and the covariance of t and x , as Table one illustrates. ||||||| 5a34fa94c20f60fb8872147c40d789ba110ff14c

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

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

|||||

HEAD

For a more concrete example, consider a simple case where the DGP depends only of the variable of interest and a covariate such that:

$$\begin{pmatrix} \gamma \\ \beta \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad (14)$$

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

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

$$\begin{pmatrix} \gamma^* \\ \beta^* \end{pmatrix} = \begin{pmatrix} 1.37 \\ 0.39 \end{pmatrix}$$

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.

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:

$$h_j = \underset{h'h=1, h'h_1=0, \dots, h'h_{j-1}=0}{\operatorname{argmax}} \operatorname{var} [h'X_k^*] \quad (15)$$

where h_j is the eigenvector of Σ associated with the j^{th} ordered eigenvalue λ_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 Σ associated with the j^{th} ordered eigenvalue λ_j of Σ .

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

$$x_{i,k} = h_k'X_k^* \quad (16)$$

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

$$X_i = HX_i^* \quad (17)$$

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

$$y_i = \gamma^{PCR} t_i + H X_i^{*'} \beta^{PCR} + \epsilon_i \quad (18)$$

where the coefficients are as follows

$$\begin{pmatrix} \gamma^{PCR} \\ \beta^{PCR} \end{pmatrix} = \begin{pmatrix} \sigma_t^2 & \Sigma_{t, HX^*} \\ \Sigma_{HX^*, t} & \Sigma_{HX^*} \end{pmatrix}^{-1} \begin{pmatrix} \Sigma_{yt} \\ \Sigma_{y, HX^*} \end{pmatrix} \quad (19)$$

$$= \begin{pmatrix} \sigma_t^2 & \Sigma_{t, HX^*} \\ \Sigma_{HX^*, t} & \Sigma_{HX^*} \end{pmatrix}^{-1} \begin{pmatrix} \sigma_t^2 & \Sigma_{tX} \\ \Sigma_{Xt} & \Sigma_X \end{pmatrix} \begin{pmatrix} \gamma \\ \beta \end{pmatrix} \quad (20)$$

$$= \begin{pmatrix} \gamma \\ \beta \end{pmatrix} \quad (21)$$

where the last equality comes from (13). ===== For a more concrete example, consider a simple case where the DGP depends only of the variable of interest and a covariate such that:

$$\begin{pmatrix} \gamma \\ \beta \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad (22)$$

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 both variable have positive correlation.

Principal Component Regression as bias-correction method

The classical solution for the measurement-error induced bias in econometrics has been the usage of instrumental variables. However, finding reliable source of exogeneity is difficult, as is proving the exclusion condition. Therefore, the use of IV as a bias-correction method should be taken with care given that its feasibility is hard.

Alternatively, we propose an alternative bias-correction method when there are several missmeasured 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 miss-measured 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

$$h_j = \underset{h'h=1, h'h_1=0, \dots, h'h_{j-1}=0}{\operatorname{argmax}} \operatorname{var}[h'X_k^*] \quad (23)$$

where h_j is the eigenvector of Σ associated with the j^{th} ordered eigenvalue λ_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 Σ associated with the j^{th} ordered eigenvalue λ_j of Σ .

Under our assumptions, the vector of missmeasured 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

$$x_{i,k} = h_k' X_k^* \quad (24)$$

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

$$X_i = H X_i^* \quad (25)$$

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 would be

$$y_i = \gamma^{PCR} t_i + H X_i^{*'} \beta^{PCR} + \epsilon_i \quad (26)$$

where the coefficients are as follows

$$\begin{pmatrix} \gamma^{PCR} \\ \beta^{PCR} \end{pmatrix} = \begin{pmatrix} \sigma_t^2 & \Sigma_{t, HX^*} \\ \Sigma_{HX^*, t} & \Sigma_{HX^*} \end{pmatrix}^{-1} \begin{pmatrix} \Sigma_{yt} \\ \Sigma_{y, HX^*} \end{pmatrix} \quad (27)$$

$$= \begin{pmatrix} \sigma_t^2 & \Sigma_{t, HX^*} \\ \Sigma_{HX^*, t} & \Sigma_{HX^*} \end{pmatrix}^{-1} \begin{pmatrix} \sigma_t^2 & \Sigma_{tX} \\ \Sigma_{Xt} & \Sigma_X \end{pmatrix} \begin{pmatrix} \gamma \\ \beta \end{pmatrix} \quad (28)$$

$$= \begin{pmatrix} \gamma \\ \beta \end{pmatrix} \quad (29)$$

where the last equality comes from (13).

Note that according to equation (17), the true variable $x_{i,k}$ is a linear combination of the missmeasured variables that the researcher may have, were the weights are such that equation (16) is satisfied. This allows us to think about other linear combination that could be used as a bias-correction method. 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 amount of missmeasured variables for $x_{i,k}$. Then, equation (17) will be

$$x_{i,k} = \begin{pmatrix} \frac{1}{n} & \frac{1}{n} & \dots & \frac{1}{n} \end{pmatrix} \begin{pmatrix} x_{i,1}^* \\ x_{i,2}^* \\ \vdots \\ x_{i,J}^* \end{pmatrix} \quad (30)$$

$$= \frac{1}{n} \sum_{j=1}^J x_{i,j}^* \quad (31)$$

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

Estimation of the Principal Component Regression

Recall that the missmeasured variables are such that $x_{i,k}^* = x_{i,k} + \eta_{i,k}$. Then for every i , X_i^* can we get

$$X_i^* = X_i + \mu \quad (32)$$

and so we could interpret the vector of missmeasured variables X_i^* as a factor model, in which the common factor is the vector of true variables X_i . This way, we could rewrite the model as follows

$$y_i = \gamma t_i + X_i' \beta + \epsilon_i \quad (33)$$

$$X_i^* = X_i + \mu \quad (34)$$

Which is a factor-augmented regression model in which the common “factors” between the missmeasured variables of $x_{i,k}$ is in fact $x_{i,k}$. We estimate the model following two-stages. First, we estimate X_i by factor regression. Then, the first-stage is the principal-components estimation

$$\hat{X}_i = \hat{D}^{-1} X_i^* \quad (35)$$

Second step is to regress Y_i on the estimated \hat{X}_i . Then, the coefficients would be

$$\begin{pmatrix} \hat{\gamma}^F \\ \hat{\beta}^F \end{pmatrix} = \begin{pmatrix} \hat{\sigma}_t^2 & \hat{\Sigma}_{t,\hat{X}} \\ \hat{\Sigma}_{\hat{X},t} & \hat{\Sigma}_{\hat{X}} \end{pmatrix}^{-1} \begin{pmatrix} \hat{\Sigma}_{yt} \\ \hat{\Sigma}_{y,\hat{X}} \end{pmatrix} \quad (36)$$

Using the *Frisch-Waugh-Lowell* decomposition

$$\hat{\gamma}^F = ((M_{\hat{X}} t)' M_{\hat{X}} t)^{-1} (M_{\hat{X}} t)' M_{\hat{X}} y \quad (37)$$

Where $M_{\hat{X}} = I - \hat{X}(\hat{X}' \hat{X})^{-1} \hat{X}'$. Then

$$\hat{\gamma}^F = ((M_{\hat{X}} t)' M_{\hat{X}} t)^{-1} (M_{\hat{X}} t)' M_{\hat{X}} (t\gamma + X\beta + \epsilon) \quad (38)$$

$$= \gamma + ((M_{\hat{X}} t)' M_{\hat{X}} t)^{-1} (M_{\hat{X}} t)' M_{\hat{X}} X\beta + ((M_{\hat{X}} t)' M_{\hat{X}} t)^{-1} (M_{\hat{X}} t)' M_{\hat{X}} \epsilon \quad (39)$$

5a34fa94c20f60fb8872147c40d789ba110ff14c

Estimation of the Principal Component Regression

Recall that the missmeasured variables are such that $x_{i,k}^* = x_{i,k} + \eta_{i,k}$. Then for every i , X_i^* can we get

$$X_i^* = X_i + \mu \quad (40)$$

and so we could interpret the vector of missmeasured variables X_i^* as a factor model, in which the common factor is the vector of true variables X_i . This way, we could rewrite the model as follows

$$y_i = \gamma t_i + X_i' \beta + \epsilon_i \quad (41)$$

$$X_i^* = X_i + \mu \quad (42)$$

Which is a factor-augmented regression model in which the common “factors” between the missmeasured variables of $x_{i,k}$ is in fact $x_{i,k}$. We estimate the model following two-stages. First, we estimate X_i by factor regression. Then, the first-stage is the principal-components estimation

$$\hat{X}_i = \hat{D}^{-1} X_i^* \quad (43)$$

Second step is to regress Y_i on the estimated \hat{X}_i . Then, the coefficients would be

$$\begin{pmatrix} \hat{\gamma}^F \\ \hat{\beta}^F \end{pmatrix} = \begin{pmatrix} \hat{\sigma}_t^2 & \hat{\Sigma}_{t,\hat{X}} \\ \hat{\Sigma}_{\hat{X},t} & \hat{\Sigma}_{\hat{X}} \end{pmatrix}^{-1} \begin{pmatrix} \hat{\Sigma}_{yt} \\ \hat{\Sigma}_{y,\hat{X}} \end{pmatrix} \quad (44)$$

Using the *Frisch-Waugh-Lowell* decomposition

$$\hat{\gamma}^F = ((M_{\hat{X}} t)' M_{\hat{X}} t)^{-1} (M_{\hat{X}} t)' M_{\hat{X}} y \quad (45)$$

Where $M_{\hat{X}} = I - \hat{X}(\hat{X}' \hat{X})^{-1} \hat{X}'$. Then

$$\hat{\gamma}^F = ((M_{\hat{X}} t)' M_{\hat{X}} t)^{-1} (M_{\hat{X}} t)' M_{\hat{X}} (t\gamma + X\beta + \epsilon) \quad (46)$$

$$= \gamma + ((M_{\hat{X}} t)' M_{\hat{X}} t)^{-1} (M_{\hat{X}} t)' M_{\hat{X}} X \beta + ((M_{\hat{X}} t)' M_{\hat{X}} t)^{-1} (M_{\hat{X}} t)' M_{\hat{X}} \epsilon \quad (47)$$

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 ρ 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 η_j is drawn from a white noise distribution $\mathcal{N}(\mathbf{0}, \Sigma)$ where $\mathbf{0}$ is a p -vector and Σ 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. Table 2 shows the results for different values of p .

Table 2: Average Coefficients for 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 β_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 β_1^* coefficient in the PCA regression gets steadily closer to the true β_1 value of 1. This trend could well continue as $p \rightarrow \infty$, but we did not simulate values greater than $p = 50$ due to the seeming implausibility of having more than 50 measurements of the same single covariate. Appendix 1 contains charts that show that this increase in performance is also true for different values of β_1 and β_2 .

However, there are certain circumstances where the PCA method does not lead to more accurate estimates of β_1^* . Table 3 contains the simulation results for different values of ρ (the covariance between the main variable of interest x_i and the true latent covariate z_i):

Table 3: Average Coefficients for Values of ρ

| | ρ 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 | |

We first note that when $\rho = 0.5$ then the coefficient on the variable of interest is artificially inflated when we use a single mismeasurement as a covariate (on average, 1.28 instead of the true value of 1). Conversely, when $\rho = -0.5$ 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 very little bias when $\rho = 0$, as predicted in the theory section. Since there is hardly any bias to correct, we do not see gains from using the PCA covariate method for that particular ρ value. Also, when $\rho = -1$ or 1, 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 β_1^* in our PCA regressions with the estimated β_1^* when we include all p measurements as separate covariates in the regression, and the β_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 are shown in Table 4.

Table 4: Average Coefficients for 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) |
| 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 |

As one can see from these results (and results for different values of β_1 , β_2 , and ρ in Appendix 2), there does not seem to be a noticeable difference between these three regression methods (across any values of p , β_1 , β_2 , and ρ). 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.

Application: Government Share of Healthcare Spending and Life Expectancy

We now examine the implications 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 effect of publicly and privately funded healthcare systems on 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) 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. Other studies find a variety of effects; considering developing nations in Africa, (Novignon) finds a larger impact for public spending while (Fillmer) finds no impact; in the

developed US and Canada, (Lichtenberg) finds an impact only for public spending on outcomes when controlling for lagged GDP, while (Cremieux) finds private spending on drugs is more effective. In the most similar work 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 important 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) finds that economic growth is associated with increased life expectancy in Malaysia, while considering the reversed causal direction (Acemoglu) 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).

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). Other sources of error include the presence of the informal economy and non-monetary but productive work, and the challenge of accurately measuring the value of digital services which often do not have visible prices (Charmes, Ahmad).¹ 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. 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, though we also make usage of OECD government health share data to fill in missing years, and average World Bank and OECD measurements when both are available (World Bank, OECD). We 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.

In Table 5 below we begin with a univariate OLS regression, which produces a large and significant coefficient indicating a one standard deviation increase in the government share of health expenditure is linked to a 0.56 standard deviation increase in life expectancy. Next we include the potentially mismeasured covariate of GDP per capita, which greatly reduces the size of the coefficient on the governments share. After this we include a full set of economic covariates listed in Table 6, again reducing the size of the coefficient. However, this result is then adjusted upwards again in the last two columns, where we use the mean of the mismeasured covariates, and the first principal component.

¹Due 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.

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. | 0.564*** (0.018) | 0.365*** (0.019) | 0.256*** (0.018) | 0.380*** (0.018) | 0.298*** (0.018) |
| Covariates | None | GDP PC | Econ Indicators | Mean | PC 1 |
| Observations | 1,995 | 1,995 | 1,995 | 1,995 | 1,995 |
| R^2 | 0.319 | 0.456 | 0.573 | 0.479 | 0.518 |
| Adjusted R^2 | 0.318 | 0.455 | 0.569 | 0.479 | 0.517 |
| Residual Std. Error | 0.826 | 0.738 | 0.656 | 0.722 | 0.695 |
| F Statistic | 931.676*** | 833.433*** | 176.737*** | 917.087*** | 1069.682*** |

Note:

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

Table 6: Econ Indicators

| Indicator Name |
|--|
| GDP Per Capita (Current USD) |
| GNP Per Capita PPP (Current International \$) |
| GNP Per Capita (Current USD) |
| Survey Mean Income or Consumption Per Capita |
| Survey Mean Income or Consumption Per Capita, Bottom 40% |
| ILO GDP Per Person Employed |
| Total Foreign Reserves |
| Poverty Headcount Ratio (International Extreme Poverty Line) |
| Poverty Headcount Ratio (National Poverty Line) |
| Poverty Gap (International Extreme Poverty Line) |
| Government Share of Health Expenditure |
| Net Official Development Assistance Per Capita (Current USD) |
| Net Official Development Assistance Per Capita, Percent of GNI |
| Net Official Development Assistance, Percent of Total Government Expenditure |

The results clearly demonstrate that all of the methods using multiple measures of the covariates produce noticeably different coefficients for government health share. Taking the mean of covariates produces a larger coefficient. The principal components estimator and the usage of many economic indicators directly instead produce smaller coefficients.

Moreover, these different coefficients behave in a manner similar to that predicted by our theoretical development and simulations. In Table 3, we saw the impact of variation in ρ , the correlation between measurements for a value of $p = 5$ and $\beta = 1$ for 100 observations. In the empirical setting it is difficult to tell what is a reasonable value of β , and we instead have nearly 2000 observations. Nevertheless, we see that for a positive ρ 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 6.

Results using fixed effects panel models (with country clustered standard errors), more principal components, and instrumental variables regression of using all the economic indicators as an instrument for GDP per capita, are in Table 7. Fixed effects coefficients greatly reduce the magnitude of effects. Including more principal components and the instrumental variables technique produce coefficients of moderate size close to that obtained by the single principal component. The results in column (3) and (4) also show reductions in the inflation of coefficients relative to column (2) of Table 6.

Table 7: Additional Regressions

| | <i>Life Expectancy at Birth (Years)</i> | | | |
|----------------------------|---|---------------------|---------------------|---------------------|
| | (1) | (2) | (3) | (4) |
| Govt. Share of Health Exp. | 0.024*** (0.008) | 0.026*** (0.007) | 0.307*** (0.018) | 0.333*** (0.018) |
| Covariates | None | PC 1 | PC 1-7 | GDP PC (IV) |
| Fixed Effects | Yes | Yes | No | No |
| Observations | 1,995 | 1,995 | 1,995 | 1,995 |
| R^2 | 0.987 | 0.987 | 0.529 | 0.490 |
| Adjusted R^2 | 0.985 | 0.985 | 0.527 | 0.490 |
| Residual Std. Error | 0.121 | 0.121 | 0.688 | 0.714 |
| F Statistic | 2458.091*** | 2357.344*** | 279.197*** | 958.731*** |

Note:

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

Conclusion

Discuss possible extensions

Appendix 1

Table 8: Average Coefficients for Values of β_1

| | <i>True β_1</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 |

Table 9: Average Coefficients for Values of β_2

| | <i>True β_2</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 |

Table 10: Average Coefficients for Values of β_1

| | <i>True β_1</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 |

Table 11: Average Coefficients for Values of β_2

| | <i>True β_2</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 |

Table 12: Average Coefficients for Values of ρ

| | ρ 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) |
| 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 | |

Appendix 2

Appendix 3

Figure 1: Correlations Between Covariates and Life Expectancy

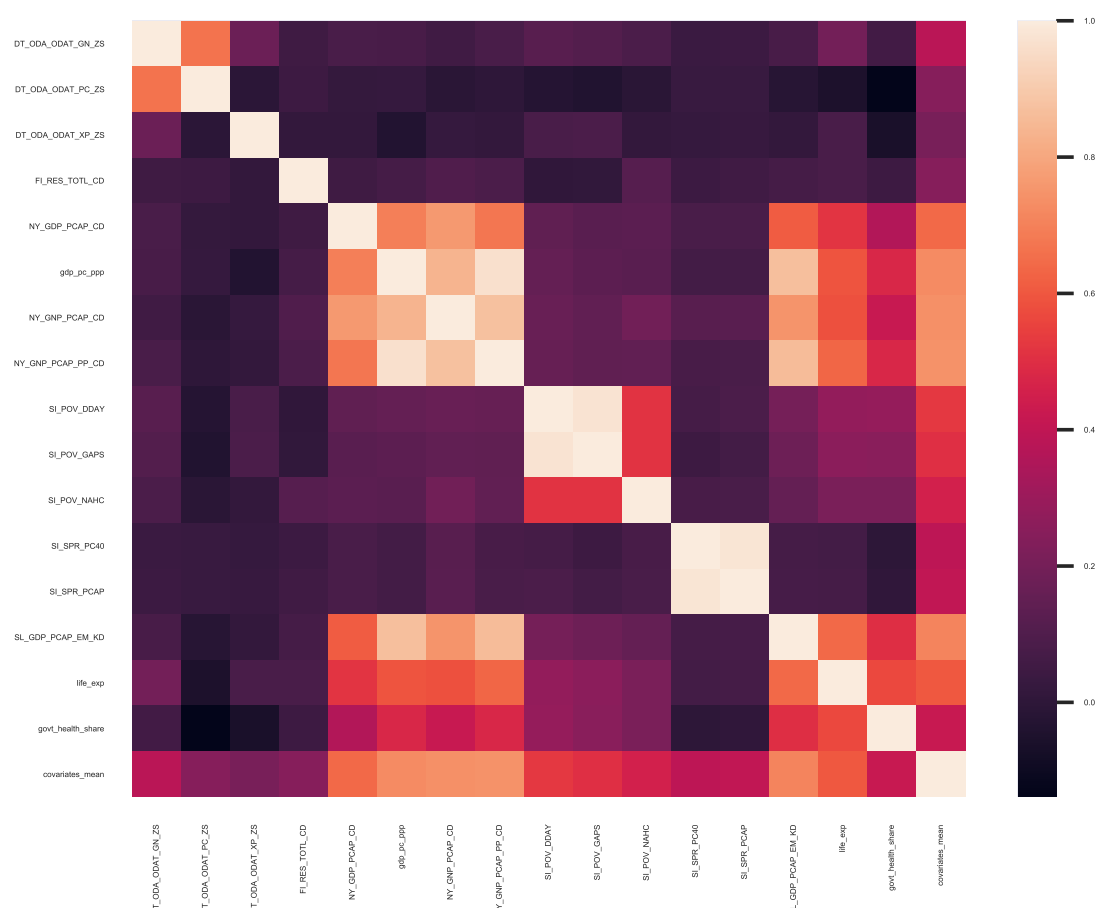


Figure 2: Economic Measures PCA Loadings

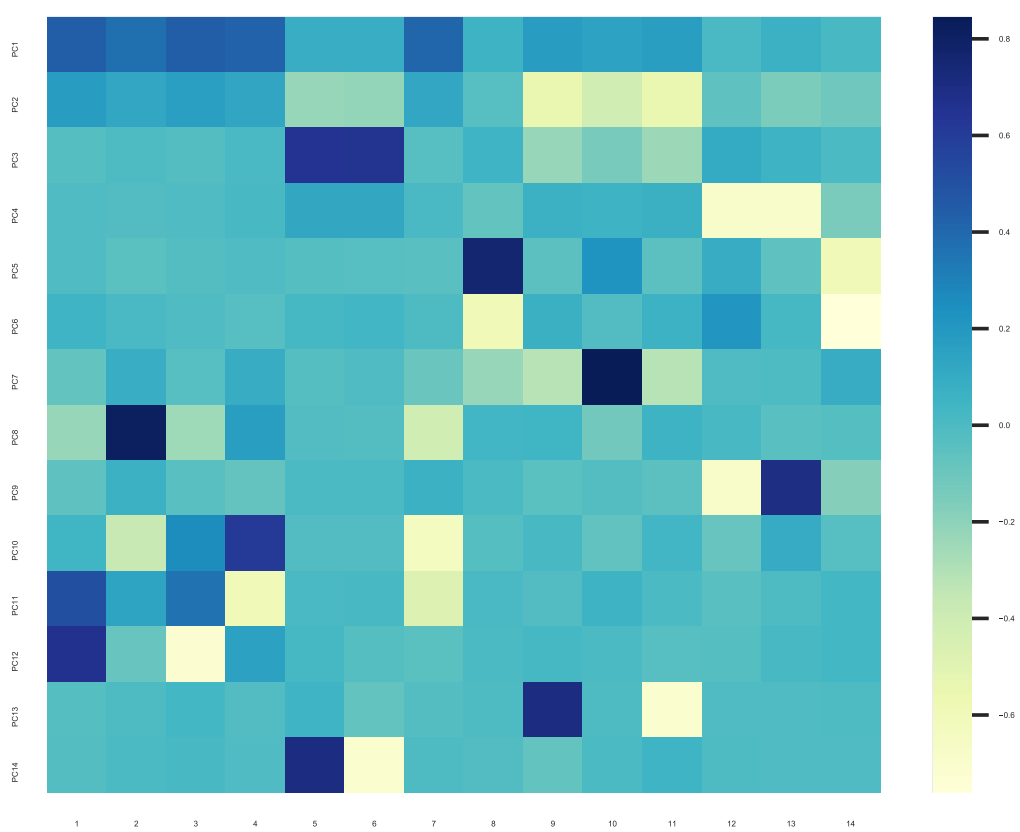


Figure 3: Economic Measures PCA Share of Variance Explained

