

# Synthetic Area Weighting for Measuring Public Opinion in Small Areas

**Shiro Kuriwaki and Soichiro Yamauchi**

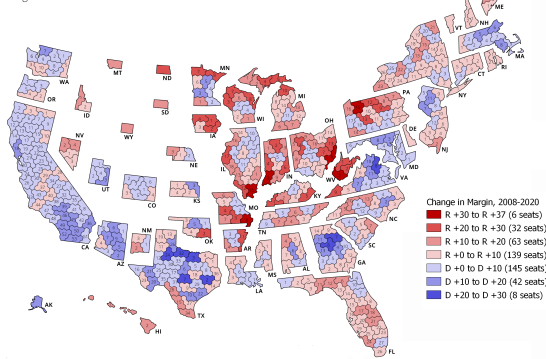
PolMeth 2021

July 2021

# Small Area Estimation

- Enables local geography / subgroup estimates

2008 to 2020 Trend in Presidential Election Margin by Congressional District  
For Congressional Districts used in 2020 Elections



DAILY KOS  
ELECTIONS

- 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:  $\mathbb{E}[Y_i | A_{ij} = 1]$

Combine

## 1. Direct estimator

(Only use FL-27 respondents)

- Weight by inverse of  $\Pr(S_i = 1 \mid \text{Covariates}_i, A_{ij} = 1)$

## 2. Indirect estimator

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

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

# Partial pooling needs two identification conditions

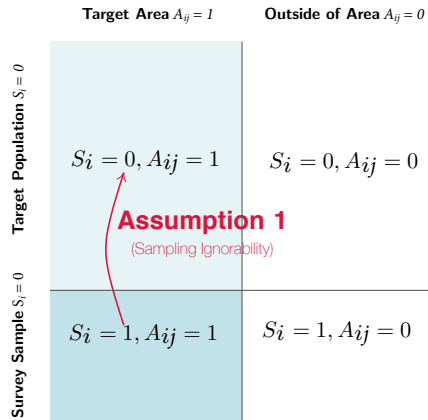
## Assumption 1 (Sampling Ignorability)

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

- $\mathbf{X}_i^P$ : Poststratification variables (with Population target)  
e.g. age, sex
  - Required for almost any survey estimate
- If satisfied, **direct estimator** using weights

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

is unbiased.



# Partial pooling needs two identification conditions

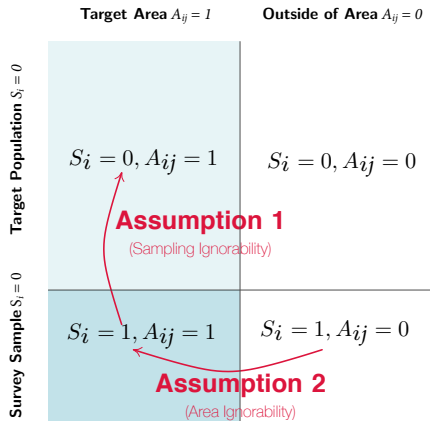
## Assumption 2 (Area Ignorability)

$$Y_i \perp\!\!\!\perp A_i \mid X_i^P, X_i^S, S_i = 1$$

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

$$\frac{1}{\Pr(S_i = 1 | X_i^P, A_{ij} = 1)} \cdot \frac{\Pr(A_{ij} = 1 | X_i^P, X_i^S, S_i = 1)}{1 - \Pr(A_{ij} = 1 | X_i^P, X_i^S, S_i = 1)}$$

is unbiased.



## Our synthetic area estimator

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

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

and take the weighted average across all respondents.

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

More generally,

$$\hat{w}_{ij}^{\text{SA}} \propto \frac{\overbrace{\hat{\Pr}(A_i = j \mid X_i^P, X_i^S, S_i = 1)}^{\text{Adjustment from new } X_i^S}}{\hat{\Pr}(A_i = j \mid X_i^P, S_i = 1)} \times \Pr(A_i = j \mid X_i^P) \times \frac{1}{\hat{\Pr}(S_i = 1 \mid X_i^P)}$$

# 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

~> Ghitza and Gelman (2013); Bisbee (2019); Ornstein (2020); Goplerud (2020) all extract as many interactions as possible from  $X_i^P$ , not 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)

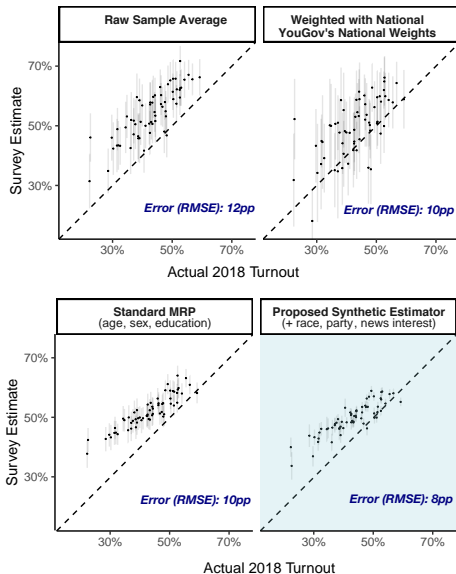
## 4. vs. “Subgroup Balancing Propensity Score”

~> Ben-Michael, Feller, and Rothstein estimate the propensity score by partial pooling, but not the outcome



# 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  $\times$  Party ID  $\times$  News interest



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

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

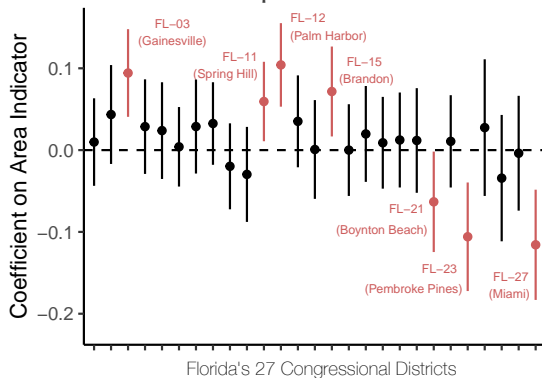
The survey weights  $w^{\text{national}}$  ...

- already given in the dataset
- adjustment for covariates beyond researcher's  $X_i^P$
- not used in MRP

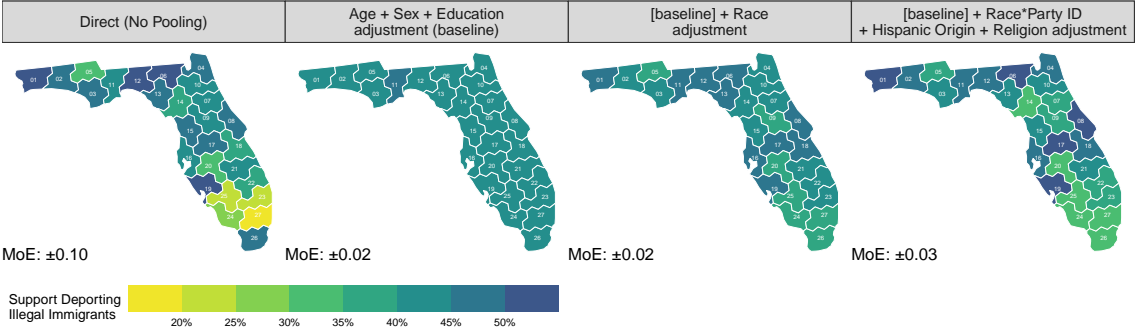
## Insight 2: Violations of area ignorability can be empirically tested

- 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.

Florida turnout example:



# Insight 3: Lack of covariates can cause partial pooling to “oversmooth”



# Takeaways

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