

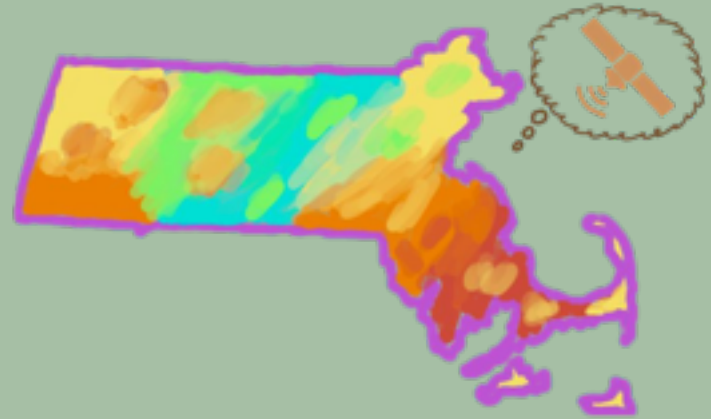
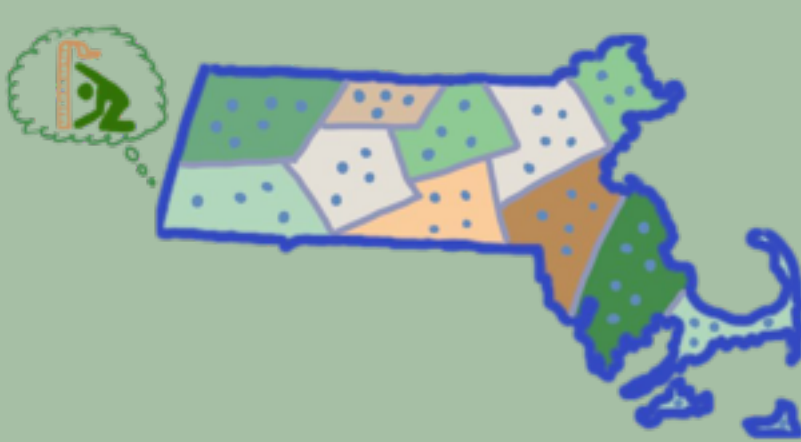


# Comparison Study of Survey Sampling Estimators

Professor Kelly McConville (Harvard)

Asteria Chilambo (Harvard Math '23),  
Jing Shang (Fudan Economics '23,  
Visiting Student at Harvard Statistics '22)

# How does typical forestry data look like?



## Ground Plot:

- Directly measures variables of interest
- Precise, but expensive and sparse!

## Remote sensor census:

- Indirectly provides related information
- Cost-effective, but not exactly what we want

**How can we use auxiliary variables  
to assist estimation for  
variables of interest?**



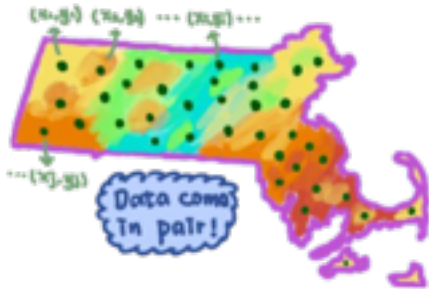
**Generalized Multivariable  
Difference Estimator (GMDE)**

**&**

**Generalized Regression  
Estimator (GREG)**

# GMDE

Step 1. Make use of Horvitz-Thompson  $\hat{t}_{y\pi}$  and auxiliary residual  $t_x - \hat{t}_{x,\pi}$



Step 2. Among all linear combination of  $\hat{t}_{y\pi}$  and  $t_x - \hat{t}_{x,\pi}$ , choose one with minimal variance

# GREG

Step 1. Start with linear regression

Step 2. Optimize regression coefficient

Step 3. Use **population** information of  $x$  to predict  $y$  **through regression**, plus **sample** information of residual

**How do GMDE and GREG perform  
under different sampling scenarios?**

# What is dominating the difference?



- **GMDE**

$$\hat{t}_{y,gmde} = \hat{t}_{y\pi} + \hat{V}_{yx,\pi} \cdot \hat{V}_{x,\pi}^{-1} \cdot (t_x - \hat{t}_{x,\pi})$$

- **GREG**

$$\hat{t}_{y,greg} = \hat{t}_{y\pi} + \left( \sum_{i \in S} \frac{x_i y_i}{\pi_i} \right)^\top \left( \sum_{j \in S} \frac{x_j x_j^\top}{\pi_j} \right)^{-1} \cdot (t_x - \hat{t}_{x\pi})$$

Pos. Corr.

If ppt. detected from sensor is bigger than est. from sample

# Simple Random Sampling

## Single Study Variable

- GMDE

$$\hat{V}_{y_{x,\pi}} \hat{V}_{x,\pi}^{-1} = \frac{\frac{1}{n} \sum_{i \in S} x_i y_i - \left( \frac{1}{n} \sum_{i \in S} y_i \right) \left( \frac{1}{n} \sum_{i \in S} x_i \right)}{\frac{1}{n} \sum_{i \in S} x_i^2 - \left( \frac{1}{n} \sum_{i \in S} x_i \right)^2}$$

Cov. of  $y$  and  $x$   
Cov. of  $x$

- GREG

$$\left( \sum_{j \in S} \frac{x_j x_j^T}{\pi_j} \right)^{-1} \left( \sum_{i \in S} \frac{x_i y_i}{\pi_i} \right) = \frac{\frac{1}{n} \sum_{i \in S} x_i y_i - \left( \frac{1}{n} \sum_{i \in S} y_i \right) \left( \frac{1}{n} \sum_{i \in S} x_i \right)}{\frac{1}{n} \sum_{i \in S} x_i^2 - \left( \frac{1}{n} \sum_{i \in S} x_i \right)^2}$$



**GMDE and GREG are the same under SRS!**

# Simple Random Sampling



**GMDE and GREG are the same under SRS!**

## Multiple Study Variables

- **GMDE**

$$(\hat{t}_{y,gmde})_m = (\hat{t}_{y\pi})_m + \underbrace{(\hat{V}_{yx,\pi})_m (\hat{V}_{x,\pi}^{-1})_m}_{\text{Cov. of } m^{th} y \text{ and all } x} (-\hat{t}_{x\pi} + t_x)$$

Cov. of all  $x$

- **GREG**

$$\hat{t}_{y,greg} = \hat{t}_{y\pi} + \underbrace{\left( \sum_{i \in S} \frac{x_i y_i}{\pi_i} \right)^\top \left( \sum_{j \in S} \frac{x_j x_j^\top}{\pi_j} \right)^{-1}}_{\text{Cov. of all } x} \cdot (t_x - \hat{t}_{x\pi})$$

Using GREG to estimate the  $m^{th} y$



# Stratified Simple Random Sampling

No simple form for two ratios, but have numerical similarities:



GMDE_DR	GREG_DR	GMDE	GREG	GMDE_Var	GREG_Var
0.601	0.591	1811.995	1814.122	0.895	0.923
0.599	0.584	1616.635	1616.157	0.721	0.815
0.595	0.580	1484.934	1483.678	1.151	1.267
0.600	0.579	1765.941	1765.211	1.146	1.178
0.608	0.595	1769.369	1766.972	1.188	1.233

Comparison of dominant ratio (also dominates var.)

Comparison of estimator and variance

## Intuition:

- When sample mean varies significantly across strata, then GMDE is more precise.
- Because  $V_{y|x,e} - V_{x,e}^{-1}$  is measured within each stratum, while reg coefficient is computed across strata. The den shrinks faster.

# Main takeaways

- A. Both GMDE and GREG make use of **auxiliary variables** and can **improve performance** compared with Horvitz-Thompson estimator
- B. Under **simple random sampling**, GMDE and GREG performs the same
- C. Under **stratified** simple random sampling, GMDE performs better

## Future work

- Take different estimators of variance into consideration (eg. Hajek-Berger, Hansen-Hurwitz)
- Extend sampling scenarios (eg. two phase sampling)

**GMDE deserves more attention!**

# Special Thanks to:

**Ray Czaplewski**

**Emeritus Scientist**

**Forest Inventory & Analysis**

**George Gaines**

**Research Mathematical Statistician**

**Forest Inventory & Analysis**

**Kelly McConville**

**Harvard University**

**Department of Statistics**