# Methods for Estimating Hard to Count Populations

Adrian E. Raftery
University of Washington
http://www.stat.washington.edu/raftery

NAS Committee on Best Practices in Assessing Mortality and Significant Morbidity Following Large-Scale Disasters Webinar

February 11, 2020

#### Outline

- Capture-recapture methods for combining multiple fragmentary data sources<sup>1</sup>
- ► Respondent-driven sampling<sup>2 3</sup>
- ► Network scale-up method<sup>4</sup> [to be presented by Tyler McCormick]

<sup>&</sup>lt;sup>1</sup>Bao, Raftery, Reddy (2015, *Statistics and its Interface*)

<sup>&</sup>lt;sup>2</sup>Baraff, McCormick, Raftery (2016, PNAS)

<sup>&</sup>lt;sup>3</sup>Green, Raftery, McCormick (2020, *Biometrika*)

<sup>&</sup>lt;sup>4</sup>Maltiel, Raftery, McCormick, Baraff (2015, Annals of Applied Statistics)

# Uncertainty About Hidden Population Size

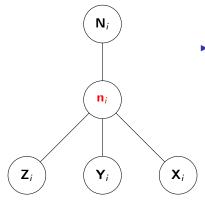
- Intrinsically hard to assess because of confounding between population size and probability of being counted:
  - It's hard to distinguish between low population size (close to the number seen) with high probability of being counted, and
  - high population size with low probability of being counted
  - ightharpoonup uncertainty about N can be large
  - Especially for estimation from one data source
  - Prior or external information can be useful, especially:
    - information about detection probabilities
    - a (possibly soft) upper bound on N when detection probabilities are small
- Bayesian approach:
  - can combine information of different types
  - can deal with missing data
  - allows one to use data from other periods, regions, countries to inform the prior distribution

# Capture-Recapture Estimation

Data Sources for Size Estimation of Intravenous Drug Users in Bangladesh, 2004

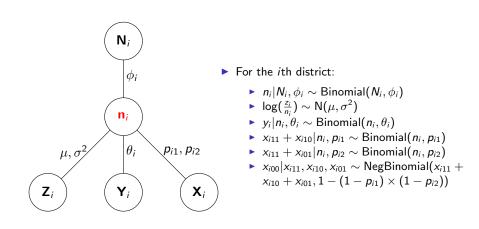
| Notation            | Data Source                    | Locations    | Year |
|---------------------|--------------------------------|--------------|------|
| $N_i$               | adult male population size     | 64 districts | 2003 |
| $Y_i$               | NASROB mapping data            | 24 districts | 2001 |
| $X_{i10} + X_{i11}$ | NEP intervention               | 3 districts  | 2001 |
| $X_{i01} + X_{i11}$ | BSS survey                     | 4 districts  | 2002 |
| $X_{i11}$           | reported NEP enrollment in BSS | 2 districts  | 2002 |
| $Z_i$               | CARE RSA estimate              | 14 districts | 2003 |

#### Data Structure

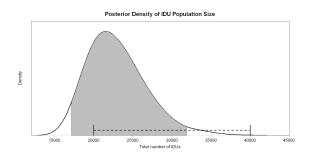


- ► For the *i*th district:
  - $ightharpoonup N_i = total population$
  - $n_i = number of IDUs$
  - $ightharpoonup Z_i = \text{estimate of } n_i \text{ from other sources}$
  - ► *Y<sub>i</sub>* = the number of IDUs in survey or intervention
  - ▶  $X_i = (X_{i10}, X_{i01}, X_{i11}) = \text{data from two}$  lists from multiplier method or capture-recapture

# Bayesian Hierarchical Model



# Population Size Estimate



- ▶ Posterior median 23,500
- ▶ Posterior interval  $18,000 \sim 33,700$
- ightharpoonup Bangladesh technical group estimate: 20,000  $\sim$  40,000

# Respondent-Driven Sampling

#### Introduction

#### Respondent-Driven Sampling (RDS):

- Problem no sampling frame for hard-to-reach or hidden populations (IDU, FSW, MSM, etc.)
- Developed by Heckathorn (1997, 2002) to add statistical rigor to snowball (chain-referral) sampling
  - known biases due to homophily and centrality
- Start with small number of initial respondents (seeds)
  - usually convenience sample
- Seeds given coupons to distribute to peers
- When coupons redeemed, respondents are given new coupons creating chain of recruitment
- Additional information gathered
  - coupon number who recruited whom
  - degree of each respondent how many could have recruited

#### **Estimation**

#### Volz-Heckathorn Estimator (2008):

ightharpoonup To estimate the mean  $\mu$  of a trait

$$\hat{\mu}_{VH} = \frac{\sum_{i=1}^{n} \frac{x_i}{d_i}}{\sum_{i=1}^{n} \frac{1}{d_i}},$$

where  $x_i$  = value of the trait in individual i and  $d_i$  = degree of individual i

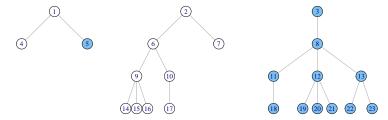
- Hansen-Hurwitz (H-H) type estimator with selection probabilities proportional to degree
- Unbiased if Markov chain model holds

**Problem:** How do we estimate  $Var(\hat{\mu}_{VH})$ ?

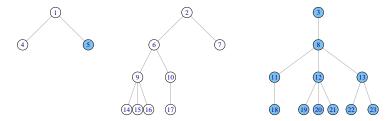
#### My idea:

- Ignore traits and focus on the tree structure of the RDS
- Multilevel bootstrap
  - 1. Resample seeds with replacement
  - 2. From each seed resample recruits with replacement
  - 3. From each recruit resample further recruits with replacement
  - 4. Repeat 3 until no recruits are available
- Assumes the parts of the network unseen "look" like the parts of the network seen
- Benefits
  - Correlation structure of entire tree preserved, not just first-order transitions
  - One bootstrap sample for all traits

#### Sample:

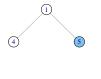


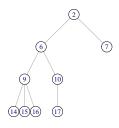
#### Sample:

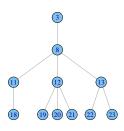




#### Sample:







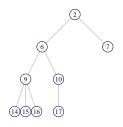


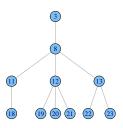




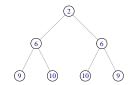
#### Sample:







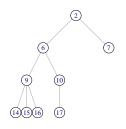


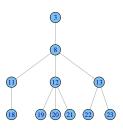




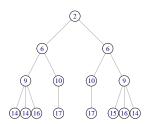
#### Sample:

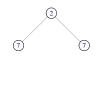












- Tree bootstrap performed better than competing methods in a simulation study<sup>5</sup>
- ► Tree bootstrap estimate of variance is asymptotically correct<sup>6</sup>

<sup>&</sup>lt;sup>5</sup>Baraff, McCormick, Raftery (2016, *PNAS*) <sup>6</sup>Green, McCormick, Raftery (2020, *Biometrika*)