

Integration of anterior samples information and modeling in repeated survey with complex rotation pattern: an illustration with CPS

Workshop on Data Integration

D. Bonnér, JPSM, UMD

Based on joint work with

Y. Cheng (Census Bureau) and **P. Lahiri** (JPSM, UMD).

August 6, 2019

Any views expressed are those of the authors and not necessarily those of the U.S. Census Bureau.

- ❶ The Current Population Survey (CPS)
- ❷ Model-assisted approach
- ❸ Maximum likelihood estimation
- ❹ Application to CPS

The Current Population Survey (CPS)

- American labor force statistics
- Stratified multi-stage sampling design
- Repeated survey with
 - 8 rotation groups (samples of households) per month
 - rotation pattern: 4-8-4

- m : time indicator (month).
- U_m : the population at time m ,
- S_m : the sample at time m ,
- $w_{m,k}$: the sampling adjusted weight individual k at time m ,
- Employment status of individual k at time m is:

$$y_{k,m} = \begin{cases} \text{status}_1 = (1, 0, 0) & \text{if } k \text{ is "employed"} \\ \text{status}_2 = (0, 1, 0) & \text{if } k \text{ is "unemployed"} \\ \text{status}_3 = (0, 0, 1) & \text{if } k \text{ is "not in the labor force"} \end{cases}$$

- $t_{y_m} = \sum_{k \in U_m} \mathbf{y}_{k,m}$.

All rotation groups G_1, \dots drawn simultaneously.

Table: CPS Rotation chart

	G_1	G_2	G_3	G_4	G_5	G_6	G_7	G_8	G_9	G_{10}	G_{11}	G_{12}	G_{13}	G_{14}	G_{15}	G_{16}	G_{17}	G_{18}	G_{19}	G_{20}	
Jan 05	S_{18}	$S_{1,7}$	S_{16}	S_{15}										S_{14}	$S_{1,3}$	$S_{1,2}$	$S_{1,1}$				
Feb 05		S_{28}	$S_{2,7}$	S_{26}	S_{25}									S_{24}	$S_{2,3}$	$S_{2,2}$	$S_{2,1}$				
Mar 05			S_{38}	$S_{3,7}$	S_{36}	S_{35}									S_{34}	$S_{3,3}$	$S_{3,2}$	$S_{3,1}$			
Apr 05				S_{48}	$S_{4,7}$	S_{46}	S_{45}									S_{44}	$S_{4,3}$	$S_{4,2}$	$S_{4,1}$		
May 05					S_{58}	$S_{5,7}$	S_{56}	S_{55}									S_{54}	$S_{5,3}$	$S_{5,2}$	$S_{5,1}$	
Jun 05						S_{68}	$S_{6,7}$	S_{66}	S_{65}									S_{64}	$S_{6,3}$	$S_{6,2}$	
Jul 05							S_{78}	$S_{7,7}$	$S_{7,6}$	$S_{7,5}$								S_{74}	$S_{7,3}$		
Aug 05								S_{88}	$S_{8,7}$	S_{86}	S_{85}									S_{84}	
Sep 05									S_{98}	$S_{9,7}$	S_{96}	S_{95}									
Oct 05										S_{108}	$S_{10,7}$	S_{106}	S_{105}								
Nov 05											S_{118}	$S_{11,7}$	S_{116}	S_{115}							
Dec 05												S_{128}	$S_{12,7}$	S_{126}	S_{125}						
Jan 06													S_{138}	$S_{13,7}$	S_{136}	S_{135}					
Feb 06														S_{148}	$S_{14,7}$	S_{146}	S_{145}				
Mar 06															S_{158}	$S_{15,7}$	S_{156}	S_{155}			
Apr 06																S_{168}	$S_{16,7}$	S_{166}	S_{165}		
May 06																	S_{178}	$S_{17,7}$	S_{176}	S_{175}	
Jun 06																		S_{188}	$S_{18,7}$	S_{186}	
Jul 06																			S_{198}	$S_{19,7}$	
Aug 06																				S_{208}	

$S_{m,r}$ be the r th rotation group for month m , $S_m = \bigcup_{r=1}^8 S_{m,r}$.

- How to integrate the information from the previous months ?
- How to account for the special rotation pattern ?

For $m \in \{1, \dots, M\}$, $r \in \{1, \dots, 8\}$, $i \in \{1, \dots, 3\}$,

$$\hat{t}_{y_m}^{\text{mis},r} = 8 \sum_{k \in S_{m,r}} w_{m,k} y_{k,m}.$$

When $\text{Cov} \left[\hat{t}_{y_{m,i}}^{\text{mis},r}, \hat{t}_{y_{m',i'}}^{\text{mis},r'} \right]$ is known, the best linear combination of the m.i.s. estimates can be defined.

Currently at the Census, the AK estimator is the empirical best linear estimator in a subclass of linear estimators.

- Unrealistic assumptions were made on the covariance matrix: 10-year stationarity. To our knowledge, the model assumptions on the covariance matrix were not studied or called into question.
- Outdated optimisation techniques were used to compute the optimal coefficients.
- The optimal AK estimator for month to month change is not optimal for level estimation.
- The estimators are sensible to the "rotation bias".

Other linear estimates (Yansaneh-Fuller) suffer from the same problem of estimation of covariance.

Regression composite estimators (Statistics Canada)

- Depend on a tuning parameter α .
- Do provide a compromise optimal estimator.
- Perform well on simulations.
- Why is it working ?

Model Assisted estimation

Consider auxiliary variables \mathbf{x} of known total $t_{\mathbf{x}}$.

In a **model-assisted** approach, Generalized Regression is associated to a **superpopulation linear model**: $\forall k \in U$, y_k, x_k distribution is such

$$E[y_k | x_k] = x_k \beta.$$

Then the Generalized Regression estimator is

$$\hat{t}_{\mathbf{y}}^{\text{GREG}} = \hat{t}_{\mathbf{y}}^{\text{Direct}} + (t_{\mathbf{x}} - \hat{t}_{\mathbf{x}}^{\text{Direct}}) \hat{\beta}.$$

It will be design consistent irrespective of the model and effective if the model fits the data.

The generalized regression estimator is also equal to

$$\hat{t}_y^c = \sum_{k \in S} w_k^* y_k,$$

where w_k^* minimizes $\sum_k (w_k^* - w_k)^2 / w_k$ under the constraint:

$$\sum_{k \in S} w_k^* x_k = t_x.$$

Since \mathbf{y}_{m-1} is correlated to \mathbf{y}_m , why not choose $\mathbf{x}_m = \mathbf{y}_{m-1}$?

Note that the linear assumption holds:

$$E [y_{k,m} \mid y_{k,m-1}] = y_{k,m-1} \beta_m,$$

where

$$\beta_m = P(y_{m,k} = \text{status}_{i'} \mid y_{m-1,k} = \text{status}_i)_{i,i'=1,\dots,3}.$$

The GREG estimator is

$$\hat{t}_{\mathbf{y}_m}^{\text{GREG}} = \hat{t}_{\mathbf{y}_m}^{\text{Direct}} + \left(t_{\mathbf{y}_{m-1}} - \hat{t}_{\mathbf{y}_{m-1}}^{\text{Direct}} \right) \hat{\beta}_m.$$

But

- the **characteristic** $y_{m-1,k}$ is unknown for $k \in S_m \setminus S_{m-1}$,
- the **total** $t_{\mathbf{y}_{m-1}}$ is unknown.

 Singh, A. and Merkouris, P. (1995).

Composite estimation by modified regression for repeated surveys.

In *ASA Proc. Surv. Res. Meth. Sec*, pages 420–5.

Solution:

- replace **characteristic** $y_{m-1,k}$ by a proxy variable,
- replace **total** $t_{y_{m-1}}$ by an estimation.

Proxy variables:

- MR1:

$$x_{m,k,MR1} = \begin{cases} y_{m-1,k} & \text{if } k \in S_m \cap S_{m-1}, \\ N_{m-1}^{-1} \hat{t}_{y,t-1}^c & \text{if } k \in S_m \setminus S_{m-1}, \end{cases}$$

- MR2:

$$x_{m,k,MR2} = \begin{cases} (\tau^{-1}[y_{m-1,k} - y_{m,k}] + y_{m,k}) & \text{if } k \in S_m \cap S_{m-1}, \\ y_{m,k} & \text{if } k \in S_m \setminus S_{m-1}, \end{cases}$$

- MR1 will carry misrepresentativeness of first sample over time.
- MR2 will chain evolutions without taking into account the levels of intrant rotation groups.

Solution: $x_{m,k,\text{MR3}} = [x_{m,k,\text{MR1}}, x_{m,k,\text{MR2}}]$

 Fuller, W. A. and Rao, J. (2001).

A regression composite estimator with application to the canadian labour force survey.

Survey Methodology, 27(1):45–52.

Solution: define a new proxy variable. Let $\alpha \in [0, 1]$:

$$x_{m,k} = \alpha x_{m,k,\text{MR1}} + (1 - \alpha) x_{m,k,\text{MR2}}$$

Why is it working ?

The calibration equation can be rewritten:

$$\begin{aligned}
 & \left(\sum_{k \in S_m \cap S_{m-1}} w_{m,k}^* y_{m-1,k} + \sum_{k \in S_m \setminus S_{m-1}} w_{m,k}^* \hat{y}_{m-1,k} \right) - \hat{t}_{\mathbf{y}, m-1}^c \\
 = & (1 - \alpha) \left((\tau^{-1} - 1) \left(\sum_{k \in S_m \cap S_{m-1}} w_{m,k}^* (y_{m,k} - y_{m-1,k}) \right) \right. \\
 & \quad \left. - \sum_{k \in S_t \setminus S_{m-1}} w_{m,k}^* (y_{m,k} - \hat{y}_{m-1,k}) \right)
 \end{aligned}$$

where $\hat{y}_{k,m-1} = N_{m-1}^{-1} \hat{t}_{\mathbf{y}, m-1}^c$, $\tau = \frac{3}{4}$

We notice that:

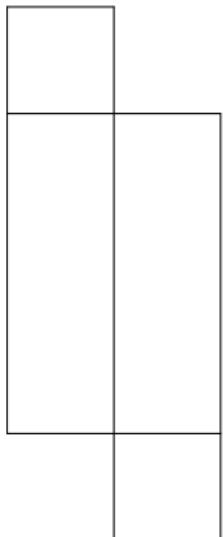
- Composite Regression calibration equation is a balance between two hard constraints: new set of weights will compromise between level and change.
- Does not require revision of $\hat{t}_{y,m-1}$

Unanswered questions:

- How to compute the variance of this estimator ?
- How to choose the best α ?

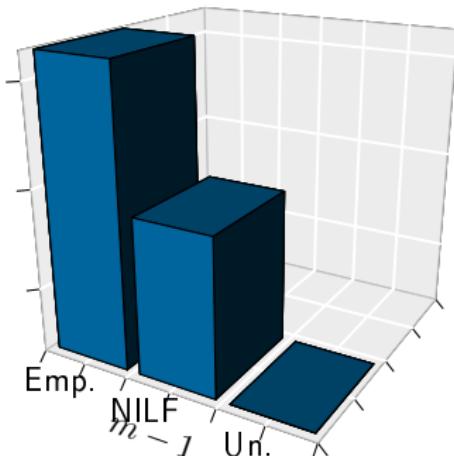
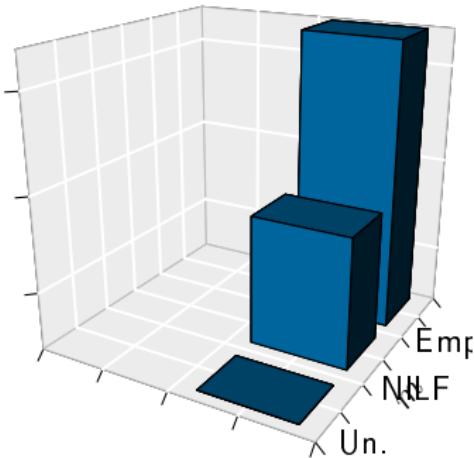
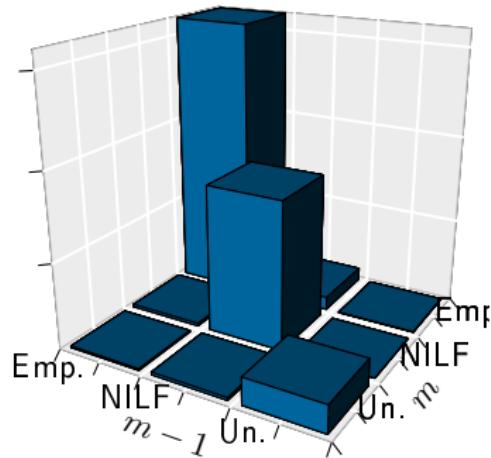
Take a step back:

S_{m-1} S_m



$(N_{1,.}, N_{2,.}, N_{3,.})$

$$\begin{pmatrix} N_{1,1} & N_{1,2} & N_{1,3} \\ N_{2,1} & N_{2,2} & N_{2,3} \\ N_{3,1} & N_{3,2} & N_{3,3} \end{pmatrix}$$
$$(N_{.,1}, N_{.,2}, N_{.,3})$$



When

$$N_{i_1,i_2} \sim \text{Multinomial}(\gamma_{i_1,i_2}),$$

the likelihood is

$$\begin{aligned} & \sum_{i_1=1}^3 N_{i_1,.} \ln (\gamma_{i_1,1} + \gamma_{i_1,2} + \gamma_{i_1,3}) \\ & + \sum_{i_1,i_2=1}^3 N_{i_1,.} \ln (\gamma_{i_1,i_2}) \\ & + \sum_{i_2=1}^3 N_{.,i_2} \ln (\gamma_{1,i_2} + \gamma_{2,i_2} + \gamma_{3,i_2}). \end{aligned}$$

We assume the following population model: $\forall k, (y_{m,k})_{m \in \{1, \dots, M\}}$ are the i.i.d realisation of a random variable $Y = (Y_1, \dots, Y_m)$. Let γ a $\underbrace{3 \times 3 \times \dots \times 3}_{M \text{ times}}$ array, such that

$$\gamma_{i_1, \dots, i_M} = P \left(\bigcap_{m=1}^M Y_m = \text{status}_{i_m} \right).$$

Example: for $M = 3$,

$$\gamma_{1,3,3} = P (\{Y_1 = \text{status}_1\} \cap \{Y_2 = \text{status}_3\} \cap \{Y_3 = \text{status}_3\})$$

$$P(Y_3 = \text{status}_1)$$

$$= \sum_{i_1=1}^3 \sum_{i_2=1}^3 P(\{Y_1 = \text{status}_{i_1}\} \cap \\ \{Y_2 = \text{status}_{i_2}\} \cap \{Y_3 = \text{status}_1\})$$

$$P(Y_M = \text{status}_1) = \sum_{i_1, \dots, i_{M-1}=1}^3 \gamma_{i_1, \dots, i_{M-1}, 1}$$

We do not observe y_k , but $y_k|_{A_k}$, for example

$$A_k = \{1, 2, 3, 4, 13, 14, 15, 16\},$$

$$A_k = \{2, 3, 4, 5, 14, 15, 16, 17\}.$$

The contribution to the log-likelihood of $y_k|_{A_k}$ is:

$$w_k \ln\left(\text{P}\left(Y|_{A_k} = y_k|_{A_k}\right)\right) = w_k \ln\left(\sum \{\gamma_{i_1, \dots, i_M} \mid \forall k \in A_k, i_m = y_{k,m}\}\right).$$

So the log-likelihood is:

$$\begin{aligned} & \mathcal{L}\left(\gamma; \left(y_k|_{A_k}\right)_{k \in \bigcup_{m=1}^M S_m}\right) \\ &= \sum_{k \in \bigcup_{m=1}^M S_m} w_k \ln\left(\text{P}\left(Y|_{A_k} = y_k|_{A_k}\right)\right) \\ &= \sum_{k \in \bigcup_{m=1}^M S_m} w_k \ln \left(\sum_{i_1, \dots, i_M \mid \forall m \in A_k, i_m = y_{k,m}} \gamma_{i_1, \dots, i_M} \right). \end{aligned}$$

If $M = 2$, then $A_k \in \{\{1\}, \{1, 2\}, \{2\}\}$. Let

- $N_{i_1,.} = \sum \left\{ w_k \mid k \in \bigcup_{m=1}^M S_m, A_k = \{1\}, y_1 = i_1 \right\}$
- $N_{i_1,i_2} = \sum \left\{ w_k \mid k \in \bigcup_{m=1}^M S_m, A_k = \{1, 2\}, y_1 = i_1, y_2 = i_2 \right\}$
- $N_{.,i_2} = \sum \left\{ w_k \mid k \in \bigcup_{m=1}^M S_m, A_k = \{1\}, y_2 = i_2 \right\}$

The likelihood is

$$\begin{aligned} & \mathcal{L} \left(\gamma; \left(y_k \mid_{A_k} \right)_{k \in \bigcup_{m=1}^M S_m} \right) \\ &= \sum_{i_1=1}^3 N_{i_1,.} \ln (\gamma_{i_1,1} + \gamma_{i_1,2} + \gamma_{i_1,3}) \\ & \quad + \sum_{i_1,i_2=1}^3 N_{i_1,i_2} \ln (\gamma_{i_1,i_2}) + \sum_{i_1=1}^3 N_{.,i_2} \ln (\gamma_{1,i_2} + \gamma_{2,i_2} + \gamma_{3,i_2}). \end{aligned}$$

If $M = 3$, then $A_k \in \{\{1\}, \{1, 2\}, \{1, 2, 3\}, \{2, 3\}, \{3\}\}$.

Let

- $N_{i_1,..} = \sum \left\{ w_k \mid k \in \bigcup_{m=1}^M S_m, A_k = \{1\}, y_1 = i_1 \right\}$
- $N_{i_1,i_2,.} = \sum \left\{ w_k \mid k \in \bigcup_{m=1}^M S_m, A_k = \{1, 2\}, y_1 = i_1, y_2 = i_2 \right\}$
- $N_{i_1,i_2,i_3} = \sum \left\{ w_k \mid k \in \bigcup_{m=1}^M S_m, A_k = \{1, 2, 3\}, y_1 = i_1, y_2 = i_2, y_3 = i_3 \right\}$
- $N_{.,i_2,i_3} = \sum \left\{ w_k \mid k \in \bigcup_{m=1}^M S_m, A_k = \{2, 3\}, y_2 = i_2, y_3 = i_3 \right\}$
- $N_{.,.,i_3} = \sum \left\{ w_k \mid k \in \bigcup_{m=1}^M S_m, A_k = \{3\}, y_3 = i_3 \right\}$

$$\begin{aligned}
 \mathcal{L} \left(\gamma; \left(y_k | A_k \right)_{k \in \bigcup_{m=1}^M S_m} \right) = & \sum_{i_1=1}^3 N_{i_1, \dots} \ln \left(\sum_{i_2, i_3=1}^3 \gamma_{i_1, i_2, i_3} \right) \\
 & + \sum_{i_1, i_2=1}^3 N_{i_1, i_2, \cdot} \ln \left(\sum_{i_3=1}^3 \gamma_{i_1, i_2, i_3} \right) \\
 & + \sum_{i_1, i_2, i_3=1}^3 N_{i_1, i_2, i_3} \ln (\gamma_{i_1, i_2, i_3}) \\
 & + \sum_{i_2, i_3=1}^3 N_{\cdot, i_2, i_3} \ln \left(\sum_{i_1=1}^3 \gamma_{i_1, i_2, i_3} \right) \\
 & + \sum_{i_3=1}^3 N_{\cdot, \cdot, i_3} \ln \left(\sum_{i_1, i_2=1}^3 \gamma_{i_1, i_2, i_3} \right).
 \end{aligned}$$

- No assumption on the stochastic structure.
- Optimality obtained for all function of γ : month to month change, month level for totals, employment rates...

Let

$$\hat{\gamma}^{\text{MLE}} \in \operatorname{argmax}_{\gamma} \left\{ \mathcal{L} \left(\gamma; \left(y_k | A_k \right)_{k \in \bigcup_{m=1}^M S_m} \right) \right\}.$$

Define

$$\hat{\gamma}^{\text{MLE}}|_{m,i} = \sum_{i'_1, \dots, i'_{m-1}, i'_{m+1}, i'_M} \gamma_{i'_1, \dots, i'_{m-1}, i, i'_{m+1}, i'_M}$$

Define

$$\hat{u}_m^{\text{MLE}} = \frac{\hat{\gamma}^{\text{MLE}}|_{m,2}}{\hat{\gamma}^{\text{MLE}}|_{m,1} + \hat{\gamma}^{\text{MLE}}|_{m,2}}$$

- Choice of number of lags, for example -12,-3,-2,-1,0,1,2,3,12.
- Define a best compromise if no revision is required.
- Integration into existing adjustment procedure: weight adjustment based on MLE estimates.

github.com/DanielBonnery/pubBonneryChengLahiri2016

- Choice of lags -2,-1,0,1,2.
- Generation of 3 synthetic populations.
- CSP design was mimicked.
- MSE was computed with Monte Carlo
- MLE performs (slightly) better than Regression Composite, for level and month to month change.

Why?

- More information is used.
- No need to compromise between level and change.

M	$\dim(\gamma)$
5	243
10	59,049

Proposition:

- For each month, select lags (ex: 12,3,2,1,0,-1,-2,-3,-12)
- Compute MLE.
- Under the no revision constraint, define a criterium for best compromise

Current work.

- Parameter dimension limits the information that can be used.
- How to reduce the model dimension when modeling the microdata

Solution

- Use machine learning methods to predict each individual outcome, with a measure of uncertainty.
- Include additional variables.
- Integrate the uncertainty.

Thank you for your attention.