



# PHI Applied Research Fellows 2021: Survey Statistics

Jessica Godwin

July 7, 2021

Data-generating Processes

Sampling Schemes

Design-based Estimation

Model-based estimation

References

# Data-generating Processes

# Data-generating Processes

- ▶ Statisticians often choose models and likelihoods based on a combination of:
  - ▶ how closely they reflect the true **data-generating process**
  - ▶ the mathematical and statistical properties
  - ▶ (hopefully) the principle of parsimony
- ▶ If our outcome is binary:

# Data-generating Processes

- ▶ Statisticians often choose models and likelihoods based on a combination of:
  - ▶ how closely they reflect the true **data-generating process**
  - ▶ the mathematical and statistical properties
  - ▶ (hopefully) the principle of parsimony
- ▶ If our outcome is binary:
  - ▶ Flipping a coin? or
  - ▶ Drawing from an urn?

# Data-generating Processes

- ▶ Statisticians often choose models and likelihoods based on a combination of:
  - ▶ how closely they reflect the true **data-generating process**
  - ▶ the mathematical and statistical properties
  - ▶ (hopefully) the principle of parsimony
- ▶ If our outcome is binary:
  - ▶ Flipping a coin? or
  - ▶ Drawing from an urn?
  - ▶ What does flipping a coin or drawing from an urn have to do with surveys with human respondents?

# Finite vs. superpopulations

For observations  $i = 1, \dots, n$ , let

$$Z = \sum_{i=1}^n y_i$$

$$E[Z] = \sum_{i=1}^n E[y_i]$$

$$Z \sim \text{Bin}(n, p)$$

$$y_i = \begin{cases} 1, & \text{success,} \\ 0, & \text{failure.} \end{cases}$$

$$E[y_i] = 0 \cdot P(y_i = 0) + 1 \cdot P(y_i = 1) = p$$

flipping coin

Superpopulation: If  $y \sim \text{Bernoulli}(p)$ ,

$$E[y] = p \quad \text{Var}(y) = p(1-p).$$

Finite population: If  $y \sim \text{Hypergeometric}(N, K, n)$ ,

$$E[y] = \frac{K}{N} \quad \text{Var}(y) = \frac{K}{N} \left(1 - \frac{K}{N}\right) \left(1 - \frac{n}{N}\right).$$

How do we say what  $\hat{p}$  means in either case? Is it the same?

drawing from an urn



## Finite vs. superpopulations, cont'd

If  $y_i \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$ ,

$$\hat{p} = \frac{\sum_{i=1}^n y_i}{n},$$

$$\widehat{\text{Var}}(\hat{p}) = \frac{\hat{p}(1 - \hat{p})}{n}.$$

$$\boxed{\text{Var}(\hat{p}) \rightarrow \text{SE}(\hat{p}) = \sqrt{\text{Var}(\hat{p})}}$$

If  $y_i \stackrel{\text{iid}}{\sim} \text{Hypergeometric}(N, K, n)$ ,

$$\hat{p} = \frac{\sum_{i=1}^n y_i}{n} = \frac{k}{n}$$

$$\widehat{\text{Var}}(\hat{p}) = \frac{\hat{p}(1 - \hat{p})}{n} \times \left(1 - \frac{n}{N}\right).$$

How do we say what  $\hat{p}$  means in either case? Is it the same?



# Sampling Schemes

# Simple random sampling(SRS)

- Under an **SRS** of  $n$  observations

$$\Pr(\text{subject } k \in \text{sample}, S) =$$

$$\boxed{\pi_k} = \frac{1}{N}$$

$$\Pr(\text{subjects } k, k' \in \text{sample}, S) =$$

$$\pi_{k,k'} = \frac{1}{N} \times \frac{1}{N}$$

$P(A \cap B)$

$\rightarrow \stackrel{\text{ind}}{=} P(A)P(B)$

$\rightarrow \pi_k \perp \pi_{k'}$

- Under an **SRSWOR** of  $n$  observation  
*w/out replacement*

$$\Pr(\text{subject } k \in \text{sample}, S) =$$

$$\pi_k = \frac{n}{N} = \pi_{k'}$$

$$\Pr(\text{subjects } k, k' \in \text{sample}, S) =$$

$$\pi_{k,k'} = \frac{n}{N} \times \frac{n-1}{N-1}$$

$P(A \cap B)$   
 $= P(A|B)P(B)$

$\rightarrow \pi_k \times \pi_{k'}$

# Systematic sampling

- ▶ Select every  $r^{\text{th}}$  sampling unit from the sampling frame of length

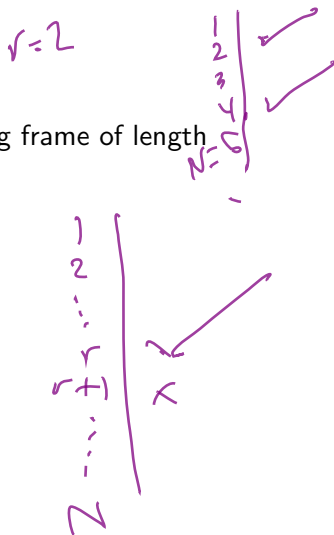
$$N: r \times n \leq N < r \times (n+1)$$

- ▶ What is  $\pi_k$  for individual  $k = r$ ?  $k = r+1$ ?

$$E[X] = 0$$

$$\pi_r = 1$$

$$\pi_{r,v+1} = 0 \quad \pi_1 = \dots = \pi_{r+1} = 0$$



# Systematic sampling

- ▶ Select every  $r^{th}$  sampling unit from the sampling frame of length  $N$ :  $r \times n \leq N < r \times (n + 1)$ 
  - ▶ What is  $\pi_k$  for individual  $k = r$ ?  $k = r + 1$ ?
  - ▶ Can a systematic sample be implemented so that it is the equivalent of an SRS?

# Systematic sampling

- ▶ Select every  $r^{th}$  sampling unit from the sampling frame of length  $N$ :  $r \times n \leq N < r \times (n + 1)$ 
  - ▶ What is  $\pi_k$  for individual  $k = r$ ?  $k = r + 1$ ?
  - ▶ Can a systematic sample be implemented so that it is the equivalent of an SRS?
  - ▶ What is  $\pi_{r,r+1}$ ?

# Systematic sampling

- ▶ Select every  $r^{\text{th}}$  sampling unit from the sampling frame of length  $N$ :  $r \times n \leq N < r \times (n + 1)$ 
  - ▶ What is  $\pi_k$  for individual  $k = r$ ?  $k = r + 1$ ?
  - ▶ Can a systematic sample be implemented so that it is the equivalent of an SRS?
  - ▶ What is  $\pi_{r,r+1}$ ?
- ▶ Random single start  $\rightarrow$  what changes?

$$1 \leq v \leq 10$$

$$\pi_r = \frac{1}{10}$$

$$\pi_{r,r+1}$$

$$\begin{array}{c} 1 \\ 2 \\ 3 \\ \vdots \\ r \\ r+1 \end{array} \Bigg|$$

# Systematic sampling

- ▶ Select every  $r^{\text{th}}$  sampling unit from the sampling frame of length  $N$ :  $r \times n \leq N < r \times (n + 1)$ 
  - ▶ What is  $\pi_k$  for individual  $k = r$ ?  $k = r + 1$ ?
  - ▶ Can a systematic sample be implemented so that it is the equivalent of an SRS?
  - ▶ What is  $\pi_{r,r+1}$ ?
- ▶ Random single start  $\rightarrow$  what changes?
- ▶ Multiple starts
  - ▶ No individual sampling probabilities are 0 or 1
  - ▶ Joint sampling probabilities defined

$$1 \leq r_1, r_2 \leq 10$$

# Stratified simple random sampling (strSRS)

- Consider  $h = 1, \dots, H$  strata from each of which you want to sample  $n_h$  individuals.

$$N = N_1 + \dots + N_H$$

$$\Pr(\text{subject } k \in S_h) = \pi_k = \frac{n_h}{N_h}$$

$$\Pr(\text{subjects } \underline{k, k'} \in S_h) = \pi_{k,k'} = \frac{n_h}{N_h} \times \frac{n_h - 1}{N_h - 1}$$

$$\Pr(\text{subjects } \underline{k \in S_h, k' \in S_{h'}}) = \pi_{k,k'} = \frac{n_h}{N_h} \times \frac{n_{h'}}{N_{h'}}$$

$$\pi_k \quad \pi_{k'}$$

Answer  
 $\pi_k \pi_{k'} | k$



## strSRS, cont'd

- ▶ Why stratify? Why not an SRS or SRSWOR?

## strSRS, cont'd

- Why stratify? Why not an SRS or SRSWOR?

- Availability of **sampling frame**

- Cost, convenience, speed

- $N_1, \dots, N_h$  vary widely

- Rare outcomes within certain strata

- We know strata are related to outcome of interest → precision gains!

- What happens if we ignore the stratification?

- Waste a lot of folks' money!!

- Implicit assumption that outcome of interest doesn't differ by strata

- → obscure differences in outcomes by strata

- → OVERESTIMATE variance/standard errors

- → worsens variability in outcomes between strata grows and within strata shrinks

- → worsens as variability in  $\pi_{k \in S_h}$  between strata grows

$$S^2 = \sum_{i=1}^n (x_i - \bar{x})^2$$

## Cluster sampling

Consider sampling  $c = 1, \dots, C$  clusters or **primary sampling units (PSU)** from your population of  $N_C$  clusters and  $N$  **units**.

Individuals  $k$  are the **observation units** contained within clusters on which we will make measurements.

### One-stage cluster sampling

$$\Pr(\text{PSU } c \in S) = \frac{C}{N_C}$$
$$\pi_{k \in S_c} = \begin{cases} 1, & \text{PSU } c \in S, \\ 0, & \text{otherwise.} \end{cases}$$

### Two-stage cluster sampling

Sample  $m_c$  from  $M_c$  units in cluster  $c$ .

$$\Pr(\text{PSU } c \in S) = \frac{C}{N_C}$$
$$\pi_{k \in S_c} = \begin{cases} \frac{m_c}{M_c}, & \text{PSU } c \in S, \\ 0, & \text{otherwise.} \end{cases}$$

## Cluster sampling, cont'd

- ▶ Probability proportional to size (PPS) sampling
  - ▶  $\pi_c \propto M_c$
  - ▶ When does this make sense?
- ▶ Why implement a cluster sample?
  - ▶ The only sampling frame we have is a list of groups of observation units
  - ▶ Cost and convenience
- ▶ What happens if we ignore clustering in our sample?
  - ▶ The  $m_c$  observation units sampled in cluster  $c$  are **not** independent samples
  - ▶ → we have LESS information than  $m_c$  observations from an SRS
  - ▶ → we will UNDERESTIMATE variances and standard errors if we ignore this dependence
  - ▶ → this underestimation worsens as the correlation between outcomes from individuals in a cluster increases

$$n = \sum_c m_c$$

# Complex surveys

## Multi-stage sampling

- ▶ **Example:** DHS (among others) stratify clusters by administrative divisions  $\times$  urban/rural  $\rightarrow$  select women within households within clusters within strata
- ▶ **Stratified two-stage cluster sampling**
- ▶ PSUs  $\rightarrow$  **secondary sampling units (SSUs)**  $\rightarrow$  observation units
- ▶ One could stratify within clusters if a sampling frame necessitates (never encountered this yet)

## Multi-phase sampling

- ▶ Fancy term for trying again to reach non-respondents!!
- ▶ Sub-sample (perhaps fully) your nonrespondents in attempts to get a response.

# Design-based Estimation

# Horvitz-Thompson estimators

- ▶ Each individual  $k$  has their responses weighted by their **sampling weight**  $w_k = \frac{1}{\pi_k}$ 
  - ▶ i.e. an individual with low chance of being sampled  $\rightarrow \pi_k$  small  $\rightarrow w_k$  big
  - ▶  $w_k$  can be interpreted as number of individuals in the finite population that individual  $k$ 's response represents
  - ▶ **Caveat:** nonresponse
- ▶ **Average or arithmetic mean**

$$\begin{aligned}
 \frac{\sum_{k=1}^n y_k}{n} &\stackrel{?}{=} \frac{\sum_{k=1}^n w_k y_k}{\sum_{k=1}^n w_k} = \frac{\sum_{k=1}^n \frac{N}{n} y_k}{\sum_{k=1}^n \frac{N}{n}} \\
 &= \frac{\frac{N}{n} \sum_{k=1}^n y_k}{\frac{N}{n} \sum_{k=1}^n 1} = \frac{N}{n} \left( \frac{\sum_{k=1}^n y_k}{\frac{N}{n} \times n} \right) = \frac{N}{n} \left( \frac{\sum_{k=1}^n y_k}{N} \right) = \frac{\sum_{k=1}^n y_k}{n}
 \end{aligned}$$

Handwritten notes:  $\frac{1}{\pi_k}$  and  $\sum w_k y_k$  are written in purple above the equation. The final result  $\frac{\sum_{k=1}^n y_k}{n}$  is circled in purple.

## Horvitz-Thompson estimators

- ▶ Each individual  $k$  has their responses weighted by their **sampling weight**  $w_k = \frac{1}{\pi_k}$ 
  - ▶ i.e. an individual with low chance of being sampled  $\rightarrow \pi_k$  small  $\rightarrow w_k$  big
  - ▶  $w_k$  can be interpreted as number of individuals in the finite population that individual  $k$ 's response represents
  - ▶ **Caveat:** nonresponse
- ▶ **Weighted average**

$$\sum_{k=1}^n w_k y_k \text{ such that } w_k \in [0, 1] \text{ and } \sum_k w_k = 1$$



## Horvitz-Thompson estimators

- ▶ Consider a population of size  $N$ , a sample of size  $n$ , where each individual has outcome  $Y_k$
- ▶  $Y_k$  is **not** random, but  $Z_k$  is

$$Z_k = \begin{cases} 1, & k \in S \\ 0, & \text{otherwise.} \end{cases}$$

- ▶ Once sample taken  $y_k = Y_k \times Z_k$  denotes an individual's observed response (may contain measurement error)
  - ▶  $E[y_k] = E[Y_k \times Z_k] = Y_k E[Z_k] = Y_k \times \pi_k$

## Horvitz-Thompson estimators

- The population total of outcomes  $Y$  is

$$T = \sum_{k=1}^N Y_k$$

$$\hat{T} = \sum_{k=1}^n \underline{w_k y_k} = \sum_{k=1}^n \frac{y_k}{\pi_k}$$
$$\widehat{Var}(\hat{T}) = \sum_{k,k'} \frac{y_k y_{k'}}{\pi_k \pi_{k'}} - \frac{y_k y_{k'}}{\pi_{kk'}}$$

## Horvitz-Thompson estimators

- The population mean of outcomes  $Y$  is

$$\bar{Y} = \frac{\sum_{k=1}^N Y_k}{N}$$

$$\hat{\bar{Y}} = \frac{\sum_{k=1}^n w_k y_k}{N} = \frac{1}{N} \sum_{k=1}^n \frac{y_k}{\pi_k}$$

$$\widehat{\widehat{Var}}(\hat{\bar{Y}}) = \frac{\widehat{Var}(\hat{T})}{N^2}$$

$$\text{Var}(aX) = a^2 \text{Var}(X)$$

## Horvitz-Thompson estimators

- The population mean of binary outcomes  $Y$  or **prevalence** is

$$P = \frac{\sum_{k=1}^N Y_k}{N}$$

$$\hat{P} = \frac{\sum_{k=1}^n w_k y_k}{N} = \frac{1}{N} \sum_{k=1}^n \frac{y_k}{\pi_k}$$

$$\widehat{Var}(\hat{P}) = \frac{\widehat{Var}(\hat{T})}{N^2}$$

# Horvitz-Thompson estimators

- Stratified sampling

$$\hat{T} = \sum_{h=1}^H \hat{T}_h = \sum_{h=1}^H \sum_{k=1}^{n_h} w_{hk} y_{hk},$$

$$\widehat{Var}(\hat{T}) = \sum_{h=1}^H \widehat{Var}(\hat{T}_h) = \sum_{h=1}^H \left( \sum_{k,k'} \frac{y_{hk} y_{hk'}}{\pi_{hk} \pi_{hk'}} - \frac{y_{hk} y_{hk'}}{\pi_{hkk'}} \right),$$

- Calculate variance in terms of each individual's difference from their respective strata total.

# Horvitz-Thompson estimators

- ▶ Cluster sampling

$$\hat{T} = \sum_{c=1}^C T_c = \sum_{c=1}^C \sum_{k=1}^{N_c} w_{ck} y_{ck} = \sum_{c=1}^C w_c \sum_{k=1}^{N_c} y_{ck},$$

- ▶ Calculate the variance in terms of each cluster total's difference from the overall population total

# Horvitz-Thompson estimators

- Stratified two-stage cluster sampling

$$\begin{aligned}
 \hat{T} &= \sum_{h=1}^H \hat{T}_h = \sum_{h=1}^H \sum_{c_1=1}^{C_{1h}} \hat{T}_{h[c_1]} \\
 &= \sum_{h=1}^H \sum_{c_1=1}^{C_{1h}} \sum_{c_2=1}^{C_{2h}} \hat{T}_{h[c_1:c_2]} = \sum_{h=1}^H \sum_{c_1=1}^{C_{1h}} \sum_{c_2=1}^{C_{2h}} \sum_{k=1}^{n_{c_2}} w_{h[c_1:c_2]k} y_{h[c_1:c_2]k} \\
 \widehat{Var}(\hat{T}) &= \sum_{h=1}^H \widehat{Var}(\hat{T}_h).
 \end{aligned}$$

Handwritten notes in purple ink:

- Next to the first equation:  $EAS$  with arrows pointing to  $\hat{T}$  and  $\hat{T}_h$ , and  $n$  with an arrow pointing to  $\hat{T}_{h[c_1]}$ .
- Below the second equation: A bracket under the innermost sum  $\sum_{k=1}^{n_{c_2}}$  with a line pointing to the  $n_{c_2}$  term.

- Apply methods from previous two in appropriate summation order

# Horvitz-Thompson estimators

- ▶ What if we don't know  $N$ ?

$$\hat{P} = \frac{\sum_{k=1}^n w_k y_k}{N} \approx \frac{\sum_{k=1}^n w_k y_k}{\hat{N}} = \frac{\sum_{k=1}^n w_k y_k}{\sum_{k=1}^n w_k}$$

$$\widehat{Var}(\hat{P}) = \widehat{Var}\left(\frac{\hat{T}}{\hat{N}}\right) = \frac{\widehat{??}}{\widehat{??}}$$

- ▶ **Linearization:** Use Taylor series expansions to approximate the variance
  - ▶ survey package in R



## Binder (1983) and regression

- ▶ Linear regression mean model

$$E[Y|\theta, X] = X\beta,$$

- ▶ The **likelihood**

$$L(\theta|y, \mathbf{X}) = \prod_{k=1}^n L(\theta|y_k, \mathbf{x}_k)$$

- ▶ The **log-likelihood**

$$l(\theta|y, \mathbf{X}) = \log L(\theta|y, \mathbf{X}) = \sum_{k=1}^n \log L(\theta|y_k, \mathbf{x}_k)$$

## Binder (1983) and regression, cont'd

- ▶ The **score function**

$$\nabla l(\theta|y, \mathbf{X}) = \left[ \frac{\partial l}{\partial \beta_0} \quad \cdots \quad \frac{\partial l}{\partial \beta_p} \right]$$

is set equal to 0 to estimate  $\hat{\beta}$  in **maximum likelihood estimation**

- ▶ Incorporates sampling weights in **pseudolikelihood method** by weighting each observation unit's contribution to the score function by  $w_k$
- ▶ `survey::svyglm` function

# Model-based estimation

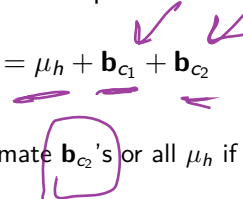
## Model-based estimation

Design-based methods have nice properties, but

- ▶ what if sampling weights not provided?
- ▶ small sample sizes → design-based standard errors too large
- ▶ especially a concern in **small area estimation**

Fixed effects for strata → different means for strata

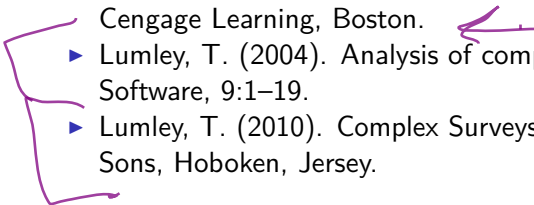
Random effects for cluster → account for dependence between observations within cluster

$$y_{h[c_1:c_2]k} = \mu_h + \mathbf{b}_{c_1} + \mathbf{b}_{c_2}$$


- ▶ Is there enough information to estimate  $\mathbf{b}_{c_2}$ 's or all  $\mu_h$  if  $H$  is large or  $m_{c_1:c_2}$  is small compared to  $m_{c_1}$ 's?

# References

## References

- ▶ Horvitz, D. G. and Thompson, D. J. (1952). A generalization of sampling without replacement from a finite universe. *Journal of the American statistical Association*, 47(260):663– 685.
  - ▶ Binder, D. (1983). On the variances of asymptotically normal estimators from complex surveys. *International Statistical Review*, 51:279–292.
  - ▶ Lohr, S. (2010). *Sampling: Design and Analysis*, Second Edition. Brooks/Cole Cengage Learning, Boston.
  - ▶ Lumley, T. (2004). Analysis of complex survey samples. *Journal of Statistical Software*, 9:1–19.
  - ▶ Lumley, T. (2010). *Complex Surveys: A Guide to Analysis using R*. John Wiley and Sons, Hoboken, Jersey.
- 

## Maximum Likelihood Estimation Examples, if needed

- ▶  $k = 1, \dots, n, y_k \overset{iid}{\sim} \text{Bernoulli}(p)$
- ▶  $k = 1, \dots, n, y_k \overset{iid}{\sim} \text{Normal}(\mu, \sigma^2)$