

Regression with an Imputed Dependent Variable

CEMMAP-UKHLS Missing Data Workshop

Thomas F. Crossley (Essex, IFS and ESCoE)

with Peter Levell (IFS and UCL) and Stavros Poupakis (UCL)

May 2019

Motivation

- Wish to estimate $y = X\beta + \epsilon$.
- β is the object of interest.
- OLS would be fine if we had complete data: $\text{plim}(X\epsilon) = 0$.
- But no data on $\frac{1}{N} \sum yX$.
- Have some data on $(y_1, Z_1), (X_2, Z_2)$.
 - these are both random samples from the population of interest
 - subscript indexes data set (or sample), absence of subscript means population.
- Z is a proxy or proxies for y .

Motivating Example

- Estimating effect of income or wealth (shocks) on consumption expenditure with:
 - a data set with food exp. and income/wealth.
 - a data set with food exp. and consumption exp.
 - e.g., PSID and CE, UKHLS and LCF, HFCS and National Budget Surveys



Set up

- data on (y_1, Z_1) and (X_2, Z_2) .

$$z = y\gamma + u.$$

$$z = X\beta\gamma + \epsilon\gamma + u.$$

- Z must depend on ϵ

$$\text{plim}\left(\frac{1}{n_j} X_j' u_j\right) = 0$$

$$\text{plim}\left(\frac{1}{n_j} y_j' u_j\right) = 0 \text{ (can relax)}$$

Alternative Imputation/Data Combination Strategies (1)

- Skinner (1987) suggested regressing y_1 on Z_1 in the CE and using the resulting coefficients to predict \hat{y}_2 in the PSID
- Then regressing \hat{y}_2 on X_2 .
- With a single spending category as the proxy, the first stage is an “inverse” Engel curve.
- **RP** procedure for “Regression Prediction”.
- Advocated by Browning, Crossley, Weber (2003).
- Employed by Attanasio and Pistaferri (2014), Arrondel et al., (2015).
- Add a first-stage residual: **RP+**.

Alternative Imputation/Data Combination Strategies (2)

- Blundell, Pistaferri, Preston (2004,2008) regress Z_1 on y_1 then predict $\hat{y}_2 = Z_2 \frac{1}{\hat{\gamma}}$: **BPP** procedure.
- Again using the CE and PSID, proxy, z , is food expenditure.
- Estimate an Engel curve and then invert it to predict consumption.
- Recently been employed Attanasio, Hurst and Pistaferri (2012).

Alternative Imputation/Data Combination Strategies (3)

- Do not impute y at the unit level at all.
- Recover β from a combination of moments taken from the two surveys.
- Arellano and Meghir (1992): **AM** procedure.
- Here:
 - Regress Z_1 on y_1 to get $\hat{\gamma}$.
 - Regress Z_2 on X_2 to get $\hat{\gamma}\beta$
 - Take ratio of the two to estimate β .

RP inconsistent for β

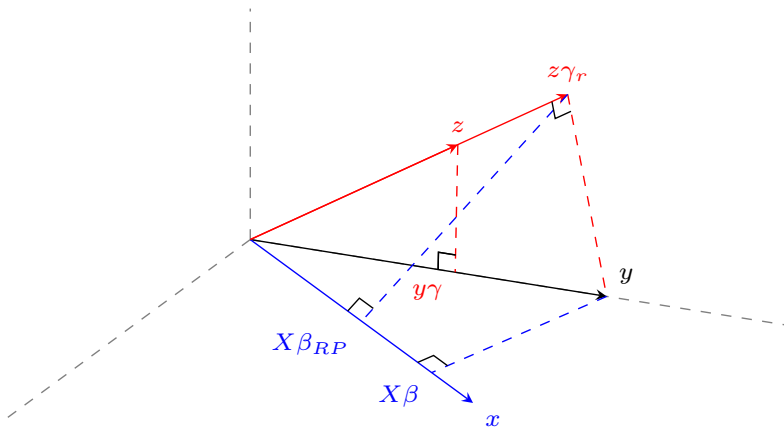
- **RP** does not consistently estimate β : $plim(\hat{\beta}^{RP}) = \beta R_{y_1, Z_1}^2$

Proof.

$$\begin{aligned}
 plim(\hat{\beta}^{RRP}) &= plim \left\{ \left(\frac{X_2' X_2}{n_2} \right)^{-1} \frac{X_2' Z_2}{n_2} \left(\frac{Z_1' Z_1}{n_1} \right)^{-1} \frac{Z_1' y_1}{n_1} \right\} \\
 &= plim \left\{ \left(\frac{X_2' X_2}{n_2} \right)^{-1} \frac{X_2' Z_2}{n_2} \left(\frac{Z_1' Z_1}{n_1} \right)^{-1} \frac{Z_1' y_1}{n_1} \frac{1}{R_{y_1, Z_1}^2} R_{y_1, Z_1}^2 \right\} \\
 &= \beta \gamma Q_{ZZ}^{-1} \gamma' \Sigma_{yy} \left[\Sigma_{yy} \gamma Q_{ZZ}^{-1} \gamma' \Sigma_{yy} \right]^{-1} \Sigma_{yy} \phi_{y, Z} \\
 &= \beta \gamma Q_{ZZ}^{-1} \gamma' \left[\gamma Q_{ZZ}^{-1} \gamma' \right]^{-1} = \beta \phi_{y, Z}
 \end{aligned}$$



Geometric Intuition



Magnitude

- First stage R^2 for food on total consumption 50 – 70%.
- Downward bias 30 – 50%.

- Similar problem with **RP+** (intuition below).

RRP

- As the bias in the RP procedure is an estimable quantity, it can be corrected.
- Rescale $\hat{\beta}^{RP}$ by the estimated first stage (centered) $R_{y,Z}^2$:
“Re-scaled Regression Prediction” (**RRP**, $\hat{\beta}^{RRP}$).

$$plim(\hat{\beta}^{RRP}) = plim\left(\frac{\hat{\beta}^{RP}}{R_{y_1, Z_1}^2}\right) = \beta$$

- Proof follows immediately from Proposition 1.
- Rescaling of $\hat{\beta}^{RP}$ is equivalent to rescaling the predicted consumption vector \hat{y}_2^{RP} by $1/R_{y,Z}^2$.

AM, BPP and RRP

- IFF there is a single proxy z , $\hat{\beta}^{RRP}$, $\hat{\beta}^{BPP}$ and $\hat{\beta}^{AM}$ are numerically identical.

Proof.

$$\hat{\beta}^{BPP} = (X_2'X_2)^{-1}X_2'z_2(y_1'z_1)^{-1}y_1'y_1 = \hat{\beta}^{RRP}$$

$$\begin{aligned}\hat{\beta}^{AM} &= \widehat{\beta\gamma} / \hat{\gamma} = (X_2'X_2)^{-1}X_2'z_2 \left[(y_1'y_1)^{-1}y_1'z_1 \right]^{-1} \\ &= (X_2'X_2)^{-1}X_2'z_2(y_1'z_1)^{-1}y_1'y_1 = \hat{\beta}^{RRP} = \hat{\beta}^{BPP}\end{aligned}$$

- Therefore **BPP**, **AM** consistently estimate β .

Other Moments: Variances and Covariances

$$\text{Asymp Var}(\hat{y}^{RP}) = \text{Asymp Var}(y) \times \phi_{y,Z}.$$

$$\text{Asymp Cov}(\hat{y}^{RP}, X) = \text{Asymp Cov}(y, X) \times \phi_{y,Z}.$$

$$\text{Asymp Var}(\hat{y}^{RRP}) = \text{Asymp Var}(\hat{y}^{BPP}) = \text{Asymp Var}(y) / \phi_{y,Z}.$$

$$\text{Asymp Cov}(\hat{y}^{RRP}, X) = \text{Asymp Cov}(\hat{y}^{BPP}, X) = \text{Asymp Cov}(y, X)$$

Other Moments: Means

- de-meaning not necessary
- $plim\left(\frac{\sum \hat{y}^{RP}}{n_2}\right) = \mu_y$,
- But then $plim\left(\frac{\sum \hat{y}^{RP} / R_{y,Z}^2}{n_2}\right) \neq \mu_y$.
- Statistical agency cannot release a single imputed \hat{y} .

Other Moments

Table: Summary of imputation methods (consistency)

	μ_y	σ_{yy}	β
Regression Prediction (RP)	✓	×	×
Regression Prediction + $\hat{\epsilon}$ (RP+)	✓	✓	×
Rescaled Regression Prediction (RRP)	×	×	✓
Blundell et al., 2004; 2008 (BPP)	✓	×	✓
Arellano and Meghir, 1992 (AM)	-	-	✓

Hot-Deck Imputation and Item Non-response

- y drawn from a matched cell is a regression prediction plus a residual (Lillard et al., 1983; David et al., 1986). .
 - Saturated regression on categorical variables
- If one or more matching variables are excluded from X this maps into **RP+**.
- Suppose fraction $\frac{n_a}{N}$ of cases missing.

$$plim(\hat{\beta}) = \frac{N - n_a}{N}\beta + \frac{n_a}{N}\beta\phi_{y,Z} = \beta \left(1 + \frac{n_a}{N}(\phi_{y,Z} - 1) \right).$$

Practicalities: Additional Covariates

- Residualize or add to both stages (Frisch-Waugh-Lovell).
- first-stage partial R^2 .

Practicalities: Measurement error in y

- Easy to see that with **AM** we need an instrument for y because we need consistent estimate of γ .
- True for **BPP** and **RRP** too.
- For **RRP**, sample R^2 will not be consistent estimate of population R^2 in presence of ME (but can estimate).

Practicalities: Panel Case

- Often wish to estimate $\Delta y = \Delta X\beta + \epsilon$
- OLS would be fine if we had complete data $E[X\epsilon] = 0$, $plim(X\epsilon) = 0$
- But no data on $\frac{1}{N} \sum \Delta y \Delta X$
- Have some data on $(y_1^1, Z_1), (y_2^0, Z_2), (\Delta X_3, Z_3)$ where $\Delta y = y^1 - y^0$
 - e.g. cross sectional budget survey and panel income/wealth survey

Practicalities: Panel Case

- **BPP** and **RRP** are identical (with one proxy) and consistent.

$$\hat{\beta} = (\Delta X' \Delta X)^{-1} \Delta X' \Delta \hat{y}$$

where $\Delta \hat{y} = \hat{y}^1 - \hat{y}^0$

Related Literature: Berkson Measurement Error

- classical measurement error on the right:

$$\tilde{X} = X + \tilde{v}, \tilde{v} \perp X, y = \tilde{X}\beta - \tilde{v}\beta + \epsilon$$

- classical measurement error on the left:

$$\tilde{y} = y + \tilde{v}, \tilde{v} \perp y, \tilde{y} = X\beta + \epsilon + \tilde{v}$$

- Berkson measurement error (prediction error) on the right:

$$X = \hat{X} + \hat{v}, \hat{v} \perp \hat{X}, y = \hat{X}\beta - \hat{v}\beta + \epsilon$$

- Here: prediction error (Berkson) on the left.

Related Literature: Berkson Measurement Error

- Hyslop and Imbens (2001) show attenuation bias in a regression of \hat{y} on X where \hat{y} is an optimal linear predictor (see also Hoderlien and Winter (2010))
- Key differences:
 - assume prediction by respondent and respondent knows Z , β and $E[X]$
 - also assume $Z = y + u$; ($\gamma = 1$)

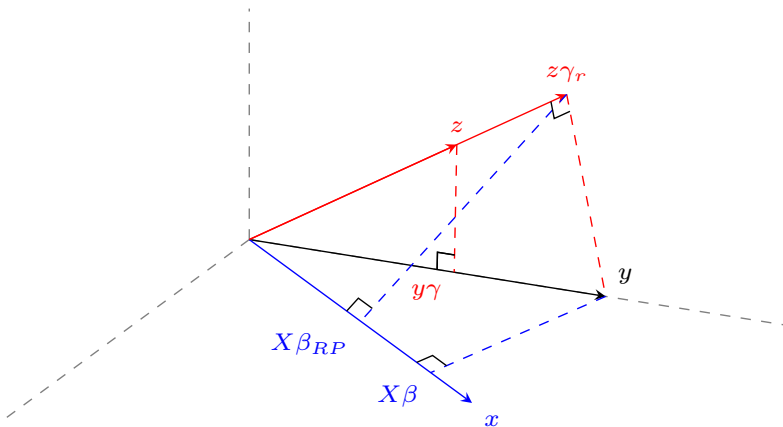
Related Literature: 2SIV

- Wish to estimate $y = X\beta + \epsilon$, have some data on $(y_1, Z_1), (X_2, Z_2)$.
- Z is a grouping variable (e.g.. birth cohort, occupation, birth cohort x education).
- In this case we effectively use **Z to impute X**.
- Two Sample IV (2SIV, Angrist & Krueger, 1992).

$$\hat{\beta}_{TSIV} = \left(\frac{Z_2'X_2}{n_2} \right)^{-1} \left(\frac{Z_1'y_1}{n_1} \right)$$

- key assumption: $Z \perp \epsilon$ (Z affects y *only* through X).
- Lusardi (1996): CE and PSID.

Geometric Intuition



Inference

- OLS with full data: asymptotic variance for $\hat{\beta}$ of $(\Sigma_{XX})^{-1} \sigma_{\epsilon}^2$.
- **two** losses of precision: imputation, data combination.
- One proxy case:
 - $\hat{\beta}^{AM}$, $\hat{\beta}^{RRP}$ and $\hat{\beta}^{BPP}$ are numerically identical
 - Start from $\hat{\beta}^{AM}$
- General case in paper.

Imputation on Full Data

- Asymptotic variance-covariance matrix for moments:

$$F = \begin{bmatrix} \sigma_u^2 \Sigma_{yy} & \beta \sigma_u^2 \Sigma_{XX} \\ \beta \sigma_u^2 \Sigma_{XX} & (\gamma^2 \sigma_u^2 + \sigma_\epsilon^2) \Sigma_{XX} \end{bmatrix}$$

- The asymptotic variance covariance matrix of (β, γ) is $(G'F^{-1}G)^{-1}$.
- The asymptotic variance of $\hat{\beta}$ is:

$$\text{Asymp. Var}(\hat{\beta}) = \frac{(\Sigma_{XX})^{-1} \sigma_\epsilon^2}{\phi_{yZ}}$$

- Loss of asymptotic precision proportional to the first stage R_{yZ}^2 .
 - Note similarity to linear IV (Shea, 1997).

OLS SE are biased

- With some algebra, can show that:

$$plim \left[\hat{V}^{OLS}(\hat{\beta}) \right] = \left[\frac{(\Sigma_{XX})^{-1} \sigma_{\epsilon}^2}{\phi_{yZ}} + \beta^2 \left(\frac{1 - \phi_{yZ}}{\phi_{yZ}} \right) \right]$$

- Too small by factor $\beta^2 \left(\frac{1 - \phi_{yZ}}{\phi_{yZ}} \right)$.
- Can be corrected using the available consistent estimates of β and ϕ_{yZ} .
- STATA command available from authors.

Monte Carlo: Design

- $x \sim N(0, 2)$
- $y = 1 + \beta x + \epsilon$ with $\sigma_\epsilon = 1$, $\beta = 1$
- $z_j = 1 + 0.5y_j + u_j$ with $\sigma_u = 1$
- First stage $R^2 = 0.56$
- Simulate population, draw samples (500) of (y_1, z_1) and (x_2, z_2) .
- 10,000 replications.

