

# Simple Random Sampling (SRS)

## Finite Population

$$\mathcal{P} = \{x_1, x_2, \dots, x_N\}$$

$$\Rightarrow \mu = \frac{1}{N} \sum_{j=1}^N x_j, \quad \sigma^2 = \frac{1}{N-1} \sum_{j=1}^N (x_j - \mu)^2, \quad \dots$$

**Sample** Draw  $n < N$  items from  $\mathcal{P}$  without replacement  $\Rightarrow \mathcal{S}$ .

e.g.,  $\mathcal{S} = \{x_3, x_7, \dots, x_{91}\}$

$$\Rightarrow \hat{\mu} = \frac{1}{n} \sum_{j \in \mathcal{S}} x_j, \quad \hat{\sigma}^2 = \frac{1}{n-1} \sum_{j \in \mathcal{S}} (x_j - \hat{\mu})^2, \quad \dots$$

## Source of Randomness

$$I_j \equiv I(j \in \mathcal{S}), \quad \pi_j \equiv \mathbb{P}(j \in \mathcal{S}), \quad \pi_{j\ell} \equiv \mathbb{P}(j \in \mathcal{S} \text{ and } \ell \in \mathcal{S})$$

$$\hat{\mu} = \frac{1}{n} \sum_{j=1}^N x_j I_j$$

$$\begin{aligned} \Rightarrow \mathbb{E}(\hat{\mu}) &= \frac{1}{n} \sum_{j=1}^N x_j \mathbb{E}(I_j) = \frac{1}{n} \sum_{j=1}^N x_j \pi_j \\ &= \frac{1}{n} \sum_{j=1}^N x_j \frac{n}{N} = \frac{1}{N} \sum_{j=1}^N x_j = \mu \end{aligned}$$

**Remark** For SRS,

$$\pi_j = \frac{\binom{1}{1} \binom{N-1}{n-1}}{\binom{N}{n}} = \frac{n}{N} \quad \forall j = 1, 2, \dots, N.$$

# Variance

$$\hat{\mu} = \frac{1}{n} \sum_{j=1}^N x_j I_j$$

$$\Rightarrow \text{Var}(\hat{\mu}) = \frac{1}{n^2} \left[ \sum_{j=1}^N x_j^2 \text{Var}(I_j) + \sum_{j \neq \ell} x_j x_\ell \text{Cov}(I_j, I_\ell) \right]$$

$$= \frac{1}{n^2} \left[ \sum_{j=1}^N x_j^2 \pi_j (1 - \pi_j) + \sum_{j \neq \ell} x_j x_\ell (\pi_{j\ell} - \pi_j \pi_\ell) \right]$$

$$= \dots = \left(1 - \frac{n}{N}\right) \frac{\sigma^2}{n}$$

**Exercise** Show that, for SRS,  $\pi_{j\ell} = (n/N)[(n-1)/(N-1)]$  (**easy**) and complete the missing steps (**not hard but unpleasant**).

# Unpleasant Details

$$\begin{aligned}
 & \frac{1}{n^2} \left[ \sum_{j=1}^N x_j^2 \binom{n}{N} \left(1 - \frac{n}{N}\right) + \sum_{j \neq \ell} x_j x_\ell \left( \frac{n}{N} \frac{n-1}{N-1} - \frac{n}{N} \frac{n}{N} \right) \right] \\
 &= \frac{1}{n^2} \binom{n}{N} \left[ \sum_{j=1}^N x_j^2 \left( \frac{N-n}{N} \right) + \underbrace{\sum_{j \neq \ell} x_j x_\ell \left( \frac{n-1}{N-1} - \frac{n}{N} \right)}_{\frac{n-N}{N(N-1)}} \right] \\
 &= \frac{1}{n^2} \binom{n}{N} \left( \frac{N-n}{N} \right) \frac{1}{N-1} \underbrace{\left[ (N-1) \sum_{j=1}^N x_j^2 - \sum_{j \neq \ell} x_j x_\ell \right]}_{\stackrel{(\dagger)}{=} N(N-1)\sigma^2} \\
 & \qquad \qquad \qquad = \left(1 - \frac{n}{N}\right) \frac{\sigma^2}{n}
 \end{aligned}$$

## Step (†)

$$\mu^2 = \left[ \frac{1}{N} \sum_{j=1}^N x_j \right]^2 = \frac{1}{N^2} \left[ \sum_{j=1}^N x_j^2 + \sum_{j \neq \ell} x_j x_\ell \right]$$

$$\sigma^2 = \frac{1}{N-1} \sum_{j=1}^N (x_j - \mu)^2 = \dots = \frac{1}{N-1} \left\{ \sum_{j=1}^N x_j^2 - N\mu^2 \right\}$$

$$= \frac{1}{N-1} \left\{ \sum_{j=1}^N x_j^2 - \frac{1}{N} \left[ \sum_{j=1}^N x_j^2 + \sum_{j \neq \ell} x_j x_\ell \right] \right\}$$

$$= \frac{1}{N(N-1)} \left\{ (N-1) \sum_{j=1}^N x_j^2 - \sum_{j \neq \ell} x_j x_\ell \right\}$$

# Stratified Sampling

## Population

$$\mathcal{P} = \underbrace{\{x_{1,1}, \dots, x_{1,N_1}\}}_{\mathcal{P}_1} \oplus \underbrace{\{x_{2,1}, \dots, x_{2,N_2}\}}_{\mathcal{P}_2} \oplus \dots \oplus \underbrace{\{x_{K,1}, \dots, x_{K,N_K}\}}_{\mathcal{P}_K}$$

$$\mu_k = \frac{1}{N_k} \sum_{j=1}^{N_k} x_{k,j}, \quad \sigma_k^2 = \frac{1}{N_k - 1} \sum_{j=1}^{N_k} (x_{k,j} - \mu_k)^2, \quad k = 1, 2, \dots, K$$

**Sample** Draw **simple random sample**  $\mathcal{S}_k$  of size  $n_k < N_k$  from  $\mathcal{P}_k$ , **independently** for  $k = 1, 2, \dots, K \Rightarrow \mathcal{S} = \mathcal{S}_1 \oplus \dots \oplus \mathcal{S}_K$ .

**Examples** Stratify by race, gender, age (18-24, 25-34, ...), etc.

# Population Quantities

$$N = N_1 + N_2 + \dots + N_K$$

$$\mu = \frac{1}{N} \sum_{k=1}^K \sum_{j=1}^{N_k} x_{k,j} = \sum_{k=1}^K \frac{N_k}{N} \mu_k$$

$$\begin{aligned} \sigma^2 &= \frac{1}{N-1} \sum_{k=1}^K \sum_{j=1}^{N_k} (x_{k,j} - \mu)^2 = \frac{1}{N-1} \sum_{k=1}^K \sum_{j=1}^{N_k} (x_{k,j} - \mu_k + \mu_k - \mu)^2 \\ &= \frac{1}{N-1} \left[ \sum_{k=1}^K (N_k - 1) \sigma_k^2 + \sum_{k=1}^K N_k (\mu_k - \mu)^2 \right] \end{aligned}$$

# Sample Quantities

Within Strata:

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{x_{k,j} \in \mathcal{S}_k} x_{k,j} = \frac{1}{n_k} \sum_{j=1}^{N_k} x_{k,j} I(x_{k,j} \in \mathcal{S}_k)$$

$$\mathbb{E}(\hat{\mu}_k) = \dots = \mu_k, \quad \text{Var}(\hat{\mu}_k) = \dots = \left(1 - \frac{n_k}{N_k}\right) \frac{\sigma_k^2}{n_k} \quad (\text{from SRS})$$

Overall:  $\hat{\mu}_{st} = \sum_{k=1}^K \frac{N_k}{N} \hat{\mu}_k \Rightarrow$

$$\left\{ \begin{array}{l} \mathbb{E}(\hat{\mu}_{st}) = \sum_{k=1}^K \frac{N_k}{N} \mathbb{E}(\hat{\mu}_k) = \sum_{k=1}^K \frac{N_k}{N} \mu_k = \mu \\ \text{Var}(\hat{\mu}_{st}) = \sum_{k=1}^K \left(\frac{N_k}{N}\right)^2 \text{Var}(\hat{\mu}_k) = \sum_{k=1}^K \left(\frac{N_k}{N}\right)^2 \left(1 - \frac{n_k}{N_k}\right) \frac{\sigma_k^2}{n_k} \end{array} \right.$$

# Proportional Allocation

If  $n_k = n(N_k/N)$ , then

$$\text{Var}(\hat{\mu}_{st}) = \frac{1}{n} \left(1 - \frac{n}{N}\right) \left[ \sum_{k=1}^K \left(\frac{N_k}{N}\right) \sigma_k^2 \right];$$

whereas, for SRS,

$$\begin{aligned} \text{Var}(\hat{\mu}) &= \left(1 - \frac{n}{N}\right) \frac{\sigma^2}{n} \\ &= \frac{1}{n} \left(1 - \frac{n}{N}\right) \left[ \underbrace{\sum_{k=1}^K \frac{(N_k - 1)}{N - 1} \sigma_k^2}_{\frac{N_k - 1}{N - 1} \approx \frac{N_k}{N}} + \underbrace{\sum_{k=1}^K \frac{N_k}{N - 1} (\mu_k - \mu)^2}_{\text{main reduction here between-strata variation}} \right]. \end{aligned}$$

almost always better; best if stratas **different** in terms of  $\mu_k$

# Horvitz-Thompson Estimator

To estimate population quantity

$$\theta = \sum_{j=1}^N h(x_j), \quad \text{e.g.,} \quad h(x_j) = \begin{cases} x_j & \text{(total)} \\ x_j/N & \text{(mean)} \\ I(x_j \in \mathcal{A})/N & \text{(proportion)} \\ \dots & \text{(...)} \end{cases},$$

use

$$\hat{\theta}_{\text{HT}} = \sum_{j \in \mathcal{S}} \frac{h(x_j)}{\pi_j}, \quad \pi_j = \mathbb{P}(j \in \mathcal{S})$$

where  $\pi_j$  depends on the specific probability sampling scheme.

**Remark** Principle = inverse probability weighting. Always unbiased (next slide). Generic way to derive estimators.

# Horvitz-Thompson Estimator

$$\hat{\theta}_{\text{HT}} = \sum_{j=1}^N \frac{h(x_j)}{\pi_j} I_j$$

$$\Rightarrow \mathbb{E}(\hat{\theta}_{\text{HT}}) = \sum_{j=1}^N \frac{h(x_j)}{\pi_j} \mathbb{E}(I_j) = \sum_{j=1}^N \frac{h(x_j)}{\pi_j} \pi_j = \sum_{j=1}^N h(x_j) = \theta$$

$$\Rightarrow \text{Var}(\hat{\theta}_{\text{HT}}) = \sum_{j=1}^N y_j^2 \text{Var}(I_j) + \sum_{j \neq \ell} y_j y_\ell \text{Cov}(I_j, I_\ell)$$

$$= \sum_{j=1}^N y_j^2 \pi_j (1 - \pi_j) + \sum_{j \neq \ell} y_j y_\ell (\pi_{j\ell} - \pi_j \pi_\ell) = \dots$$

where  $y_j = h(x_j)/\pi_j$

## Example: SRS

$$\mathbb{P}(j \in \mathcal{S}) = \frac{n}{N} \quad \forall j$$

$$\mu = \frac{1}{N} \sum_{j=1}^N x_j \quad \Rightarrow \quad h(x_j) = x_j/N$$

$$\hat{\mu} = \sum_{j \in \mathcal{S}} \frac{h(x_j)}{\mathbb{P}(j \in \mathcal{S})} = \sum_{j \in \mathcal{S}} \frac{x_j/N}{n/N} = \frac{1}{n} \sum_{j \in \mathcal{S}} x_j$$

# Cluster Sampling

## Population

$$\mathcal{P} = \underbrace{\{x_{1,1}, \dots, x_{1,N_1}\}}_{\mathcal{C}_1} \oplus \underbrace{\{x_{2,1}, \dots, x_{2,N_2}\}}_{\mathcal{C}_2} \oplus \dots \oplus \underbrace{\{x_{M,1}, \dots, x_{M,N_M}\}}_{\mathcal{C}_M}$$

$$\mu_g = \frac{1}{N_g} \sum_{j=1}^{N_g} x_{g,j}, \quad \sigma_g^2 = \frac{1}{N_g - 1} \sum_{j=1}^{N_g} (x_{g,j} - \mu_g)^2, \quad g = 1, 2, \dots, M$$

**Sample** Draw **simple random sample** of  $m < M$  clusters  $\Rightarrow$  e.g.,  
 $\mathcal{S} = \mathcal{C}_1 \oplus \mathcal{C}_{23} \oplus \dots \oplus \mathcal{C}_{97}$ .

**Example** Cluster = household, classroom, postal code, etc.

# Similar Analysis

(Details Omitted)

- apply HT principle to obtain unbiased estimator

$$\mu = \sum_{g=1}^M \frac{N_g}{N} \mu_g \quad \Rightarrow \quad \hat{\mu}_{\text{cl}} = \sum_{\mathcal{C}_g \in \mathcal{S}} \frac{(N_g/N) \mu_g}{\mathbb{P}(\mathcal{C}_g \in \mathcal{S})} = \sum_{\mathcal{C}_g \in \mathcal{S}} \frac{(N_g/N) \mu_g}{m/M}$$

- analyze  $\text{Var}(\hat{\mu}_{\text{cl}})$
- messier algebra ... slight simplification by assuming  $N_1 = N_2 = \dots = N_M$  (equal cluster size)



not necessarily better; best if clusters **similar** in terms of  $\mu_g$

# Illustration

$$\mathcal{P} = \{1, 2, 3, 1, 2, 3, 1, 2, 3\}$$

ideal stratification

$$\mathcal{P}_1 = \{1, 1, \mathbf{1}\}$$

$$\mathcal{P}_2 = \{\mathbf{2}, 2, 2\}$$

$$\mathcal{P}_3 = \{3, 3, \mathbf{3}\}$$

$$\sigma_k^2 = 0$$
$$\sum_k (\mu_k - \mu)^2 \text{ large}$$

ideal clustering

$$\mathcal{C}_1 = \{1, 2, 3\}$$

$$\mathcal{C}_2 = \{\mathbf{1}, \mathbf{2}, \mathbf{3}\}$$

$$\mathcal{C}_3 = \{1, 2, 3\}$$

$$\sigma_g^2 \text{ large}$$
$$\sum_g (\mu_g - \mu)^2 = 0$$

analysis so far “merely” saying these obvious truths

# Other Topics

- multi-stage sampling
  - e.g., cluster sampling, then SRS again within cluster
- network sampling (**especially relevant at the moment**)
  - e.g., snowball sampling (i.e., ask a friend)
- and so on ...

## Key to Analysis

probability of ending up in the sample