Bayesian Subnational Estimation using Complex Survey Data: Overview, Motivation and Survey Sampling

#### Jon Wakefield

Departments of Statistics and Biostatistics University of Washington

#### Outline

#### Overview Motivating Data Smoothing and Bayes

#### Survey Sampling

Design-Based Inference Complex Sampling Schemes

Discussion

# Overview

# Terminology

- Charactering and understanding subnational variation in health and demographic outcomes is an important public health endeavor.
- Many outcomes are binary, or public health targets are binary.
- For example, in the Sustainable Development Goals (SDGs), Goal 3.2 states, "By 2030, end preventable deaths of newborns and children under 5 years of age, with all countries aiming to reduce neonatal mortality to at least as low as 12 per 1,000 live births and under-5 mortality to at least as low as 25 per 1,000 live births".
- With respect to binary objectives, prevalence is defined as the proportion of a population who have a specific characteristic in a given time period.
- Examination of these proportions across space, is known as prevalence mapping – we may map continuously in space, or across discrete administrative areas.

# Terminology

- "The problem of small area estimation (SAE) is how to produce reliable estimates of characteristics of interest such as means, counts, quantiles, etc., for areas or domains for which only small samples or no samples are available, and how to assess their precision." (Pfeffermann, 2013).
- SAE methods provide one approach to performing prevalence mapping, for administrative areas.
- "The term geostatistics is a short-hand for the collection of statistical methods relevant to the analysis of geolocated data, in which the aim is to study geographical variation throughout a region of interest, but the available data are limited to observations from a finite number of sampled locations." (Diggle and Giorgi, 2019)
- Model-based geostatistics (MBG) provide another approach to performing prevalence mapping, over continuous space, though these continuous surfaces can be averaged for area-level inference.

#### **Overview of Lecture Series**

- Data: We consider the situation in which the available data arise from surveys with a complex design.
- A Problem: If small sample sizes in some areas/time periods, there is high instability. In the limit, there may be no data...
- Survey Sampling Methodology: Required for design and analysis.
- Shrinkage and Spatial Smoothing: To reduce instability, use the totality of data to smooth both locally and globally over space.
- Bayesian Modeling: Is convenient for encoding notions of smoothing, and for carrying out inference.
- Implementation: In R programming environment, using the SUMMER package.
- Visualization: Maps of uncertainty, accompanied with uncertainty, produced using the GIS capabilities of R.

#### **Overview of Lecture Series**

#### Lectures:

- Complex Survey Data.
- Bayesian Smoothing Models.
- Prevalence Mapping.
- Implementation, with examples, via the SUMMER package lectures by Zehang Richard Li.

#### Website:

http://faculty.washington.edu/jonno/space-station.html

The examples presented will mostly concern subnational estimation of under-5 mortality risk (U5MR).

### **Demographic Health Surveys**

- Motivation: In many developing world countries, vital registration is not carried out, so that births and deaths go unreported.
- Objective: To provide reliable estimates of demographic/health indicators at the (say) Admin1 or Admin2 level<sup>1</sup>, at which policy interventions are often carried out.
- We will illustrate using data from Demographic Health Surveys (DHS).
- DHS Program: Typically stratified cluster sampling to collect information on population, health, HIV and nutrition; more than 300 surveys carried out in over 90 countries, beginning in 1984.
- The Problem: Data are sparse, at the Admin2 level in particular.
- SAE: Leverage space-time similarity to construct a Bayesian smoothing model.

<sup>&</sup>lt;sup>1</sup>Admin0 = country level boundaries, Admin1 = first level administrative boundaries (states in US), Admin 2 = second level administrative boundaries (counties in US)

# 2014 Kenyan DHS

- The 3 most recent Kenya DHS were carried out in 2003, 2008 and 2014.
- The DHS use stratified two-stage cluster sampling. The strata consist of urban/rural crossed with geographic administrative strata.
- In each strata, enumeration areas (EAs) are selected with probability proportional to size using a sampling frame developed from the most recent census.
- In each of the clusters, households are selected. Within each household, women between the ages of 15 and 49 are interviewed.



Figure 1: Cluster locations in three Kenya DHS, with county boundaries.

# 2014 Kenya DHS

- We will focus on the 2014 Kenya DHS, in which the stratification was county (47) and urban/rural (2).
- Nairobi and Mombasa are entirely urban, so there are 92 strata in total.
- We have data from a total of 1584 EAs across the 92 strata. In the second stage, 40,300 households are sampled.
- DHS provides sampling (design) weights, assigned to each individual in the dataset.



Figure 2: Counties of Kenya.

### Aim: Inference for U5MR over Counties and Years



Figure 3: SAE estimates of under-5 mortality risk, across time, and Kenyan counties. These estimates were obtained using the SUMMER package.

# 2013 Nigeria DHS

- As a second DHS example, we consider measles vaccination rates in Nigeria, from the 2013 Nigerian DHS.
- Across African countries, there is great variability in the number of Admin2 areas.
- In Nigeria, the Admin2 areas correspond to Local Government Areas (LGAs) and there are 774 in total – with such a large number there are many LGAs with little/no data.
- There are no clusters in 255 LGAs.



Figure 4: Vaccination prevalence for LGAs in Nigeria. LGAs with no data are in white.

Specific methods are required for spatial data due to the dependence between points in space.

Within public and global health data different spatial methods are available for different endeavors:

- Disease Mapping: Spatial dependence is a virtue that we can exploit.
- Spatial Regression: Spatial dependence is a nuisance confounding by location.
- Cluster Detection: Spatial pattern of data is of primary interest.
- Assessment of Clustering: Spatial pattern of data is of primary interest.
- Small Area Estimation: Spatial dependence is a virtue that we can exploit.

Spatial methods often hinge upon some form of smoothing.

# Smoothing/Penalization

- When looking at estimates over space or time, we want to know if the differences we see are "real", or simply reflecting sampling variability.
- In data sparse situations, when one expects similarity, smoothing local patterns (in time, space, or both) can be highly beneficial.
- This can equivalently be thought of penalization, in which large deviations from "neighbors", suitably defined, are discouraged.
- We start with a temporal example, since time is easier to think about! One dimensional and an obvious direction...



Figure 5: Nile data with random walk of order 1 fits under different smoothing parameter choice.

### Temporal Smoothing for Ecuador U5MR



Figure 6: Yearly weighted estimates of under-5 mortality in Ecuador, with 95% uncertainty intervals for weighted and IHME, and 90% for UN IGME.

# Two Approaches to Prevalence Mapping

 Model at the area level using a discrete spatial model. These are the SAE models that are implemented in the SUMMER package.

• Model at the point level using a continuous spatial model. Model-based geostatistics is a popular approach.



# 2013 Nigeria DHS

- Recall that almost a third of the LGAs in Nigeria have no data (left plot below).
- We fit a discrete spatial model in which the rates in neighboring areas (as defined by sharing a boundary) are "encouraged" to be similar (right plot below).



Figure 7: Vaccination prevalences in Nigeria in 2013. Left: Weighted estimates. Right: Estimates from a discrete spatial smoothing model.

# Survey Sampling

# Outline

Many national surveys employ stratified cluster sampling, also known as multistage sampling, so that's where we'd like to get to.

We will discuss:

- Simple Random Sampling (SRS).
- Stratified SRS.
- · Cluster sampling.
- Multistage sampling.

First, we briefly explain why taking account of the survey design (data collection process) is important.

# Acknowledging the Design: Stratification



Figure 8: In the DHS, stratification is based on counties (the solid lines) and on a binary urban/rural variable (urban indicated in blue, the white is rural).

- Suppose we are interested in the proportion of women aged 20–29 who complete secondary education – this is much higher in urban areas
- If we oversample urban areas but ignore this when we analyze the data we will overestimate the fraction of women who complete secondary education, i.e., we will introduce bias.
- Taking into account of the stratification also reduces the variance of the estimator.
- In the design-based approach to inference, the stratification is accounted for via design weights.
- In the model-based approach to inference, the stratification is accounted for in the mean model.

# Acknowledging the Design: Cluster Sampling

- The DHS also employs cluster sampling, in which multiple units (individuals) within the same cluster are interviewed.
- Units within the same cluster tend to be more similar than units in different clusters, which reduces the information content of the clustered sample, relative to independently sampled units.
- The dependence can be measured via the intraclass correlation coefficient.
- In the design-based approach to inference, the clustering is accounted for in the variance calculation that is carried out.
- In the model-based approach to inference, the clustering is accounted for by including a cluster-specific random effect in the model.

### Modes of Inference

- Surveys can be analyzed using design- and model-based inference. In this lecture, the former will be focused upon.
- The target of inference are the set of means for areas indexed by *i* (e.g., Admin2 regions).
- Let y<sub>ik</sub> be the binary indicator on the k-th unit sampled in area i, for k ∈ S<sub>i</sub> (the set of selected individuals) and i = 1,..., n.

#### **Design-Based Inference**

- Labels *S<sub>i</sub>* of sampled units are random.
- Responses *y*<sub>ik</sub> are fixed.
- Asymptotic inference, perhaps using resampling.

#### Model-Based Inference

- Condition on units that are actually sampled.
- Responses *Y<sub>ik</sub>* are random.
- Exact inference, conditional on model.

Suppose we carry out stratified cluster sampling, with one-stage of clusters, and the outcome is continuous.

Let  $y_{ck}$  be the outcome from sampling unit *k* in sampled cluster *c*, and  $s_c$  the location of cluster *c*,

Suppose the data were collected within two strata, urban and rural.

A model-based approach to inference might begin with

$$Y_{ck} = \alpha + \gamma I(\mathbf{s}_c \in \text{rural}) + \epsilon_c + \upsilon_{ck},$$

where

- $\alpha$  is the mean for urban and  $\alpha + \gamma$  is the mean for rural.
- within-cluster dependence is modeled via the random effect  $\epsilon_c \sim_{iid} N(0, \sigma_{\epsilon}^2)$ .
- Measurement error is  $v_{ck} \sim_{iid} N(0, \sigma_v^2)$ .

#### **Design-Based Inference**

• We will focus on design-based inference: in this approach the population values of the variable of interest:

 $y_1,\ldots,y_N$ 

are viewed as fixed, while the indices of the individuals who are sampled are random.

- Imagine a population of size N = 4 and we sample n = 2
- Possible samples, with sampled unit indices in red and non-sampled in blue:

 $y_1, y_2, y_3, y_4$  $y_1, y_2, y_3, y_4$ 

 Different designs follow from which probabilities we assign to each of these possibilities. Design-based inference is frequentist, so that properties are based on hypothetical replications of the data collection process; hence, we require a formal description of the replication process.

A complex random sample may be:

- Better than a simple random sample (SRS) in the sense of obtaining the same precision at lower cost.
- May be worse in the sense of precision, but be required logistically.

# **Probability Samples**

Notation for random sampling, in a single population (and not distinguishing areas):

- N is population size.
- n is sample size.
- π<sub>k</sub> is the sampling probability for a unit (which will often correspond to a person) k, k = 1,..., N.

Random does not mean "equal chance", but means that the choice does not depend on variables/characteristics (either measured or unmeasured), except as explicitly stated via known sampling probabilities.

For example, in stratified random sampling, the probabilities of selection differ, in different strata.

# Common sampling designs

- Simple random sampling: Select each individual with probability  $\pi_k = n/N$ .
- Stratified random sampling: Use information on each individual in the population to define strata *h*, and then sample *n<sub>h</sub>* units independently within each stratum.
- Probability-proportional-to-size sampling: Given a variable related to the size of the sampling unit,  $Z_k$ , on each unit in the population, sample with probabilities  $\pi_k \propto Z_k$ .
- Cluster sampling: All units in the population are aggregated into larger units called clusters, known as primary sampling units (PSUs). Clusters are then sampled from this the set of PSUs, with units within these clusters being subsequently sampled.
- Multistage sampling: Stratified cluster sampling, with multiple levels of clustering.

# **Probability Samples**

• The label probability sample is often used instead of random sample.



Probability Sampling Vs Non-Probability Sampling

• Non-probability samples cannot be analyzed with design-based approaches, because there are no  $\pi_k$ .

#### Non-probability sampling approaches include:

• Convenience sampling (e.g., asking for volunteers). Also known as accidental or haphazard sampling.



- Purposive (also known as judgmental) sampling in which a researcher uses their subject knowledge to select participants (e.g, selecting an "average" looking individual).
- Quota sampling in which quotas in different groups are satisfied (but unlike stratified sampling, probability sampling is not carried out, for example, the interviewer may choose friendly looking people!).

#### For design-based inference:

- To obtain an unbiased estimator, every individual k in the population needs to have a non-zero probability  $\pi_k$  of being sampled, k = 1, ..., N.
- To carry out inference, this probability π<sub>k</sub> must be known only for every individual in the sample.
- So not needed for the unsampled individuals, which is key to implementation, since we will usually not know the sampling probabilities for those not sampled.

#### For design-based inference:

- To obtain a form for the variance of an estimator: for every pair of units, *k* and *l*, in the sample, there must a non-zero probability of being sampled together, call this probability, π<sub>kl</sub> for units *k* and *l*, *k* = 1,..., N, *l* = 1,..., N, *k* ≠ *l*.
- The probability  $\pi_{kl}$  must be known for every pair in the sample.
- in practice, these are often approximated, or the variance is calculated via a resampling technique such as the jackknife.

#### Inference

- Suppose we are interested in a variable denoted *y*, with the population values being *y*<sub>1</sub>,..., *y*<sub>N</sub>.
- Random variables will be represented by upper case letters, and constants by lower case letters.
- Finite population view: We have a population of size *N* and we are interested in characteristics of this population, for example, the mean:

$$\overline{y}_U = \frac{1}{N} \sum_{k=1}^N y_k.$$

- Infinite population view: The population variables are drawn from a hypothetical distribution, with mean  $\mu$ .
- In the model-based view, Y<sub>1</sub>,..., Y<sub>N</sub> are random variables and properties are defined with respect to p(·); often we say Y<sub>k</sub> are independent and identically distributed (iid) from p(·).
- As an estimator of  $\mu$ , we may take the sample mean:

$$\widehat{\mu} = \frac{1}{n} \sum_{k=1}^{n} Y_k.$$

- $\widehat{\mu}$  is a random variable because  $Y_1, \ldots, Y_n$  are each random variables.
- Assume Y<sub>k</sub> are iid observations from a distribution, p(·), with mean μ and variance σ<sup>2</sup>.
- The sample mean is an unbiased estimator, and has variance  $\sigma^2/n$ .

• Unbiased estimator:

$$E[\widehat{\mu}] = E\left[\frac{1}{n}\sum_{k=1}^{n}Y_{k}\right] = \frac{1}{n}\sum_{k=1}^{n}\underbrace{E[Y_{k}]}_{=\mu}$$
$$= \frac{1}{n}\sum_{k=1}^{n}\mu = \mu$$

• Variance:

$$\operatorname{var}(\widehat{\mu}) = \operatorname{var}\left(\frac{1}{n}\sum_{k=1}^{n}Y_{k}\right) \underbrace{=}_{\operatorname{iid}} \frac{1}{n^{2}}\sum_{k=1}^{n}\underbrace{\operatorname{var}(Y_{k})}_{=\sigma^{2}}$$
$$= \frac{1}{n^{2}}\sum_{k=1}^{n}\sigma^{2} = \frac{\sigma^{2}}{n}$$

• The variance  $\sigma^2$  is unknown so we estimate by the unbiased estimator

$$s^2 = \frac{1}{n-1} \sum_{k=1}^n (y_k - \widehat{\mu})^2.$$

A 95% asymptotic confidence interval is

$$\widehat{\mu} \pm 1.96 imes rac{s}{\sqrt{n}}.$$

 In practice, "asymptotic" means that *n* is sufficiently large that the sampling distribution of μ̂ (i.e., it's distribution in hypothetical repeated samples) is close to normal.

### **Design-Based Inference**

- In the design-based approach to inference the *y* values are treated as unknown but fixed.
- To emphasize: the y's are not viewed as random variables, so we write

 $y_1,\ldots,y_N,$ 

and the randomness, with respect to which all procedures are assessed, is associated with the particular sample of individuals that is selected, call the random set of indices *S*.

- Minimal reliance on distributional assumptions.
- Sometimes referred to as inference under the randomization distribution.
- In general, the procedure for selecting the sample is under the control of the researcher.

#### **Design-Based Inference**

• Define design weights as

$$w_k = \frac{1}{\pi_k}.$$

• The basic estimator is the weighted mean (Horvitz and Thompson, 1952; Hájek, 1971)

$$\widehat{\overline{y}}_U = \frac{\sum_{k \in S} w_k y_k}{\sum_{k \in S} w_k}$$

- This is an estimator of the finite population mean  $\overline{y}_U$ .
- So long as the weights are correctly calculated, and the sample size is not small, this estimator is appealing, though it may have high variance, if *n* is small.

### Simple Random Sample (SRS)

- The simplest probability sampling technique is simple random sampling without replacement.
- Suppose we wish to estimate the population mean in a particular population of size *N*.
- In everyday language: consider a population of size *N*; a random sample of size *n* ≤ *N* means that any subset of *n* people from the total number *N* is equally likely to be selected.

### Simple Random Sample (SRS)

• We sample *n* people from *N*, choosing each person independently at random and with the same probability of being chosen:

$$\pi_k = \frac{n}{N},$$

 $k = 1, \ldots, N.$ 

 Since sampling without replacement the joint sampling probabilities are

$$\pi_{kl}=\frac{n}{N}\times\frac{n-1}{N-1}$$

for  $k, l = 1, ..., N, k \neq l$ .

- In this situation:
  - The sample mean is an unbiased estimator.
  - The uncertainty, i.e. the variance, of the estimator can be easily estimated.
  - Unless *n* is quite close to *N*, the uncertainty does not depend on *N*, only on *n*.

# The Indices are Random!

• **Example:** N = 4, n = 2 with SRS. There are 6 possibilities:

 $\{y_1, y_2\}, \{y_1, y_3\}, \{y_1, y_4\}, \{y_2, y_3\}, \{y_2, y_4\}, \{y_3, y_4\}.$ 

- The random variable describing this design is *S*, the set of indices of those selected.
- The sample space of S is

 $\{(1,2), (1,3), (1,4), (2,3), (2,4), (3,4)\}$ 

and under SRS, the probability of sampling one of these possibilities is 1/6.

· The selection probabilities are

 $\pi_k = \Pr(\text{ individual } i \text{ in sample }) = \frac{3}{6} = \frac{1}{2}$ 

which is of course  $\frac{n}{N}$ .

 In general, we can work out the selection probabilities without enumerating all the possibilities!

#### **Design-Based Inference**

- Fundamental idea behind design-based inference: An individual with a sampling probability of  $\pi_k$  can be thought of as representing  $w_k = 1/\pi_k$  individuals in the population.
- **Example:** in SRS each person selected represents  $\frac{N}{n}$  people.
- The sum of the design weights,

$$\sum_{k\in\mathcal{S}}w_k=n\times\frac{N}{n}=N,$$

is the total population.

- Sometimes the population size may be unknown and the sum of the weights provides an unbiased estimator.
- In general, examination of the sum of the weights can be useful as if it far from the population size (if known) then it can be indicative of a problem with the calculation of the weights.

• The weighted estimator is

$$\widehat{\overline{y}}_{U} = \frac{\sum_{k \in S} w_{k} y_{k}}{\sum_{k \in S} w_{k}}$$
$$= \frac{\sum_{k \in S} \frac{N}{n} y_{k}}{\sum_{k \in S} \frac{N}{n}}$$
$$= \frac{\sum_{k \in S} y_{k}}{n} = \overline{y},$$

the sample mean, which is reassuring under SRS!

• This is an unbiased estimator, i.e.,

$$\mathsf{E}\left[\widehat{\overline{y}}_{U}\right] = \overline{y}_{U},$$

where we average over all possible samples we could have drawn, i.e., over S.

#### Unbiasedness

• For many designs:  $\sum_{k \in S} w_k = N$  so we examine the estimator

$$\widehat{\overline{y}}_U = \frac{1}{N} \sum_{k \in S} w_k y_k.$$

There's a neat trick in here, we introduce an indicator random variable of selection *l<sub>k</sub>* ~ Bernoulli(*π<sub>k</sub>*):

$$E\left[\widehat{\overline{y}}_{U}\right] = \underbrace{E\left[\frac{1}{N}\sum_{k\in\mathcal{S}}w_{k}y_{k}\right]}_{S \text{ is random in here}} = \underbrace{E\left[\frac{1}{N}\sum_{k=1}^{N}I_{k}w_{k}y_{k}\right]}_{I_{k} \text{ are random in here}}$$
$$= \frac{1}{N}\sum_{k=1}^{N}E\left[I_{k}\right]w_{k}y_{k} = \frac{1}{N}\sum_{i=1}^{N}\pi_{k}\frac{1}{\pi_{k}}y_{k} = \frac{1}{N}\sum_{i=1}^{N}y_{k} = \overline{y}_{U}$$

• It can be shown that the variance is

$$\operatorname{var}(\overline{y}) = \left(1 - \frac{n}{N}\right) \frac{S^2}{n},\tag{1}$$

where,

$$S^2 = \frac{1}{N-1}\sum_{k=1}^N (y_k - \overline{y}_U)^2.$$

- Contrast (1) with the model-based variance which is  $\sigma^2/n$ .
- The factor

$$1-\frac{n}{N}$$

is the famous finite population correction (fpc) factor.

- Because we are estimating a finite population mean, the greater the sample size relative to the population size, the more information we have (relatively speaking), and so the smaller the variance.
- In the limit, if *n* = *N* we have no uncertainty, because we know the population mean!

• The variance of the estimator depends on the population variance  $S^2$ , is unknown, and we estimate using the unbiased estimator:

$$s^2 = \frac{1}{n-1} \sum_{k \in S} (y_k - \overline{y})^2.$$

• Substitution into (1) gives an unbiased estimator of the variance:

$$\widehat{\operatorname{var}}(\overline{y}) = \left(1 - \frac{n}{N}\right) \frac{s^2}{n}.$$
 (2)

• The standard error is

$$\mathsf{SE}(\overline{y}) = \sqrt{\left(1 - \frac{n}{N}\right)\frac{s^2}{n}}.$$

• Note:  $S^2$  is not a random variable but  $s^2$  is.

If n, N and N − n are "sufficiently large"<sup>2</sup>, a 95% asymptotic confidence interval for y
<sub>U</sub> is

$$\overline{y} \pm 1.96 \times \sqrt{1 - \frac{n}{N}} \frac{s}{\sqrt{n}}.$$
 (3)

- The interval given by (3) is random (across samples) because  $\overline{y}$  and  $s^2$  (the estimate of the variance) are random.
- In practice therefore, if *n* ≪ *N*, we obtain the same confidence interval whether we take a design- or a model-based approach to inference (though the interpretation is different).

<sup>&</sup>lt;sup>2</sup>so that the normal distribution provides a good approximation to the sampling distribution of the estimator

# Stratified Sampling

- Simple random samples are rarely taken in surveys because they are logistically difficult and there are more efficient designs for gaining the same precision at lower cost.
- Stratified random sampling is one way of increasing precision and involves dividing the population into groups called strata and drawing probability samples from within each one, with sampling from different strata being carried out independently.
- An important practical consideration of whether stratified sampling can be carried out is whether stratum membership is known for every individual in the population, i.e., we need a sampling frame containing the strata variable.

Lohr (2010, Section 3.1) provides a good discussion of the benefits of stratified sampling, we summarize here.

- Protection from the possibility of a "really bad sample", i.e., very few or zero samples in certain stratum giving an unrepresentative sample.
- Obtain known precision required for subgroups (domains) of the population this is usual for the DHS.
- For example, from the Kenya DHS sampling manual (Kenya National Bureau of Statistics, 2015):

"The 2014 KDHS was designed to produce representative estimates for most of the survey indicators at the national level, for urban and rural areas separately, at the regional (former provincial) level, and for selected indicators at the county level."

### Rationale for Stratified Sampling

- Flexible since sampling frames can be constructed differently in different strata.
- For example, one may carry out different sampling in urban and rural areas.
- More precise estimates can be obtained if stratum can be found that are associated with the response of interest, for example, age and gender in studies of human disease.
- In a national study, the most natural form of sampling may be based on geographical regions.
- Due to the independent sampling in different stratum, variance estimation is straightforward, as long as within-stratum sampling variance estimators are available.

# Example: Washington State

- According to the census there were 2,629,126 households in Washington State in the period 2009–2013.
- Consider a simple random sample (SRS) of 2000 households, so that each household has a

$$\frac{2000}{2629126} = 0.00076,$$

chance of selection.

• Suppose we wish to estimate characteristics of household in all 39 counties of WA.

# Example: Washington State





- King (highlighted left) and Garfield (highlighted right) counties had 802,606 and 970 households so that under SRS we will have, on average, about 610 households sampled from King County and about 0.74 from Garfield county.
- The probability of having no-one from Garfield County is about 22% (binomial experiment), and the probability of having more than one is about 45%.
- If we took exactly 610 from King and 1 (rounding up) from Garfield we have an example of proportional allocation, which would not be a good idea given the objective here.
- Stratified sampling would allow control of the number of samples in each county.

# Notation

- Stratum levels are denoted h = 1, ..., H, so H in total.
- Let  $N_1, \ldots, N_H$  be the known population totals in the stratum with

 $N_1 + \cdots + N_H = N$ 

so that *N* is the total size of the population.

 In stratified simple random sampling, the simplest from of stratified sampling, we take a SRS from each stratum with n<sub>h</sub> samples being randomly taken from stratum h, so that the total sample size is

#### Figure 1: Comparison of Simple Random Sampling to Stratified Random Sampling





Visual of Simple Random Sampling Selection of 6 out of 18 People

Visual of Stratified Random Sampling: Selection of 3 out of 9 Men and 3 out of 9 Women

- We can view stratified SRS as carrying out SRS in each of the *H* stratum; we let *S<sub>h</sub>* represent the probability sample in stratum *h*.
- We also let *S* refer to the overall probability sample.

#### **Estimators**

• The sampling probabilities for unit k in strata h are

$$\pi_{hk} = \frac{n_h}{N_h},$$

which do not depend on k.

• Therefore the design weights are

$$w_{hk}=rac{N_h}{n_h}.$$

• Note that:

$$\sum_{h=1}^{H}\sum_{k\in \mathcal{S}_h}w_{hk}=\sum_{h=1}^{H}\sum_{k\in \mathcal{S}_h}\frac{N_h}{n_h}=\sum_{h=1}^{H}n_h\frac{N_h}{n_h}=N,$$

so that summing over the weights recovers the population size.

#### **Estimators**

• Weighted estimator:

$$\widehat{\overline{y}}_{U} = \frac{\sum_{h=1}^{H} \sum_{k \in S_{h}} w_{hk} y_{hk}}{\sum_{h=1}^{H} \sum_{k \in S_{h}} w_{hk}} = \sum_{h=1}^{H} \frac{N_{h}}{N} \overline{y}_{h}$$

where

$$\overline{y}_h = \frac{\sum_{k \in S_h} y_{hk}}{n_h}.$$

 Since we are sampling independently from each stratum using SRS, we have<sup>3</sup>

$$\operatorname{var}(\widehat{\overline{y}}_U) = \sum_{h=1}^{H} \left(\frac{N_h}{N}\right)^2 \left(1 - \frac{n_h}{N_h}\right) \frac{s_h^2}{n_h},\tag{4}$$

where the within stratum variances are:

$$s_h^2 = \frac{1}{n_h - 1} \sum_{k \in S_h} (y_{hk} - \overline{y}_h)^2.$$

<sup>&</sup>lt;sup>3</sup>using the variance formula for SRS, (2)

### Weighted Estimation

Recall: The weight  $w_k$  can be thought of as the number of people in the population represented by sampled person k.

#### Example 1: Simple Random Sampling

Suppose an area contains 1000 people:

- Using simple random sampling (SRS), 100 people are sampled.
- Sampled individuals have weight  $w_k = 1/\pi_k = 1000/100 = 10$ .

#### Example 2: Stratified Simple Random Sampling

Suppose an area contains 1000 people, 200 urban and 800 rural.

- Using stratified SRS, 50 urban and 50 rural individuals are sampled.
- Urban sampled individuals have weight

 $w_k = 1/\pi_k = 200/50 = 4.$ 

• Rural sampled individuals have weight  $w_k = 1/\pi_k = 800/50 = 16.$ 

#### Example 2: Stratified Simple Random Sampling

Suppose an area contains 1000 people, 200 urban and 800 rural.

- Urban risk = 0.1.
- Rural risk = 0.2.
- True risk = 0.18.

Take a stratified SRS, 50 urban and 50 rural individuals sampled:

- Urban sampled individuals have weight 4; 5 cases out of 50.
- Rural sampled individuals have weight 16; 10 cases out of 50.
- Simple mean is  $15/100 = 0.15 \neq 0.18$ .
- Weighted mean is

$$\frac{4 \times 5 + 16 \times 10}{4 \times 50 + 16 \times 50} = \frac{180}{1000} = 0.18.$$

### Motivation for Cluster Sampling

For logistical reasons, cluster sampling is an extremely common design that is often used for government surveys.

Two main reasons for the use of cluster sampling:

- A sampling frame for the population of interest does not exist, i.e., no list of population units.
- The population units have a large geographical spread and so direct sampling is not logistically feasible to implement for in-person interviews.
- It is far more cost effective (in terms of travel costs, etc.) to cluster sample.

# Terminology

- In single-stage cluster sampling or one-stage cluster sampling, the population is grouped into subpopulations (as with stratified sampling) and a probability sample of these clusters is taken, and every unit within the selected clusters is surveyed.
- In one-stage cluster sampling either all or none of the elements that compose a cluster (PSU) are in the sample.
- The subpopulations are known as clusters or primary sampling units (PSUs).
- In two-stage cluster sampling, rather than sample all units within a PSU, a further cluster sample is taken; the possible groups to select within clusters are known as secondary sampling units (SSUs).
- This can clearly be extended to multistage cluster sampling.

# Differences Between Cluster and Stratified sampling

Stratified Random Sampling	One-Stage Cluster Sampling
A sample is taken from every stratum	Observe all elements only within the sampled clusters
Variance of estimate of $\overline{y}_U$ depends on within strata variability	The cluster is the sampling unit and the more clusters sampled the smaller the variance – which depends primarily on between cluster means
For greatest precision, we want low within-strata variability but large between-strata variability	For greatest precision, high within-cluster variability and similar cluster means.
Precision generally better than SRS	Precision generally worse than SRS



Stratified Sampling Vs Cluster Sampling

# Heterogeneity

- The reason that cluster sampling loses efficiency over SRS is that within clusters we only gain partial information from additional sampling within the same cluster, since within clusters two individuals tend to be more similar than two individuals within different clusters.
- The similarity of elements within clusters is due to unobserved (or unmodeled) variables.
- The design effect (deff) is often to summarize the effect on the variance of the design:

 $\label{eq:deff} \mbox{deff} = \frac{\mbox{Variance of estimator under design}}{\mbox{Variance of estimator under SRS}},$ 

where in the denominator we use the same number of observations as in the complex design in the numerator.

# Estimation for One-Stage Cluster Sampling

- We suppose that a SRS of *n* PSUs is taken.
- The probability of sampling a PSU is *n*/*N*, and since all the SSUs are sampled in each selected PSU we have selection probabilities and design weights:

$$\pi_{ik} = \Pr(\text{ SSU } k \text{ in cluster } i \text{ is selected }) = \frac{n}{N}$$
$$w_{ik} = \text{Design weight for SSU } k \text{ in cluster } i = \frac{N}{n}.$$

Let S represent the set of sampled clusters.

### Estimation for One-Stage Cluster Sampling

• Let  $M_0 = \sum_{i=1}^{N} M_i$  be the total number of secondary sampling units (SSUs), i.e., elements in the population, so the population mean is

$$\overline{y}_U = \frac{1}{M_0} \sum_{i=1}^N \sum_{k=1}^{M_i} y_{ik}$$

• An unbiased estimator is

$$\widehat{\overline{y}}_{U} = \frac{\sum_{i \in S} \sum_{k \in S_{i}} w_{ik} y_{ik}}{M_{0}}$$

• Then,

$$\widehat{\operatorname{var}}(\widehat{\overline{y}}_U) = \frac{N^2}{M_0^2} \left(1 - \frac{n}{N}\right) \frac{s_T^2}{n}$$

where  $s_T^2$  is the estimated variance of the PSU totals.

# Two-Stage Cluster Sampling with Equal-Probability Sampling

It may be wasteful to measure all SSUs in the selected PSUs, since the units may be very similar and so there are diminishing returns on the amount of information we obtain.

We discuss the equal-probability two stage cluster design:

- 1. Select a SRS of *n* PSUs from the population of *N* PSUs.
- 2. Select a SRS of *m<sub>i</sub>* SSUs from each selected PSU, the probability sample collected will be denoted *S<sub>i</sub>*.

### Two-Stage Cluster Sampling Weights

• The selection probabilities are:

Pr(k-th SSU in i-th PSU selected) = Pr(i-th PSU selected)

× Pr(*k*-th SSU | *i*-th PSU selected)  
= 
$$\frac{n}{N} \times \frac{m_i}{M_i}$$

· Hence, the weights are

$$w_{ik} = \pi_{ik}^{-1} = \frac{N}{n} \times \frac{M_i}{m_i}.$$

An unbiased estimator is

$$\widehat{\overline{y}}_U = \frac{\sum_{i \in S} \sum_{k \in S_i} w_{ik} y_{ik}}{M_0}.$$

• Variance calculation is not trivial, and requires more than knowledge of the weights.

# Variance Estimation for Two-Stage Cluster Sampling

- In contrast to one-stage cluster sampling we have to acknowledge the uncertainty in both stages of sampling; in one-stage cluster sampling the totals t<sub>i</sub> are known in the sampled PSUs, whereas in two stage sampling we have estimates t<sub>i</sub>.
- In Lohr (2010, Chapter 6) it is shown that

$$\operatorname{var}(\widehat{\overline{y}}_{U}) = \frac{1}{M_{0}^{2}} \left[ \underbrace{N^{2}\left(1 - \frac{n}{N}\right) \frac{s_{T}^{2}}{n}}_{\operatorname{One-stage cluster variance}} + \underbrace{\frac{N}{n} \sum_{i \in S} \left(1 - \frac{m_{i}}{M_{i}}\right) M_{i}^{2} \frac{s_{i}^{2}}{m_{i}}}_{\operatorname{Two-stage cluster variance}} \right]$$
(5)

where

- $s_T^2$  are the estimated variance of the cluster totals,
- $s_i^2$  is the estimated variance within the *i*-th PSU.
- In most software packages, the second term in (5) is ignored, since it is small when compared to the first term, when *N* is large.

### The Jackknife

- The jackknife is a very general technique for calculating the variance of an estimator.
- The basic idea is to delete portions of the data, and then fit the model on the remainder if one repeats this process for different portions, one can empirically obtain the distribution of the estimator.
- The key is to select the carefully select the portion of the data so that the design is respected.
- We describe in the context of multistage cluster sampling.
- Observations within a PSU should be kept together when constructing the data portions, which preserves the dependence among observations in the same PSU.

# The Jackknife for Multistage Cluster Sampling

- Assume we have *H* strata and *n<sub>h</sub>* PSUs in strata *h*, and assume PSUs are chosen with replacement.
- To apply the jackknife, delete one PSU at a time.
- Let  $\hat{\mu}_{(hi)}$  be the estimator when PSU *i* of stratum *h* is omitted.
- To calculate  $\hat{\mu}_{(hi)}$  we define a new weight variable:

$$w_{k(hi)} = \begin{cases} w_{k(hi)} & \text{if observation } k \text{ is not in stratum } h \\ 0 & \text{if observation } k \text{ is in PSU } i \text{ of stratum } h \\ \frac{n_h}{n_h - 1} w_k & \text{if observation } k \text{ is not in PSU } i \text{ but in stratum } h \end{cases}$$

Then we can use the weights  $w_{k(hi)}$  to calculate  $\hat{\mu}_{(hi)}$  and

$$\widehat{V}_{\mathsf{JK}}(\widehat{\mu}) = \sum_{h=1}^{H} \frac{n_h - 1}{n_h} \sum_{i=1}^{n_h} (\widehat{\mu}_{(hi)} - \widehat{\mu})^2.$$

# Multistage Sampling in the DHS

- A common design in national surveys is multistage sampling, in which cluster sampling is carried out within strata.
- DHS Program: Typically, 2-stage stratified cluster sampling:
  - Strata are urban/rural and region.
  - Enumeration Areas (EAs) sampled within strata (PSUs).
  - Households within EAs (SSUs).
- Weighted estimators are used and a common approach to variance estimation is the jackknife (Pedersen and Liu, 2012)
- In later lectures, we will show how model-based inference can be carried out for the DHS.

# Discussion

### Discussion

- The majority of survey sampling texts take a design-based view of inference this is a different paradigm to model-based inference, for which most spatial statistical models were developed!
- Later we will see how spatial models can incorporate the survey design.
- Variance estimation that accounts for the design has been a topic of much research.
- However, for the major designs (e.g., SRS, stratified SRS, cluster sampling, multistage sampling), weighted estimates and their variances are available within all the major statistical packages.
- When the variance is large, because of small sample sizes, we would like to use smoothing methods, with Bayes being a convenient way to do this this is the topic of the next lecture.

### Acknowledgments

This lecture series was supported by the Hewlett Foundation and the International Union for the Scientific Study of Population (IUSSP).

The research reported in this series has grown out of a longstanding collaboration between Jon Wakefield, Zehang Richard Li and Sam Clark.

Many other people have contributed, however, for full details and for links to other aspects of this work, check out:

http://faculty.washington.edu/jonno/space-station.html

#### References

- Diggle, P. J. and Giorgi, E. (2019). *Model-based Geostatistics for Global Public Health: Methods and Applications*. Chapman and Hall/CRC.
- Hájek, J. (1971). Discussion of, "An essay on the logical foundations of survey sampling, part I", by D. Basu. In V. Godambe and D. Sprott, editors, *Foundations of Statistical Inference*. Holt, Rinehart and Winston, Toronto.
- Horvitz, D. and Thompson, D. (1952). A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association*, **47**, 663–685.
- Kenya National Bureau of Statistics (2015). Kenya Demographic and Health Survey 2014. Technical report, Kenya National Bureau of Statistics.
- Lohr, S. (2010). *Sampling: Design and Analysis, Second Edition.* Brooks/Cole Cengage Learning, Boston.
- Pedersen, J. and Liu, J. (2012). Child mortality estimation: Appropriate time periods for child mortality estimates from full birth histories. *PLoS Medicine*, **9**, e1001289.
- Pfeffermann, D. (2013). New important developments in small area estimation. *Statistical Science*, **28**, 40–68.