

Towards a Cautious Modelling of Missing Data in Small Area Estimation

Statistische Woche 2017, Rostock

Julia Plass¹, Aziz Omar^{1,2}, Thomas Augustin¹

¹ Department of Statistics, Ludwig-Maximilians University and

² Department of Mathematics, Insurance and Appl. Statistics, Helwan University

21st of September 2017



Aziz Omar

“Towards a Cautious
Modelling of Missing Data
in Small Area Estimation”

Thomas Augustin
Julia Plass



Aziz Omar

“Towards a Cautious
Modelling of **Missing Data**
in **Small Area Estimation**”

Thomas Augustin
Julia Plass



Aziz Omar

“Towards a Cautious
Modelling of Missing Data
in Small Area Estimation”

Thomas Augustin
Julia Plass



Aziz Omar

“Towards a **Cautious**
Modelling of **Missing Data**
in **Small Area Estimation**”

Thomas Augustin
Julia Plass

- Existing approaches for dealing with nonresponse in SAE are based on strong assumptions on the missingness process
- Such assumptions are usually **not testable**, and wrongly imposing them may lead to **biased** results.



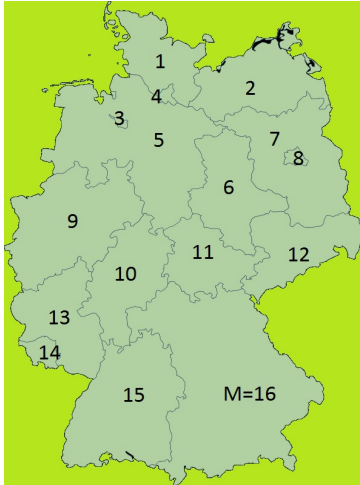
(Manski, 2003, Partial Identification of Probability Distributions, Jaeger, 2006, ECML,...)

What's the problem? \Rightarrow 1. Small Area Estimation (SAE)



- Population with N individuals

What's the problem? \Rightarrow 1. Small Area Estimation (SAE)



- Population with N individuals
- M areas, each contains N_i individuals, $i = 1, \dots, M$

What's the problem? \Rightarrow 1. Small Area Estimation (SAE)



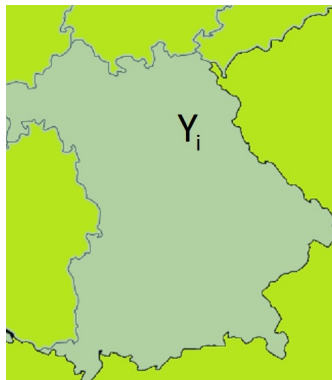
- Population with N individuals
- M areas, each contains N_i individuals, $i = 1, \dots, M$
- **Of interest:**
Area-specific mean \bar{Y}_i

What's the problem? \Rightarrow 1. Small Area Estimation (SAE)



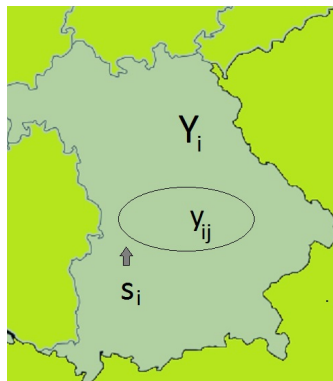
- Population with N individuals
 - M areas, each contains N_i individuals, $i = 1, \dots, M$
 - **Of interest:**
Area-specific mean \bar{Y}_i
 - **Problem:**
For each area, only sample s_i with small sample size n_i available
- \Rightarrow Using auxiliary variables (covariates) X_1, \dots, X_k
- \Rightarrow “borrowing strength”

What's the problem? \Rightarrow 1. Small Area Estimation (SAE)



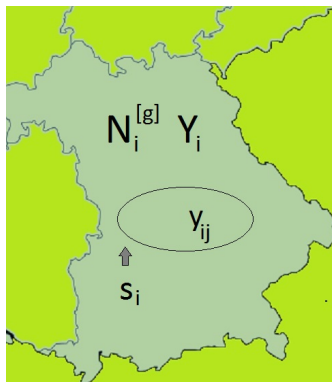
- Binary variable of interest
 \Rightarrow probability that Y_i is equal to 1
:= π_i (poverty rate)

What's the problem? \Rightarrow 1. Small Area Estimation (SAE)



- Binary variable of interest
 \Rightarrow probability that Y_i is equal to 1
 $:= \pi_i$ (poverty rate)
- $1/w_{ij}$ is the probability that individual j in area i is selected in s_i
- Sample values y_{ij} known for $j \in s_i$
- Sample data from German General Social Survey (GESIS Leibniz Institute for the Social Sciences, 2016), $y_{ij} = 1$: 'poor', $y_{ij} = 0$: 'rich'

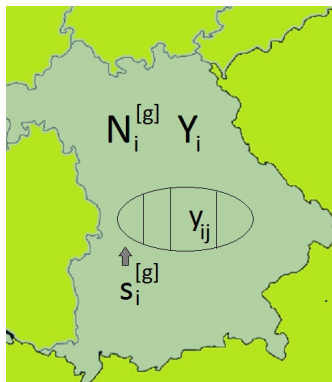
What's the problem? \Rightarrow 1. Small Area Estimation (SAE)



- Binary covariates (Abitur, sex)
- Cross classifications of the covariates
 \Rightarrow subgroup g , $g = 1, \dots, v$
- Known absolute frequencies $N_i^{[g]}$
Federal Statistical Office's data report:

		Abitur	
		no	yes
sex	male	$N_i^{[1]}$	$N_i^{[2]}$
	female	$N_i^{[3]}$	$N_i^{[4]}$

What's the problem? \Rightarrow 1. Small Area Estimation (SAE)

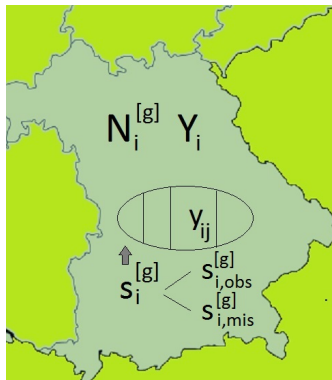


- Binary covariates (Abitur, sex)
- Cross classifications of the covariates
 \Rightarrow subgroup g , $g = 1, \dots, v$
- Known absolute frequencies $N_i^{[g]}$
Federal Statistical Office's data report:

		Abitur	
		no	yes
sex	male	$N_i^{[1]}$	$N_i^{[2]}$
	female	$N_i^{[3]}$	$N_i^{[4]}$

- Joint information about x_{ij} and y_{ij}
 \Rightarrow We know y_{ij} for $j \in s_i^{[g]}$

What's the problem? \Rightarrow 2. Missing data



- some sample values y_{ij} are missing
- $s_i^{[g]}$ is partitioned into $s_{i,obs}^{[g]}$ and $s_{i,mis}^{[g]}$

Cautious Approach for Dealing with Nonresponse

(ISIPTA '15, Plass, Augustin, Cattaneo, Schollmeyer)

- An observation model is determined by the missingness parameters $q_{na|y}^{[g]}$ ($:=$ probability to refuse the answer (“na”), given subgroup g and the true value y)

- Maximizing the log-likelihood

$$\begin{aligned} \ell(\pi^{[g]}, q_{na|0}^{[g]}, q_{na|1}^{[g]}) &= n_1^{[g]} \left(\ln(\pi^{[g]}) + \ln(1 - q_{na|1}^{[g]}) \right) \\ &+ n_0^{[g]} \left(\ln(1 - \pi^{[g]}) + \ln(1 - q_{na|0}^{[g]}) \right) + n_{na}^{[g]} \left(\ln(\pi^{[g]} q_{na|1}^{[g]} + (1 - \pi^{[g]}) q_{na|0}^{[g]}) \right) \end{aligned}$$

gives set-valued estimator.

- Resulting bounds of $\hat{\pi}^{[g]}$ under **no assumptions about $q_{na|y}^{[g]}$** :

$$\hat{\pi}^{[g]} = \frac{n_1^{[g]}}{n_{na}^{[g]} + n_1^{[g]} + n_0^{[g]}} \quad \text{and} \quad \bar{\pi}^{[g]} = \frac{n_1^{[g]} + n_{na}^{[g]}}{n_{na}^{[g]} + n_1^{[g]} + n_0^{[g]}}.$$

Cautious Approach for Dealing with Nonresponse

(ISIPTA '15, Plass, Augustin, Cattaneo, Schollmeyer)

- **Incorporate assumptions** by missingness ratio (Nordheim, 1984)

$$R = q_{na|1}^{[g]} / q_{na|0}^{[g]}, \quad \text{with } R \in \mathcal{R} \subseteq \mathbb{R}_0^+$$

- Specific values of R point-identify $\pi^{[g]}$
- Partial assumptions, expressed by $\mathcal{R} = [\underline{R}, \overline{R}]$, refine the result without any missingness assumptions ($R \in [0, 1]$)
 \Rightarrow Bounds for $\hat{\pi}^{[g], \mathcal{R}}$, $\hat{q}_{na|0}^{[g], \mathcal{R}}$ and $\hat{q}_{na|1}^{[g], \mathcal{R}}$ obtained under \underline{R} and \overline{R}

Illustration of the cautious approach (survey data only)

observation

sex	Abitur	poverty class		
		rich (0)	poor (1)	na
male (0)	no (0)	777	217	131
	yes (1)	492	61	81
female (1)	no (0)	668	284	144
	yes (1)	441	72	98

Illustration of the cautious approach (survey data only)

observation

sex	Abitur	poverty class		
		rich (0)	poor (1)	na
male (0)	no (0)	777	217	131
	yes (1)	492	61	81
female (1)	no (0)	668	284	144
	yes (1)	441	72	98

underlying truth

sex	Abitur	poverty class	
		0	1
0	0	777+□	217+□
	1	492+□	61+□
1	0	668+□	284+□
	1	441+□	72+□

no asmt. about $q_{naly}^{[g]}$

Illustration of the cautious approach (survey data only)

observation

		poverty class		
sex	Abitur	rich (0)	poor (1)	na
male (0)	no (0)	777	217	131
	yes (1)	492	61	81
female (1)	no (0)	668	284	144
	yes (1)	441	72	98

underlying truth

		poverty class	
sex	Abitur	0	1
0	0	777+131	217
	1	492+81	61
1	0	668+144	284
	1	441+98	72

no asmt. about $q_{naly}^{[g]}$

$$\hat{\pi}^{[1]} \approx 0.19$$

$$\hat{\pi}^{[2]} \approx 0.10$$

$$\hat{\pi}^{[3]} \approx 0.27$$

$$\hat{\pi}^{[4]} \approx 0.12$$

Illustration of the cautious approach (survey data only)

observation

		poverty class		
sex	Abitur	rich (0)	poor (1)	na
male (0)	no (0)	777	217	131
	yes (1)	492	61	81
female (1)	no (0)	668	284	144
	yes (1)	441	72	98

underlying truth

		poverty class	
sex	Abitur	0	1
0	0	777	217+131
	1	492	61+81
1	0	668	284+144
	1	441	72+98

no asmt. about $q_{naly}^{[g]}$

$$\hat{\pi}^{[1]} \approx 0.19, \quad \bar{\pi}^{[1]} \approx 0.31$$
$$\hat{\pi}^{[2]} \approx 0.10, \quad \bar{\pi}^{[2]} \approx 0.22$$
$$\hat{\pi}^{[3]} \approx 0.27, \quad \bar{\pi}^{[3]} \approx 0.37$$
$$\hat{\pi}^{[4]} \approx 0.12, \quad \bar{\pi}^{[4]} \approx 0.28$$

Illustration of the cautious approach (survey data only)

observation

		poverty class		
sex	Abitur	rich (0)	poor (1)	na
male (0)	no (0)	777	217	131
	yes (1)	492	61	81
female (1)	no (0)	668	284	144
	yes (1)	441	72	98

underlying truth

		poverty class	
sex	Abitur	0	1
0	0	777+102	217+29
	1	492+72	61+9
1	0	668+101	284+43
	1	441+84	72+14

no asmt. about $q_{naly}^{[g]}$

$$\hat{\pi}^{[1]} \approx 0.19, \quad \bar{\pi}^{[1]} \approx 0.31$$

$$\hat{\pi}^{[2]} \approx 0.10, \quad \bar{\pi}^{[2]} \approx 0.22$$

$$\hat{\pi}^{[3]} \approx 0.27, \quad \bar{\pi}^{[3]} \approx 0.37$$

$$\hat{\pi}^{[4]} \approx 0.12, \quad \bar{\pi}^{[4]} \approx 0.28$$

Specific value about R , here: $R = 1$ (MAR):

$$\hat{\pi}^{[1]} \approx 0.22, \quad \hat{\pi}^{[2]} \approx 0.11, \quad \hat{\pi}^{[3]} \approx 0.30, \quad \hat{\pi}^{[4]} \approx 0.14$$

Illustration of the cautious approach (survey data only)

observation

		poverty class		
sex	Abitur	rich (0)	poor (1)	na
male (0)	no (0)	777	217	131
	yes (1)	492	61	81
female (1)	no (0)	668	284	144
	yes (1)	441	72	98

underlying truth

		poverty class	
sex	Abitur	0	1
0	0	777+□	217+□
	1	492+□	61+□
1	0	668+□	284+□
	1	441+□	72+□

weak info $R \in [0, 1]$

$$\hat{\pi}^{[1]} \approx 0.19, \quad \bar{\pi}^{[1]} \approx 0.22$$

$$\hat{\pi}^{[2]} \approx 0.10, \quad \bar{\pi}^{[2]} \approx 0.11$$

$$\hat{\pi}^{[3]} \approx 0.27, \quad \bar{\pi}^{[3]} \approx 0.30$$

$$\hat{\pi}^{[4]} \approx 0.12, \quad \bar{\pi}^{[4]} \approx 0.14$$

Specific value about R , here: $R = 1$ (MAR):

$$\hat{\pi}^{[1]} \approx 0.22, \quad \hat{\pi}^{[2]} \approx 0.11, \quad \hat{\pi}^{[3]} \approx 0.30, \quad \hat{\pi}^{[4]} \approx 0.14$$

The synthetic estimator (without nonresponse)

- Horvitz-Thompson (HT) estimator
(Horvitz and Thompson, 1952, JASA)

$$\hat{\pi}_{HT,i} = \frac{1}{N_i} \sum_{j \in s_i} w_{ij} y_{ij}$$

- The synthetic estimator (González, 1973, JASA)

$$\hat{\pi}_{SYN} \equiv \hat{\pi}_{SYN,i} = \frac{1}{N} \sum_{i=1}^M \sum_{j \in s_i} w_{ij} y_{ij} = \frac{1}{N} \sum_{i=1}^M N_i \cdot \hat{\pi}_{HT,i}$$

- **No assumptions:**

$$\hat{\pi}_{SYN} = \frac{1}{N} \sum_{i=1}^M \left(\sum_{j \in \mathcal{S}_{i,obs}} w_{ij} y_{ij} + \sum_{j \in \mathcal{S}_{i,mis}} w_{ij} \cdot y_{ij} \right)$$

$$\hat{\pi}_{SYN} = \dots \left(\dots + \sum_{j \in \mathcal{S}_{i,mis}} w_{ij} \cdot 0 \right), \quad \bar{\pi}_{SYN} = \dots \left(\dots + \sum_{j \in \mathcal{S}_{i,mis}} w_{ij} \cdot 1 \right)$$

- **No assumptions:**

$$\hat{\pi}_{SYN} = \frac{1}{N} \sum_{i=1}^M \left(\sum_{j \in \mathcal{S}_{i,obs}} w_{ij} y_{ij} + \sum_{j \in \mathcal{S}_{i,mis}} w_{ij} \cdot y_{ij} \right)$$

$$\hat{\underline{\pi}}_{SYN} = \dots \left(\dots + \sum_{j \in \mathcal{S}_{i,mis}} w_{ij} \cdot 0 \right), \quad \hat{\overline{\pi}}_{SYN} = \dots \left(\dots + \sum_{j \in \mathcal{S}_{i,mis}} w_{ij} \cdot 1 \right)$$

- **Partial assumptions:**

$$\hat{\underline{\pi}}_{SYN}^{\mathcal{R}} = \frac{1}{N} \sum_{i=1}^M \left(\sum_{j \in \mathcal{S}_{i,obs}} w_{ij} y_{ij} + \underbrace{\hat{q}_{na|1i}^{\mathcal{R}} \cdot \hat{\underline{\pi}}_i^{\mathcal{R}} \cdot \sum_{j \in \mathcal{S}_i} w_{ij}} \right)$$

smallest est. weighted number of nonrespondents
with $y_{ij} = 1$, under the assumption in focus.

Analogously, $\hat{\overline{\pi}}_{SYN}^{\mathcal{R}}$ is achieved by using $\hat{q}_{na|1i}^{\mathcal{R}}$ and $\hat{\overline{\pi}}_i^{\mathcal{R}}$.

The LGREG estimator (without nonresponse)...

(Lehtonen and Veijanen, 1998, Surv. Methodol.)

- ... in its representation how we need it:

$$\hat{\pi}_{LGREG,i} = \sum_{g=1}^v \left(\overbrace{\sum_{j \in s_i^{[g]}} w_{ij} y_{ij}}^{\text{HT-part}} + \overbrace{\hat{\pi}^{[g]} (N_i^{[g]} - \sum_{j \in s_i^{[g]}} w_{ij})}^{\text{correction term}} \right) / N_i$$

with $\hat{\pi}^{[g]} = \sum_{i=1}^M \sum_{j \in s_i^{[g]}} \frac{y_{ij}}{n^{[g]}}$

- The **correction term** accounts for under/overrepresentation of certain constellations of covariates in the sample
- In most cases: $w_{ij} = w_i, \forall j = 1, \dots, n_i, i = 1, \dots, M$.

No assumptions: Cautious LGREG estimator

Breaking the summation over all areas into a term for area i^* of interest and areas $i \neq i^*$ leads to

$$\sum_{g=1}^v \left(\left(\frac{1}{n^{[g]}} \sum_{\substack{i=1 \\ i \neq i^*}}^M \left(\sum_{j \in S_{i,obs}^{[g]}} y_{ij} + \sum_{j \in S_{i,mis}^{[g]}} y_{ij} \right) \right) (N_{i^*}^{[g]} - n_{i^*}^{[g]} w_{i^*}) \right. \\ \left. + \frac{1}{n^{[g]}} \left(\sum_{j \in S_{i^*,obs}^{[g]}} y_{i^*j} + \sum_{j \in S_{i^*,mis}^{[g]}} y_{i^*j} \right) (N_{i^*}^{[g]} - w_{i^*} (n_{i^*}^{[g]} + n^{[g]})) \right) / N_{i^*}$$

No assumptions: Cautious LGREG estimator

Breaking the summation over all areas into a term for area i^* of interest and areas $i \neq i^*$ leads to

$$\sum_{g=1}^v \left(\left(\frac{1}{n^{[g]}} \sum_{\substack{i=1 \\ i \neq i^*}}^M \left(\sum_{j \in S_{i,obs}^{[g]}} y_{ij} + \sum_{j \in S_{i,mis}^{[g]}} y_{ij} \right) \right) (N_{i^*}^{[g]} - n_{i^*}^{[g]} w_{i^*}) \right. \\ \left. + \frac{1}{n^{[g]}} \left(\sum_{j \in S_{i^*,obs}^{[g]}} y_{i^*j} + \sum_{j \in S_{i^*,mis}^{[g]}} y_{i^*j} \right) (N_{i^*}^{[g]} - w_{i^*} (n_{i^*}^{[g]} + n^{[g]})) \right) / N_{i^*}$$

To determine $\hat{\pi}_{LGREG,i^*}$:

$N_{i^*}^{[g]} \geq w_{i^*} (n_{i^*}^{[g]} + n^{[g]})$		$N_{i^*}^{[g]} < w_{i^*} (n_{i^*}^{[g]} + n^{[g]})$
$N_{i^*}^{[g]} \geq n_{i^*}^{[g]} w_{i^*}$	$y_{ij} = 0, \forall j \in S_{i,mis}$	$y_{ij} = \begin{cases} 0 & \forall j \in S_{i,mis}, i \neq i^* \\ 1 & \forall j \in S_{i,mis}, i = i^* \end{cases}$
$N_{i^*}^{[g]} < n_{i^*}^{[g]} w_{i^*}$	$y_{ij} = \begin{cases} 1 & \forall j \in S_{i,mis}, i \neq i^* \\ 0 & \forall j \in S_{i,mis}, i = i^* \end{cases}$	$y_{ij} = 1, \forall j \in S_{i,mis}$

- 1.) Regard $\hat{\pi}_{LGREG, i^*}$ as a combination of two estimators:
 - \Rightarrow a global one that borrows strength and
 - \Rightarrow a specific one associated to area i^* .
- 2.) Maximize the two log-likelihoods under \underline{R} and \overline{R} :
 - $\ell(\pi^{[g], \mathcal{R}}, q_{na|0}^{[g], \mathcal{R}}, q_{na|1}^{[g], \mathcal{R}})$ and
 - $\ell(\pi_{i^*}^{[g], \mathcal{R}}, q_{na|0i^*}^{[g], \mathcal{R}}, q_{na|1i^*}^{[g], \mathcal{R}})$
- 3.) Include the estimators that minimize

$$\sum_{g=1}^v \left(\overbrace{\sum_{j \in S_{i^*, obs}^{[g]}} w_{i^*} y_{i^*j} + \hat{q}_{na|1i^*}^{[g], \mathcal{R}} \hat{\pi}_{i^*}^{[g], \mathcal{R}} \sum_{j \in S_{i^*}^{[g]}} w_{i^*j}}^{\text{HT-part}} + \overbrace{\hat{\pi}_{i^*}^{[g], \mathcal{R}} (N_{i^*}^{[g]} - n_{i^*}^{[g]} w_{i^*})}^{\text{correction term}} \right) / N_{i^*}$$

\Rightarrow Since $\pi^{[g]}$ and $\pi_{i^*}^{[g]}$ are estimated distinctively, interrelation between them should be considered.

Some results (example)

- Intervals for the synthetic estimator

no assumption	$\mathcal{R} = [0, 1]$
[0.167, 0.300]	[0.167, 0.193]

- Intervals for the LGREG estimator

Federal state	no assumption	$\mathcal{R} = [0, 1]$
BW	[0.129, 0.366]	[0.129, 0.210]
BY	[0.088, 0.233]	[0.088, 0.133]
HB	[0.077, 0.405]	[0.115, 0.193]
...

Next steps:

- Optimization of one overall likelihood, instead of two, to obtain the cautious LGREG-estimator
- Comparison of the magnitude of both principally differing kinds of uncertainty induced by the two problems in focus

Conclusion:

	no assumptions	weak auxiliary information
synthetic estimator	✓	✓
LGREG estimator	✓	~

Furthermore (not in this talk):

First sensitivity analysis for model-based SAE estimators