Synthetic Area Weighting for Measuring Public Opinion in Small Areas

Shiro Kuriwaki and Soichiro Yamauchi

PolMeth 2021

July 2021

Small Area Estimation

Overview

• Enables local geography / subgroup estimates



 Dominant method: Multilevel Regression + Poststratification (MRP)

But at what cost?

"

I think MRP is good. I think it's overrated also ... [I understand MRP as asking] what the polls would look like in a particular state or district without having any polls of that state or district... it's a good approximation, but it loses a lot of local variation."

Nate Silver (2018)



What this paper does

- 1. Derive identification conditions for unbiasedness when borrowing information across small areas
- 2. A weighting approach clarifies conditions
 - · compared to an outcome-based approach like MRP
- 3. We leverage existing data not used in MRP
 - initial (non-small area) survey weights
 - · covariates only measured in the survey

Our method in a nutshell: Estimating FL-27 (Miami)

- Index people by *i*, areas of by *j*
- Survey inclusion for person $i: S_i \in \{0, 1\}$
- Person *i* is in Area *j*: $A_{ij} \in \{0, 1\}$
- Estimand: $\mathbb{E}\left[Y_i|A_{ij}=1\right]$

Combine

1. Direct estimator

(Only use FL-27 respondents)

• Weight by inverse of $Pr(S_i = 1 | Covariates_i, A_{ij} = 1)$

2. Indirect estimator

(Use non-FL-27 respondents, reweight to look like FL-27)

- Weight by $\Pr(A_{ij} = 1 | \text{Covariates}_i, S_i = 1)$
- Weight them again to look like population

Paper at arXiv: 2105.05829

Partial pooling needs two identification conditions

Assumption 1 (Sampling Ignorability)

$$Y_i \perp S_i \mid \mathbf{X}_i^P, A_{ij} = 1$$

X^P_i: <u>Poststratification variables</u> (with <u>Population</u> target)

e.g. age, sex

- · Required for almost any survey estimate
- → If satisfied, **direct estimator** using weights

$$\frac{1}{\Pr\left(S_i = 1 | \mathbf{X}_i^P, A_{ij} = 1\right)}$$

is unbiased.



Partial pooling needs two identification conditions

Assumption 2 (Area Ignorability) For each area of interest *i*,

$$Y_i \perp A_{ii} \mid X_i^P, \mathbf{X}_i^S, S_i = 1$$

- X_i^S: Covariates only in the <u>Survey Sample</u> e.g., party ID, news interest
- Unique to small area estimation
- → If satisfied, **indirect estimator** using weights

$$\frac{1}{\Pr\left(S_i = 1 | \mathbf{X}_i^P, A_{ij} = 1\right)} \cdot \frac{\Pr(A_{ij} = 1 | \mathbf{X}_i^P, \mathbf{X}_i^S, S_i = 1)}{1 - \Pr(A_{ij} = 1 | \mathbf{X}_i^P, \mathbf{X}_i^S, S_i = 1)}$$

Target Area $A_{ii} = 1$ Outside of Area $A_{ii} = 0$ Target Population $S_i = 0$ $S_i = 0, A_{ij} = 1$ $S_i = 0, A_{ij} = 0$ **Assumption 1** (Sampling Ignorability) Survey Sample $S_i = 0$ $S_i = 1, A_{ij} = 0$ $S_{i} = 1, A_{ij} = 1$ Assumption 2 (Area Ignorability)

is unbiased.

Our synthetic area estimator

Simplest case: if post-stratification covariates (X_i^P) were sufficient, then compute weights

$$\widehat{w}_{ij}^{\text{SA}} \propto \Pr(A_{ij} = 1 \mid X_i^P) \times \underbrace{\frac{1}{\widehat{\Pr}(S_i = 1 \mid X_i^P)}}_{\text{Area adjustment}} \times \underbrace{\frac{1}{\widehat{\Pr}(S_i = 1 \mid X_i^P)}}_{\text{Selection adjustment}}$$

and take the weighted average across all respondents.

Result: With Assumptions 1 and 2, $\sum_{i=1}^{n} \widehat{w}_{ij}^{SA} Y_i$ is unbiased for $\mathbb{E}[Y_i | A_{ij} = 1]$.

More generally,

$$\widehat{w}_{ij}^{\text{SA}} \propto \frac{\widehat{\Pr}(A_{ij} = 1 \mid X_i^P, X_i^S, S_i = 1)}{\widehat{\Pr}(A_{ij} = 1 \mid X_i^P, S_i = 1)} \times \Pr(A_{ij} = 1 \mid X_i^P) \times \frac{1}{\widehat{\Pr}(S_i = 1 \mid X_i^P)}$$

Difference of our approach vs. other methods

1. vs. Traditional MRP

- ↔ Fay and Herriot (1979) style estimator partially pools by random effects (global + local shrinkage estimator)
- 2. vs. Machine learning (MR)P
 - \rightsquigarrow Ghitza and Gelman (2013); Bisbee (2019); Ornstein (2020); Goplerud (2020) all extract deep interactions from X_i^P , less on partially pooling
- 3. vs. "Multilevel regression with synthetic poststratification"
 - Leemann and Wasserfallen (2017) expand X_i^P through missing data estimation.
 (We estimate the *area*, not *population*, synthetically)
- 4. vs. "Subgroup Balancing Propensity Score"
 - → Ben-Michael, Feller, and Rothstein estimate the propensity score by partial pooling, but not the outcome

Paper at arXiv: 2105.05829

Combined estimator reduces variance and bias

• CCES 2018

• 63 Congressional Districts in Texas + Florida, $n \approx 140$ each

- MRP vs. Proposed estimator (synthArea)
 - X_i^P : Age group + sex + education
 - X_i^S : Race × Party ID × News interest



Insight 1: The combined estimator is a rescaling of national weights

$$\widehat{w}_{ij}^{\mathsf{SA}} \propto \frac{\widehat{\Pr}(A_{ij} = 1 \mid X_i^P, X_i^S, S_i = 1)}{\widehat{\Pr}(A_{ij} = 1 \mid X_i^P, S_i = 1)} \times \Pr(A_{ij} = 1 \mid X_i^P) \times \underbrace{\frac{1}{\widehat{\Pr}(S_i = 1 \mid X_i^P)}}_{\mathsf{Plug-in } w_i^{\mathsf{national}}}$$

The survey weights w^{national}...

- · are already given in the dataset
- adjust for covariates beyond researcher's X_i^P
- but not used in MRP

Insight 2: Violations of area ignorability can be empirically tested

- For each area *j*, regress $Y_i \sim A_{ij} + X_i^P + X_i^S$
- Area ignorability is valid for area *j* when the coefficient on *A_{ij}* is 0.



Insight 3: Lack of covariates can cause partial pooling to "oversmooth"



Takeaways

- 1. Small Area Estimation is not assumption-free
- 2. "Area Ignorability" is key assumption in Small Area Estimation
- 3. Synthetic weighting (R package synthArea) can incorporate existing pollster's survey weights and survey-only variables