

# Select synthetic microdata for survey data

Hang J. Kim

Division of Statistics and Data Science  
University of Cincinnati

presented at 2024 Joint Statistical Meeting, Portland, OR  
August 6, 2024

## Select data synthesis approach

“Synthesize a **select set of variables** and **select records** with high disclosure risks” [ from the session abstract ]

1. Fully synthetic data (Rahunathan et al. 2003)

2. Partially synthetic data (Reiter 2003)

- ▶ Often relying on a (Bayesian) conditional probability model to generate the **select** set of variables/records  $Y$  in the original data  $(X, Y)$ :

$$\tilde{Y}_{\text{synt}} \sim f(\tilde{Y}|X, Y) = \int f(\tilde{Y}|X, \theta)f(\theta|X, Y)d\theta \propto \int f(\tilde{Y}|X, \theta)f(Y|X, \theta)f(\theta)d\theta$$

- ▶ Using the combining rule for correct uncertainty quantification
  - ▶ under both Bayesian and repeated sampling frameworks
- ⇒ Some deep-learning approaches (e.g., table-GAN) overlook the uncertainty quantification

## Select data synthesis approach

“Synthesize a **select set of variables** and **select records** with high disclosure risks” [ from the session abstract ]

1. Fully synthetic data (Rahunathan et al. 2003)
2. Partially synthetic data (Reiter 2003)
  - ▶ Often relying on a (Bayesian) conditional probability model to generate the **select** set of variables/records  $Y$  in the original data  $(X, Y)$ :

$$\tilde{Y}_{\text{synt}} \sim f(\tilde{Y}|X, Y) = \int f(\tilde{Y}|X, \theta)f(\theta|X, Y)d\theta \propto \int f(\tilde{Y}|X, \theta)f(Y|X, \theta)f(\theta)d\theta$$

- ▶ Using the combining rule for correct uncertainty quantification
  - ▶ under both Bayesian and repeated sampling frameworks
- ⇒ Some deep-learning approaches (e.g., table-GAN) overlook the uncertainty quantification

## Select data synthesis approach

“Synthesize a **select set of variables** and **select records** with high disclosure risks” [ from the session abstract ]

1. Fully synthetic data (Rahunathan et al. 2003)
2. Partially synthetic data (Reiter 2003)
  - ▶ Often relying on a (Bayesian) conditional probability model to generate the **select** set of variables/records  $Y$  in the original data  $(X, Y)$ :

$$\tilde{Y}_{\text{synt}} \sim f(\tilde{Y}|X, Y) = \int f(\tilde{Y}|X, \theta) f(\theta|X, Y) d\theta \propto \int f(\tilde{Y}|X, \theta) f(Y|X, \theta) f(\theta) d\theta$$

- ▶ Using the combining rule for correct uncertainty quantification
  - ▶ under both Bayesian and repeated sampling frameworks
- ⇒ Some deep-learning approaches (e.g., table-GAN) overlook the uncertainty quantification

## Select data synthesis approach

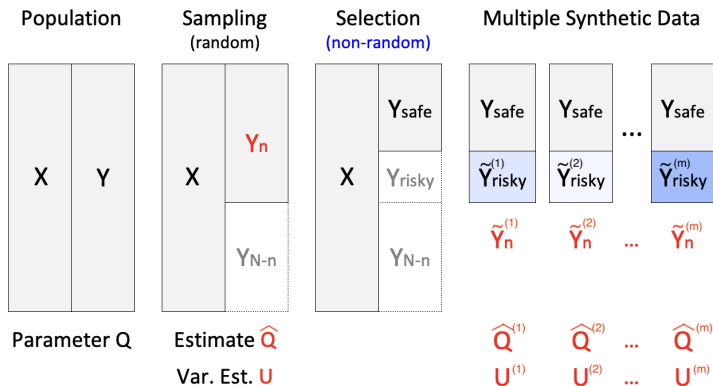
“Synthesize a **select set of variables** and **select records** with high disclosure risks” [ from the session abstract ]

1. Fully synthetic data (Rahunathan et al. 2003)
2. Partially synthetic data (Reiter 2003)
  - ▶ Often relying on a (Bayesian) conditional probability model to generate the **select** set of variables/records  $Y$  in the original data  $(X, Y)$ :

$$\tilde{Y}_{\text{synt}} \sim f(\tilde{Y}|X, Y) = \int f(\tilde{Y}|X, \theta) f(\theta|X, Y) d\theta \propto \int f(\tilde{Y}|X, \theta) f(Y|X, \theta) f(\theta) d\theta$$

- ▶ Using the combining rule for correct uncertainty quantification
  - ▶ under both Bayesian and repeated sampling frameworks
- ⇒ Some deep-learning approaches (e.g., table-GAN) overlook the uncertainty quantification

► Combining rule for uncertainty quantification (under repeated sampling)



$$\frac{\sqrt{n}(\hat{Q} - Q)}{\sqrt{U}} \rightarrow N(0, 1)$$

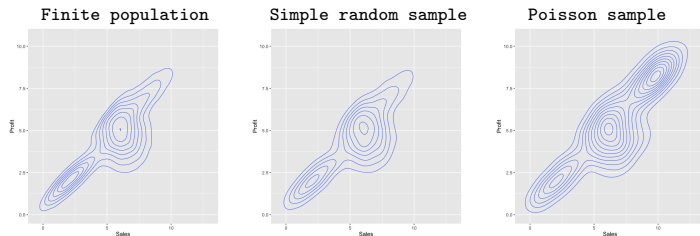
$$\frac{\sqrt{n}(\bar{Q} - Q)}{\sqrt{T}} \rightarrow N(0, 1)$$

$$\bar{Q} = \text{Avg. of } \hat{Q}^{(l)}$$

$$T = \text{Avg. of } U^{(l)} + \frac{1}{m} \left( \text{Var. of } \hat{Q}^{(l)} \right)$$

# Synthesis for survey sampling data

- ▶ The sampling distribution of survey data often differs with the distribution of finite population.

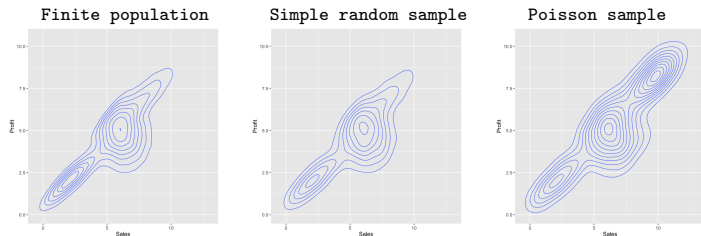


- ▶ Some model-based approaches with survey weights

1. Disregarding the survey weights,  $\prod_{i=1}^n f(y_i|\theta) = ?$
2. Reconstruct the finite pop. (bootstrap),  $f(y_1, \dots, y_N|\theta) = \prod_{i=1}^N f(\tilde{y}_i|\theta)$ ,  
where  $\tilde{y}_i = y_i$  for sampled units and other  $\tilde{y}_i$  are estimated/resampled.
3. Using the pseudo likelihood,  $f(y_1, \dots, y_N|\theta) \approx \prod_{i=1}^n f(y_i|\theta)^{w_i}$ .

# Synthesis for survey sampling data

- ▶ The sampling distribution of survey data often differs with the distribution of finite population.

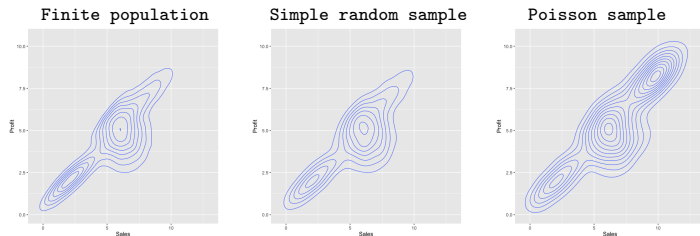


- ▶ Some model-based approaches with survey weights

1. Disregarding the survey weights,  $\prod_{i=1}^n f(y_i|\theta) = ?$
2. Reconstruct the finite pop. (bootstrap),  $f(y_1, \dots, y_N|\theta) = \prod_{i=1}^N f(\tilde{y}_i|\theta)$ ,  
where  $\tilde{y}_i = y_i$  for sampled units and other  $\tilde{y}_i$  are estimated/resampled.
3. Using the pseudo likelihood,  $f(y_1, \dots, y_N|\theta) \approx \prod_{i=1}^n f(y_i|\theta)^{w_i}$ .

# Synthesis for survey sampling data

- ▶ The sampling distribution of survey data often differs with the distribution of finite population.

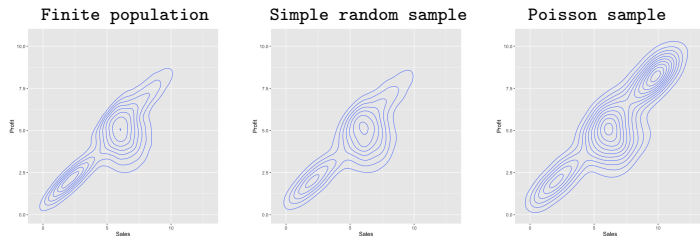


- ▶ Some model-based approaches with survey weights

1. Disregarding the survey weights,  $\prod_{i=1}^n f(y_i|\theta) = ?$
2. Reconstruct the finite pop. (bootstrap),  $f(y_1, \dots, y_N|\theta) = \prod_{i=1}^N f(\tilde{y}_i|\theta)$ ,  
where  $\tilde{y}_i = y_i$  for sampled units and other  $\tilde{y}_i$  are estimated/resampled.
3. Using the pseudo likelihood,  $f(y_1, \dots, y_N|\theta) \approx \prod_{i=1}^n f(y_i|\theta)^{w_i}$ .

# Synthesis for survey sampling data

- ▶ The sampling distribution of survey data often differs with the distribution of finite population.



- ▶ Some model-based approaches with survey weights

1. Disregarding the survey weights,  $\prod_{i=1}^n f(y_i|\theta) = ?$
2. Reconstruct the finite pop. (bootstrap),  $f(y_1, \dots, y_N|\theta) = \prod_{i=1}^N f(\tilde{y}_i|\theta)$ ,  
where  $\tilde{y}_i = y_i$  for sampled units and other  $\tilde{y}_i$  are estimated/resampled.
3. Using the pseudo likelihood,  $f(y_1, \dots, y_N|\theta) \approx \prod_{i=1}^n f(y_i|\theta)^{w_i}$ .

## Bayesian pseudo posterior approach (Savitsky and Toth, 2016)

Assuming that  $(w_i - 1)$  non-sampled units have the same values as a sampled unit  $y_i$  in evaluating the (pseudo) likelihood fn.  $l^{\text{pse}}(\theta) = \prod_{i=1}^n f(y_i|\theta)^{w_i}$ ,

$$f^{\text{pse}}(\theta|\mathbf{y}_n, \mathbf{w}_n) = f^{\text{pse}}(\theta|y_1, \dots, y_n, w_1, \dots, w_n) \propto \prod_{i=1}^n f(y_i|\theta)^{w_i} \cdot f(\theta)$$

We proved that

1.  $E(\theta|\text{Data})$  with  $f^{\text{pse}}$  is asymptotically unbiased. [Bernstein–Von Mises]

$$(n\mathbf{Q}_0)^{1/2} f^{\text{pse}}(\theta|\mathbf{y}) \rightarrow \mathcal{N}(\theta_0, \mathbf{I}) \text{ as } n \rightarrow \infty \text{ where } \mathbf{Q}_0 = -E_0 [\nabla^2 l^{\text{pse}}(\theta)]$$

2. Posterior variance of  $\theta$  is not close to the variance of the posterior mean for repeated sampling, i.e.,  $E(\hat{V}(\theta|\text{Data})) \neq V(\hat{E}(\theta|\text{Data}))$  [Godambe information]

$$(n\mathbf{Q}_0\mathbf{P}_0^{-1}\mathbf{Q}_0)^{1/2} (\hat{\theta}_n^{\text{pse}} - \theta_0) \rightarrow \mathcal{N}(\mathbf{0}, \mathbf{I}) \text{ where } \mathbf{P}_0 = E_0 [\nabla l^{\text{pse}}(\theta)\nabla l^{\text{pse}}(\theta)^{\top}]$$

## Suggestion: Variance-adjusted pseudo posterior

With the original pseudo posterior approach,  $Q_0 \neq P_0$ .

We suggest to use the power of the adjusted weights  $\kappa w_i$ ,

$$f^{\text{adj}}(\theta|y_1, \dots, y_n, w_1, \dots, w_n) \propto \prod_{i=1}^n f(y_i|\theta)^{\kappa w_i} \cdot f(\theta) \quad \text{where } \kappa = \frac{\sum_{j=1}^n w_j}{\sum_{j=1}^n w_j^2}.$$

1. With the adjusted weights,  $Q_0 = P_0$ , so the posterior mean with the adjusted pseudo likelihood follows

$$\sqrt{n} \left( \hat{\theta}_n^{\text{adj}} - \theta_0 \right) \rightarrow \mathcal{N}(\mathbf{0}, \mathbf{I})$$

2. In SRS, the adjusted pseudo posterior becomes the posterior distribution disregarding the survey weights, i.e.,

$$\kappa w_i = \frac{\sum_{j=1}^n \frac{N}{n}}{\sum_{j=1}^n \frac{N^2}{n^2}} \frac{N}{n} = 1 \quad \Rightarrow \quad \prod_{i=1}^n f(y_i|\theta)^{\kappa w_i} \cdot f(\theta) = \prod_{i=1}^n f(y_i|\theta) \cdot f(\theta)$$

# Simulation study result summary

Three synthetic data methods in comparison

1. **No weight**, ignoring survey weights,  $\prod_{i=1}^n f(y_i|\theta) \cdot f(\theta)$ .
2. **Pseudo** posterior with the original survey weights,  $\prod_{i=1}^n f(y_i|\theta)^{w_i} \cdot f(\theta)$ .
3. **Adjusted** pseudo posterior,  $\prod_{i=1}^n f(y_i|\theta)^{\kappa w_i} \cdot f(\theta)$ .

Sampling Methods		No weight	Pseudo	Adjusted
Simple Random Sampling	$E(\hat{Y}_1) - \bar{Y}_1$	0.00	0.00	0.00
	$V(\hat{Y}_1)$	0.027	0.027	0.028
	$E(\hat{V}(\hat{Y}_1))$	0.025	0.001	0.025
	95% <i>C.I</i> coverage	0.928	0.286	0.922
Poisson Sampling	$E(\hat{Y}_1) - \bar{Y}_1$	2.02	0.00	0.00
	$V(\hat{Y}_1)$	0.030	0.031	0.031
	$E(\hat{V}(\hat{Y}_1))$	0.025	0.001	0.027
	95% <i>C.I</i> coverage	0.000	0.298	0.924

## Concluding remarks

1. Disregarding sampling weights results in biased estimation when the sample is collected with unequal probability sampling.
  2. The (original) pseudo posterior approach results in variance underestimation.
  3. The suggested pseudo likelihood approach with **the adjusted weight** results in correct estimation with imputed (and synthetic) data.
- Note: William and Savitsky (2021) suggested a post-processing **after** MCMC.
- In case of needing imputation, the adjustment needs to be given **during** MCMC as suggested here.

# Thank you!

## Contact Information

**Hang Kim (hang.kim@uc.edu)**

Division of Statistics and Data Science  
Department of Mathematical Sciences  
University of Cincinnati