# Imputing Consumptiom In The PSID Using Food Demand Estimates From The Cex

Richard Blundell    Luigi Pistaferri    Ian Preston

## December (2006)

Amin Shirazian

EF9905 Presentation

## Background

- Lack of panel data on total consumption restricted researchers to use the scanty food expenditure information in the PSID.
- Food is a necessity (i.e. the budget share for food falls as total expenditure rises), using food instead of consumption:
  1. prevent us from estimating a range of elasticity terms.
  2. generally underestimates total consumption volatility.
  3. prevent us from explaining price shocks.
- We need some type of a panel data for consumption for many researches.

## Outline

- **Motivation**: Help researchers with the lack of panel data on consumption (total expenditure) , using the available data:
  1. PSID : a panel data that contains food information, but not consumption.
  2. CEX : a detailed data set on household expenditures, including food expenditure, but it is a cross sectional data.

- **Strategy**: The idea is to find a relationship between consumption and food expenditure using CEX. Then, use the relationship ro impute consumption for PSID.

## Imputing Consumption

- Demand for food (notice that it captures price effect):

$$\tau\left(f_{i,x}\right) = D_{i,x}'\beta + \gamma\eta\left(c_{i,x}\right) + e_{i,x} \tag{1}$$

- $x \rightarrow$ observation from the CEX. (v.s. $p \rightarrow$ observation from PSID.)
- $f \rightarrow$ food expenditure (available in both CEX and PSID).
- $D \rightarrow$ prices and a set of conditioning variables (available in both data sets.)
- $c \rightarrow$ non-durable expenditure (available only in CEX)
- $e \rightarrow$ unobserved heterogeneity in the demand for food (including measurement error in food expenditure)

## Imputing Consumption

- Demand for food:

$$\tau\left(f_{i,x}\right) = D'_{i,x}\beta + \gamma\eta\left(c_{i,x}\right) + e_{i,x}$$

- **Assumption 1**: functions $\tau(.)$ and $\eta(.)$ are known monotonic increasing transformations of their arguments.
- **Assumption 2**: Food is a normal good ($\gamma \geq 0$).
- **Assumption 3**: Both data sets have the same underlying population.

## Imputing Consumption

- Imputed consumption in the CEX, **assuming** $\hat{\gamma} \neq 0$:

$$\widehat{c}_{i,x} = \eta^{-1} \left( \frac{\tau\left(f_{i,x}\right) - D'_{i,x}\widehat{\beta}}{\widehat{\gamma}} \right) \qquad (2)$$

- Imputed measure of consumption in the PSID:

$$\widehat{c}_{i,p} = \eta^{-1} \left( \frac{\tau\left(f_{i,p}\right) - D'_{i,p}\widehat{\beta}}{\widehat{\gamma}} \right) \qquad (3)$$

## Imputing Consumption

- To see how well our imputations are working, we need to compare them with the data points.

- Rewriting the imputed data we get the measurment error of the form:

$$\eta\left(\widehat{c}_{i,x}\right) = D'_{i,x}\frac{(\beta - \widehat{\beta})}{\widehat{\gamma}} + \frac{\gamma}{\widehat{\gamma}}\eta\left(c_{i,x}\right) + v_{i,x} \tag{4}$$

where $v_{i,x} = \frac{e_{i,x}}{\widehat{\gamma}}$.

## Simple Case

- Hence, the imputed data is simply an error riden data (with drift).
- **Assumption 4**: $\tau(x) = x$, and $\eta(x) = x$
- To see this better, unde the above assumption, we get:

$$\widehat{c}_{i,x} = \frac{(\beta - \widehat{\beta})}{\widehat{\gamma}} + \frac{\gamma}{\widehat{\gamma}} c_{i,x} + v_{i,x} \tag{5}$$

## Simple case

- assume that $c_{i,x}$ is potentially measured with classical error:

$$c_{i,x}^* = c_{i,x} + u_{i,x}$$

- under assumption 4, we get:

$$f_{i,x} = \beta + \gamma c_{i,x}^* + e_{i,x} - \gamma u_{i,x} \qquad (6)$$

- Notice that if, total expenditure decisions were made jointly with decisions on individual commodities, such as food, $\text{Cov}(c_x, e_x) \neq 0$

## Estimator

- Let $\widehat{\gamma}(y) = \frac{\text{Cov}(f_x, y_x)}{\text{Cov}(c_x^*, y_x)}$ be an estimator of $\gamma$. We have:

$$\text{plim} \, \widehat{\gamma}(y) = \gamma + B^e(y) + B^m(y)$$
$$\text{plim} \, \widehat{\beta}(y) = \beta - (B^e(y) + B^m(y)) \, \text{plim} \, M(c_x)$$

- where

$$B^e(y) = \frac{\text{plim} \, \text{Cov}(e_x, y_x)}{\text{plim} \, \text{Cov}(c_x^*, y_x)}$$

$$B^m(y) = -\gamma \frac{\text{plim} \, \text{Cov}(u_x, y_x)}{\text{plim} \, \text{Cov}(c_x^*, y_x)}$$

## Sample moments - CEX

- Let $\widehat{c}_x(y)$ denote the imputation with $\widehat{\beta}(y)$ and $\widehat{\gamma}(y)$, we have:

$$\text{plim } M\left(\widehat{c}_x(y)\right) = \text{ plim } M\left(c_x\right) \tag{7}$$

- And the variance term becomes:

$$\text{plim } V\left(\widehat{c}_x(y)\right) = \left(\frac{\gamma}{\gamma + B^e(y) + B^m(y)}\right)^2 \left(\text{plim } V\left(c_x\right) + \frac{1}{\gamma^2} \text{ plim } V\left(e_x\right) + \frac{2}{\gamma} \text{ plim Cov}\left(c_x, e_x\right)\right) \tag{8}$$

## Sample moments - CEX

- $M(\hat{c}_x(y))$ converges in probability to $M(c_x)$, regardless of **measurement error** or **heterogeneity in food spending**.
- Keeping the assumption that expenditure decision is made jointly, notice the last two terms on the variance term would bias our variance estimation.
- Also notice that the trend in variance is independent of those terms.

# IV case

- a valid instrument $z_x$ satisfies:

$$\text{plim Cov}(c_x^*, z_x) \neq 0$$
$$\text{plim Cov}(e_x, z_x) = 0$$
$$\text{plim Cov}(u_x, z_x) = 0$$

- then $B^e(z) = B^m(z) = 0$:

$\text{plim } \widehat{\gamma}(z) = \gamma$

$\text{plim } \widehat{\beta}(z) = \beta$

$\text{plim } M(\widehat{c}_x(z)) = \text{plim } M(c_x)$

$\text{plim } V(\widehat{c}_x(z)) = \text{plim } V(c_x) + \frac{1}{\gamma^2} \text{plim } V(e_x) + \frac{2}{\gamma} \text{plim Cov}(c_x, e_x)$

## OLS case

- In OLS case, $c^\star$ satisfies:

$$\text{Inseparability:} \qquad \text{plim Cov}\,(e_x, c_x^*) \neq 0$$
$$\text{Measurement:} \qquad \text{plim Cov}\,(u_x, c_x^*) \neq 0$$
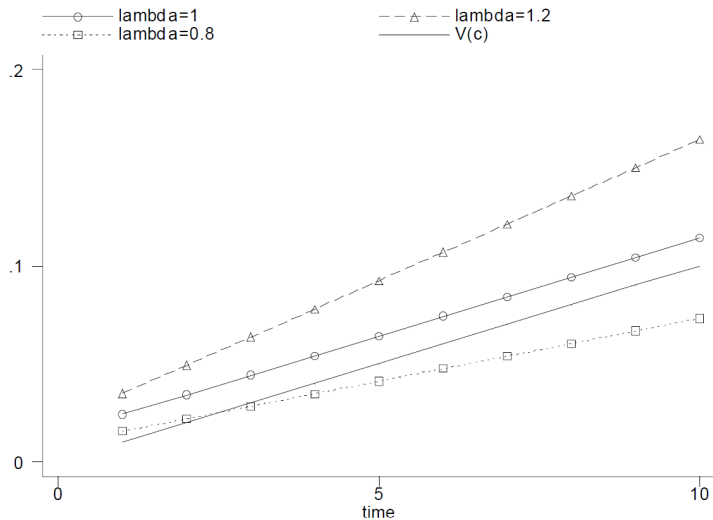
- then $B^e(z) \neq 0$, $B^m(z) \neq 0$:

$$\text{plim}\,\widehat{\gamma}(c^*) = \gamma + B^e(c^*) + B^m(c^*)$$
$$\text{plim}\,\widehat{\beta}(c^*) = \beta - (B^e(c^*) + B^m(c^*))\,\text{plim}\,M\,(c_x)$$
$$\text{plim}\,M\,(\widehat{c}_x(c^*)) = \text{plim}\,M\,(c_x)$$
$$\text{plim}\,V\,(\widehat{c}_x(c^*)) = \left(\frac{\gamma}{\gamma + B^e(y) + B^m(y)}\right)^2 \text{plim}\,V\,(\widehat{c}_x(z))$$

# OLS-IV comparison

- We want to use CEX sample moments to compute PSID's, for the mean, we have:

$$\text{plim } M\left(\widehat{c}_p(y)\right) = \text{p}\lim M\left(c_x\right) + \frac{1}{\gamma + B^e(y) + B^m(y)} \left[\text{p}\lim M\left(f_p\right) - \text{p}\lim M\left(f_x\right)\right] \tag{9}$$

- If food consumption is on average the same in the two data sets, the second term on the right hand side vanishes.

- Otherwise, the sample mean of imputed PSID consumption is potentially biased. (e.g. the two samples are not random samples drawn from the same underlying population)

# Sample moments - PSID

- For the Variance, we have:

$$
\text{plim } V\left(\hat{c}_p(y)\right) = \left(\frac{\gamma}{\gamma + B^e(y) + B^m(y)}\right)^2
\begin{pmatrix}
\text{plim } V\left(c_x\right) + \frac{1}{\gamma^2}\text{plim } V\left(e_x\right) + \\
\frac{2}{\gamma}\text{plim Cov}\left(c_x, e_x\right) + \\
\frac{1}{\gamma^2}\left(\text{plim } V\left(f_p\right) - \text{plim } V\left(f_x\right)\right)
\end{pmatrix}
\tag{10}
$$

- Notice that the slope term is not affected comparing to (8).
- There is an additional reason for $V(\hat{c}_p(y))$ be different from the population variance.

## Covariates

- Adding covariates to the simple model, we get:

$$f_{i,x} = D'_{i,x}\beta + \gamma c^*_{i,x} + e_{i,x} - \gamma u_{i,x}$$

$$D_{i,x} = \begin{pmatrix} 1 & _2d_{i,x} & _3d_{i,x} & \dots & _{k-1}d_{i,x} \end{pmatrix}'$$

- To fix ideas, let us consider a simple case:

$$f_{i,x} = \beta_0 + \beta_1 d_{i,x} + \gamma c^*_{i,x} + e_{i,x} - \gamma u_{i,x} \tag{11}$$

## Covariates

- It is easy to show:

$$\operatorname{p lim} \widehat{\beta}_0(d,y) = \beta_0 - (B^e(d,y) + B^m(d,y))(\operatorname{p lim} M(c) - \rho\operatorname{p lim} M(d))$$
$$\operatorname{p lim} \widehat{\beta}_1(d,y) = \beta_1 - \rho(B^e(d,y) + B^m(d,y))$$
$$\operatorname{p lim} \hat{\gamma}(d,y) = \gamma + B^e(d,y) + B^m(d,y)$$

- Where the bias terms are adjusted so that they can take the covariate into account:

$$B^e(d,y) = \frac{\operatorname{Cov}(e,y)V(d)}{V(d)\operatorname{Cov}(c^*,y) - \operatorname{Cov}(c^*,d)\operatorname{Cov}(d,y)}$$

$$B^m(d,y) = -\gamma\frac{\operatorname{Cov}(u,y)V(d)}{V(d)\operatorname{Cov}(c^*,y) - \operatorname{Cov}(c^*,d)\operatorname{Cov}(d,y)}$$

- and $\rho$ is defined as:

$$\rho = \frac{p\operatorname{lim}\operatorname{Cov}(c^*,d)}{p\operatorname{lim}V(d)}$$

## Covariates

- Sample means become:

$$\text{plim } M\left(\widehat{c}_x(d, y)\right) = \text{plim } M\left(c_x\right)$$

$$\text{plim } M\left(\widehat{c}_p(d, y)\right) = \text{plim } M\left(c_x\right) + \frac{1}{\gamma + B^e(d, y) + B^m(d, y)}$$

$$\left[\begin{array}{c} \text{pim } M\left(f_p - \widehat{\beta}_1(d, y)d_p\right) \\ -\text{plim } M\left(f_x - \widehat{\beta}_1(d, y)d_x\right) \end{array}\right]$$

- Convergence of the sample mean in CEX is assured.

- In PSID, sample mean may converge to a different value either because of discrepancy in the mean of the input variable ($f$) or the mean of the covariates ($d$) in the two data sets.
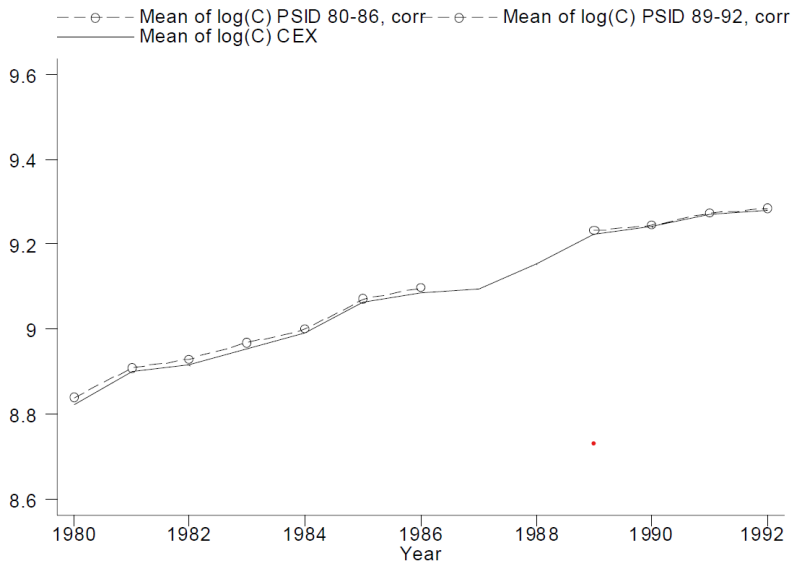
## Covariates

- For sample variance, we have:

$$\text{plim } V\left(\widehat{c}_x(d,y)\right) = \left(\frac{\gamma}{\gamma + B^e(d,y) + B^m(d,y)}\right)^2 \left[\begin{array}{c} p\lim V\left(c_x\right) + \frac{1}{\gamma^2} p\lim V\left(e_x\right) \\ + \frac{2}{\gamma}\text{plimCov}\left(c_x, e_x\right) \\ + \left(\frac{\rho(B^e(d,y) + B^m(d,y))}{\gamma}\right)^2 \text{plim } V\left(d_x\right) \\ + \frac{2\rho(B^e(d,y) + B^m(d,y))}{\gamma} p\lim Cov\left(c_x, d_x\right) \end{array}\right]$$

$$\text{plim } V\left(\widehat{c}_p(d,y)\right) = \text{plim } V\left(\widehat{c}_x(d,y)\right) + \left(\frac{1}{\gamma + B^e(d,y) + B^m(d,y)}\right)^2 \left[\begin{array}{c} \text{plim } V\left(f_p - \widehat{\beta}_1(y)d_p\right) \\ - \text{plim } V\left(f_x - \widehat{\beta}_1(y)d_x\right) \end{array}\right]$$
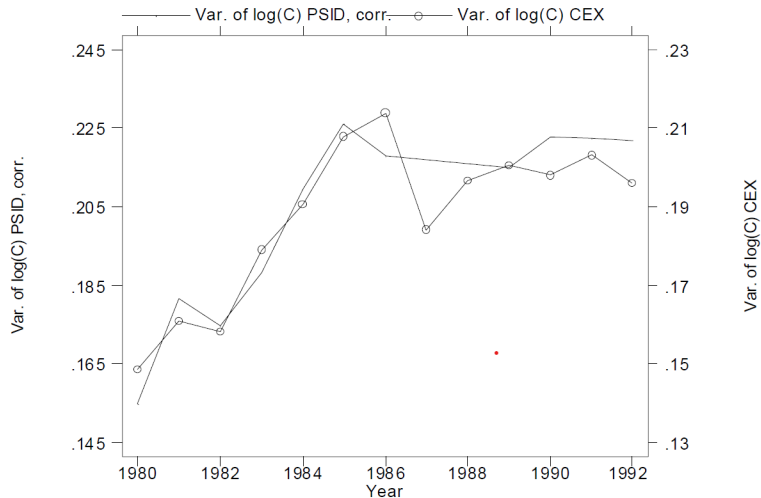
- Note also that if $y = z$ is a valid instrument, then covariates have no role in determining the asymptotic expression $V\left(\widehat{c}_x(d,y)\right)$.

# Results

- The paper expands for the case of having non-linear functional forms for $\tau(.)$ and $\eta(.)$, and having budget heterogeneity.
- Hourly wages are considered as instruments for total expenditures.
- After estimating the parameters of the model, figure 1 describes $(M(\widehat{c}_p) - M(c_x))$, and $\frac{M(f_p) - M(f_x)}{\hat{\gamma}}$.
- The main source of difference in the imputed and the true consumption is for the latter term.
- After correcting for the difference in mean food expenditure, we see the model performs well.

# Results-Corrected Mean

# Results-Corrected Variance

## Discussion

- We saw how once can use a combination of panel-Cross data to impute one variable from one to another.
- The idea is applicable to other datasets. (e.g. MEPS and HRS)
- Figuring if there is any other way that we can contribute to precision of imputation.