

# Fundament of Inference - Michael Elliott

Different adjustment estimators. E.g.,  
post-stratification;  
model-assisted estimator;  
model-based estimator

2018-2019 Spring

## Auxiliary information

### Point I : What's poststratification

**Poststratification** (often based on auxiliary information), developing weights so that the weighted mean of the sample matches that of population (external information)

1
2
:
H

} Population  $N$  has  $H$  strata, and  $N_h$  are known.

For an unit  $i$  that belongs to stratum  $h$  :  $w_{hi} = \frac{N_h}{n_h}$

Poststratified estimator of population total :

$$\hat{T}_{ps} = \sum_{i=1}^n w_{hi} y_{hi} = \sum_{h=1}^H \sum_{i=1}^{n_h} w_{hi} y_{hi} = \sum_{h=1}^H \sum_{i=1}^{n_h} \frac{N_h}{n_h} \cdot y_{hi} = \sum_{h=1}^H N_h \sum_{i=1}^{n_h} \frac{y_{hi}}{n_h}$$

$$= \sum_{h=1}^H N_h \cdot \bar{y}_h$$

Poststratified estimator of population mean :

$$\hat{Y}_{ps} = \frac{\hat{T}_{ps}}{N} = \frac{\sum_{h=1}^H N_h \cdot \bar{y}_h}{N} = \sum_{h=1}^H \frac{N_h}{N} \bar{y}_h$$

$= \sum_{h=1}^H p_h \cdot \bar{y}_h$  ← note that this is also  $\bar{y}_{ps}$  (poststratified sample mean)

Point II.

**Poststratification usefulness 1:** Poststratified estimator has reduced variance

Poststratification is useful bc it reduces variance of the estimator poststratified estimator's variance: dataset is a random variable, thus means

$$V(\bar{y}_{ps}) = E(V(\bar{y}_{ps}|n_1, n_2, \dots, n_h)) + V(E(\bar{y}_{ps}|n_1, n_2, \dots, n_h))$$

in which  $E(\bar{y}_{ps}|n_1, n_2, \dots, n_h) = E\left(\sum_{h=1}^H P_h \bar{Y}_h\right) = \sum_{h=1}^H P_h E(\bar{Y}_h) = \sum_{h=1}^H P_h \bar{Y}_h$  (i.e.) current data

$$= \bar{Y}$$

$$V(\bar{y}_{ps}|n_1, n_2, \dots, n_h) = V\left(\sum_{h=1}^H P_h \bar{Y}_h\right) = \sum_{h=1}^H P_h^2 S_h^2$$

$$\therefore V(\bar{y}_{ps}) = \underbrace{E\left(\sum_{h=1}^H P_h^2 \frac{S_h^2}{n_h}\right)}_{\text{approximately}} + V(\bar{Y})^2 \approx \sum_{h=1}^H P_h^2 \frac{S_h^2}{E(n_h)} = \sum_{h=1}^H P_h^2 \frac{S_h^2}{n \cdot P_h}$$

approximately  $\underbrace{\sum_{h=1}^H P_h S_h^2 / n}$

$$E\left(\frac{1}{n_h}\right) > \frac{1}{E(n_h)} : \left(\frac{1}{3} + \frac{1}{2}\right)/2 > \frac{1}{(3+2)/2}$$

while we know that  $E\left(\frac{1}{n_h}\right) > \frac{1}{E(n_h)}$ , there is no close equation to calculate by how much.

$\therefore$  Approximate a term to compensate

$$V(\bar{y}_{ps}) \approx \frac{1}{n} \left[ \sum_{h=1}^H P_h S_h^2 + \underbrace{\sum_{h=1}^H P_h (1-P_h) \frac{S_h^2}{n_h}}_{\text{term to compensate}}$$

SRS's variance: ignoring finite population correction

$$V(\bar{y}) = \frac{S^2}{n}$$

where  $S^2$  is population variance  $= \frac{1}{N-1} \sum_{h=1}^H \sum_{i=1}^{N_h} (Y_{hi} - \bar{Y}_h)^2$

$$= \frac{1}{N-1} \sum_{h=1}^H \sum_{i=1}^{N_h} (Y_{hi} - \bar{Y}_h + \bar{Y}_h - \bar{Y})^2$$

$$= \frac{1}{N-1} \sum_{h=1}^H \sum_{i=1}^{N_h} [(Y_{hi} - \bar{Y}_h)^2 + (\bar{Y}_h - \bar{Y})^2 + 2(Y_{hi} - \bar{Y}_h)(\bar{Y}_h - \bar{Y})]$$

$$= \frac{1}{N-1} \left[ \sum_{h=1}^H \sum_{i=1}^{N_h} (Y_{hi} - \bar{Y}_h)^2 + \sum_{h=1}^H \sum_{i=1}^{N_h} (\bar{Y}_h - \bar{Y})^2 + 2 \sum_{h=1}^H \sum_{i=1}^{N_h} (Y_{hi} - \bar{Y}_h)(\bar{Y}_h - \bar{Y}) \right]$$

$$= \frac{1}{N-1} \left[ \sum_{h=1}^H \sum_{i=1}^{N_h} (Y_{hi} - \bar{Y}_h)^2 + \sum_{h=1}^H \sum_{i=1}^{N_h} (\bar{Y}_h - \bar{Y})^2 + 2 \sum_{h=1}^H (\bar{Y}_h - \bar{Y}) \sum_{i=1}^{N_h} (Y_{hi} - \bar{Y}_h) \right]$$

$$= \frac{1}{N-1} \left[ \sum_{h=1}^H (N_h - 1) S_h^2 + \sum_{h=1}^H N_h (\bar{Y}_h - \bar{Y})^2 \right]$$

$$\approx \frac{1}{N} \left[ \sum_{h=1}^H N_h \cdot S_h^2 + \sum_{h=1}^H N_h (\bar{Y}_h - \bar{Y})^2 \right]$$

$$= \sum_{h=1}^H P_h \cdot S_h^2 + \sum_{h=1}^H P_h (\bar{Y}_h - \bar{Y})^2$$

$$\therefore V(\bar{y}) = \underbrace{\frac{1}{n} \sum_{h=1}^H P_h \cdot S_h^2}_{V(\bar{y}_{ps})} + \underbrace{\frac{1}{n} \sum_{h=1}^H P_h (\bar{Y}_h - \bar{Y})^2}_{V(\bar{Y})}$$

$$= V(\bar{y}_{ps}) + \frac{1}{n} \sum_{h=1}^H P_h (\bar{Y}_h - \bar{Y})^2$$

$\Rightarrow V(\bar{y})$  under SRS is  $\geq V(\bar{y}_{ps})$

### Point III.

#### Poststratification usefulness 2: Correct for non-response bias

Poststratification is originally designed for reducing estimator's variance, but the post-stratified estimator also correct bias in the setting of non-response.

First, a little bit about non-response error

Deterministic approach : population consists of respondents ( $\bar{Y}_R$ )

non-respondents ( $\bar{Y}_{NR}$ )

$$\text{Then Bias} = \underbrace{(1-R)}_{\text{proportion of respondents}} (\bar{Y}_R - \bar{Y}_{NR})$$

proportion of respondents

Stochastic approach : Each element in the population has a response propensity  $P_i$

$$\text{Bias} = \frac{C(Y, P)}{\bar{P}} \leftarrow \begin{array}{l} \text{covariance between} \\ Y \text{ and } P \end{array}$$

$\bar{P} \leftarrow \text{average } P$

Assuming deterministic approach :

$$\begin{aligned} \bar{Y} &= \sum_h^H P_h \cdot \bar{Y}_h = \sum_h^H P_h \left( \frac{N_R}{N} \bar{Y}_{hR} + \frac{N-N_R}{N} \bar{Y}_{hNR} \right) \\ &= \sum_h^H P_h [R_h \bar{Y}_{hR} + (1-R_h) \bar{Y}_{hNR}] \\ &= \sum_h^H P_h \cdot R_h \bar{Y}_{hR} + \sum_h^H P_h \cdot \bar{Y}_{hNR} - \sum_h^H P_h \cdot R_h \bar{Y}_{hNR} \\ &= \sum_h^H P_h \cdot \bar{Y}_{hR} - \sum_h^H P_h \bar{Y}_{hR} + \sum_h^H P_h \cdot R_h \bar{Y}_{hR} + \sum_h^H P_h \cdot \bar{Y}_{hNR} - \sum_h^H P_h \cdot R_h \bar{Y}_{hNR} \\ &= \sum_h^H P_h \cdot \bar{Y}_{hR} - \sum_h^H (1-R_h) P_h \cdot \bar{Y}_{hR} + \sum_h^H (1-R_h) P_h \cdot \bar{Y}_{hNR} \\ &\Downarrow \\ E(\bar{Y}_R) &= \sum_h^H P_h \cdot \bar{Y}_{hR} - \left[ \sum_h^H (1-R_h) P_h (\bar{Y}_{hR} - \bar{Y}_{hNR}) \right] \end{aligned}$$

$$\begin{aligned} \text{Bias} &= \bar{Y}_R - \bar{Y} = \frac{\sum_h^H N_{hR} Y_{hRi}}{N_R} - \bar{Y} = \sum_h^H \frac{1}{N_R} \sum_i^N Y_{hRi} - \bar{Y} \\ &= \sum_h^H \frac{N_{hR} \cdot \bar{Y}_{hR}}{N_R} - \bar{Y} = \sum_h^H \frac{N_h R_h}{N \cdot R} \bar{Y}_{hR} - \bar{Y} \\ &= \sum_h^H \frac{R_h}{R} P_h \cdot \bar{Y}_{hR} - \sum_h^H P_h \cdot \bar{Y}_{hR} + \left[ \sum_h^H (1-R_h) P_h (\bar{Y}_{hR} - \bar{Y}_{hNR}) \right] \\ &= \sum_h^H \left( \frac{R_h}{R} - 1 \right) P_h \cdot \bar{Y}_{hR} + \sum_h^H (1-R_h) P_h (\bar{Y}_{hR} - \bar{Y}_{hNR}) \\ &= \sum_h^H \frac{R_h - R}{R} P_h \cdot \bar{Y}_{hR} - \underbrace{\sum_h^H \frac{R_h - R}{R} P_h \cdot \bar{Y}_R}_{\text{this is } 0} + \sum_h^H (1-R_h) P_h (\bar{Y}_{hR} - \bar{Y}_{hNR}) \\ &= \underbrace{\sum_h^H \frac{R_h - R}{R} P_h (\bar{Y}_{hR} - \bar{Y}_R)}_{\text{This part can be estimated based on the observed data.}} + \underbrace{\sum_h^H (1-R_h) P_h (\bar{Y}_{hR} - \bar{Y}_{hNR})}_{\sum_h^H P_h (1-R_h) (\bar{Y}_{hR} - \bar{Y}_{hNR}) \text{ should be smaller than } (1-R)(\bar{Y}_R - \bar{Y}_{NR})} \end{aligned}$$

This part can be estimated based on the observed data.

$\sum_h^H P_h (1-R_h) (\bar{Y}_{hR} - \bar{Y}_{hNR})$  should be smaller than  $(1-R)(\bar{Y}_R - \bar{Y}_{NR})$

Continuing the previous page :

$$\begin{aligned} \text{Bias} &= \bar{Y}_R - \bar{Y} \\ &= \underbrace{\sum_{h=1}^H \frac{P_h - R}{R} P_h (\bar{Y}_{hR} - \bar{Y}_R)}_{A} + \underbrace{\sum_{h=1}^H P_h (1 - P_h) (\bar{Y}_{hR} - \bar{Y}_{hNR})}_{B} \end{aligned}$$

Non-response missingness can be divided into 3 situations:

1. MCAR (missing completely at random).

$$\text{Then } P_h = R \quad \therefore A = 0$$

$$\bar{Y}_{hR} = \bar{Y}_{hNR} \quad \therefore B = 0 \quad \Rightarrow \text{no bias}$$

2. MAR (missing at random)

$$\text{Conditioning on stratum, } \bar{Y}_{hR} = \bar{Y}_{hNR} \quad \therefore B = 0$$

$$\text{Bias} = A$$

3. MNAR (missing not at random)

$$\text{Bias} = A + B$$

Finally get to the main point, what can poststratification do?

Poststratified estimation of a sample

$$\begin{aligned} \bar{y}_{ps} &= \sum_{h=1}^H P_h \cdot \bar{y}_{Rh} \quad \leftarrow \text{ignore the response rate } (\bar{Y}_R = \sum_{h=1}^H \frac{P_h}{R} P_h \cdot \bar{Y}_{Rh}) \\ &= \frac{\sum_{h=1}^H N_h \cdot \bar{y}_{Rh}}{\sum_{h=1}^H N_h} \quad \text{and use respondents' mean in a stratum to} \\ &= \frac{\sum_{h=1}^H \sum_{i=1}^{N_h} N_h \cdot \bar{y}_{Rhi} / n_{Rh}}{\sum_{h=1}^H N_h} = \frac{\sum_{h=1}^H \sum_{i=1}^{n_{Rh}} \frac{n_R}{N} \frac{N_h \cdot \bar{y}_{Rhi}}{n_{Rh}}}{\sum_{h=1}^H \frac{n_R}{N} N_h} \\ &= \frac{\sum_{h=1}^H \sum_{i=1}^{n_{Rh}} w_{hi} \bar{y}_{Rhi}}{\sum_{h=1}^H \sum_{i=1}^{n_{Rh}} w_{hi}} \quad \text{where } w_{hi} = \frac{N_h/N}{n_{Rh}/n_R} \end{aligned}$$

$$\begin{aligned} \text{Bias}_{ps} &= E(\bar{y}_{ps}) - \bar{Y} = E\left(\sum_{h=1}^H P_h \bar{y}_{Rh}\right) - \left[\sum_{h=1}^H P_h \cdot \bar{Y}_{hR} - \sum_{h=1}^H (1 - P_h) P_h (\bar{Y}_{hR} - \bar{Y}_{hNR})\right] \\ &= \sum_{h=1}^H P_h \bar{E}(\bar{y}_{Rh}) - \sum_{h=1}^H P_h \cdot \bar{Y}_{hR} + \sum_{h=1}^H (1 - P_h) P_h (\bar{Y}_{hR} - \bar{Y}_{hNR}) \\ &\stackrel{0}{=} \sum_{h=1}^H (1 - P_h) P_h (\bar{Y}_{hR} - \bar{Y}_{hNR}) \end{aligned}$$

Comparing  $\text{Bias}_{ps}$  to  $\text{Bias}$ . A part is gone.

Poststratification is a reasonable practice under MAR assumption,

$\Rightarrow$  Both A and B are gone.  $\therefore$  with MAR could have unbiased estimator with poststratification.

## Following Point II

For large enough sample size, variance of poststratified mean will be less than or equal to SRS mean.

But this does NOT hold if the sample size shrinks:

$$V(\bar{y}_{ps}) \approx \frac{1}{n} \sum_{h=1}^H P_h S_h^2 \left( 1 + \frac{1-P_h}{n P_h} \right)$$

$$\therefore V(\bar{y}_{ps}) - V(\bar{y}) \approx \frac{1}{n^2} \sum_{h=1}^H S_h^2 (1-P_h) - \frac{1}{n} \sum_{h=1}^H P_h (\bar{Y}_h - \bar{Y})^2$$

This could be positive or negative, so

$V(\bar{y}_{ps})$  could be larger or smaller than  $V(\bar{y})$

## Following Point III

For bias induced by non-response, poststratified estimator always improve over SRS estimator if there is a relationship between

a) the probability of response and strata

$$R_h \neq R$$

b) the respondents' mean in the strata and the overall mean

$$\bar{Y}_{hR} \neq \bar{Y}_R$$

$\Rightarrow$  There is an A part in bias for the poststratification to remove  
If MAR, B part is gone

$$\therefore |B(\bar{y}_{r,ps}) - B(\bar{y}_r)| = |A| \quad \leftarrow \text{difference between the 2 biases is } |A|$$

## Point IV. Summing up Point II and Point III.

$$\text{Mean Square Error} = \text{Bias}^2 + \text{Variance}$$

Difference between the 2 MSE:

$$\begin{aligned} \text{MSE}(\bar{y}_{r,ps}) - \text{MSE}(\bar{y}_r) &= \text{Bias}^2(\bar{y}_{r,ps}) - \text{Bias}^2(\bar{y}_r) + V(\bar{y}_{r,ps}) - V(\bar{y}_r) \\ &= \underbrace{\frac{1}{n^2} \sum_{h=1}^H S_h^2 (1-P_h)}_{\text{above}} - \underbrace{\frac{1}{n} \sum_{h=1}^H P_h (\bar{Y}_h - \bar{Y})^2}_{\text{B}} - \underbrace{\left[ \sum_{h=1}^H \frac{P_h - R}{R} P_h (\bar{Y}_{hR} - \bar{Y}_R) \right]^2}_{\text{A}} \end{aligned}$$

## Point V. Raking (in contrast to all discussion on poststratification above)

Iteratively match with marginal distributions of auxiliary variables one at a time.

Better to see an example. Don't summarize here.

## Point VI. Poststratification and raking under GREG (generalized regression)

Both poststratification and raking fit within a framework — calibration estimation.

The idea is to recalibrate the survey weights  $d_i$  to new values  $w_i$ , where  $w_i$  satisfy the constraint  $\sum w_i \vec{x}_i = \vec{x}$  where  $\vec{x}$  is a vector of known population totals for auxiliary variable  $x_i$ .

(Poststratification can be think of one example of  $w_i$ , that works on top of the original weight  $d_i=1$  for all individuals)

if the original weights are  $d_i$  (e.g., based on selection probability)

$$w_i = d_i \left[ 1 + (\vec{x} - \hat{\vec{x}})' \vec{T}^{-1} \vec{x}_i \right]$$

$$\text{where } \vec{T} = \sum d_i \vec{x}_i \vec{x}_i' \text{ and } \hat{\vec{x}} = \sum d_i \vec{x}_i$$

To break down the formula above:

- $\vec{x}$  observed auxiliary variables to be matched.
- $\hat{\vec{x}}$  estimated values of the auxiliary variables based on  $d_i$
- $\vec{x} - \hat{\vec{x}}$  difference between actual  $x$  and  $d_i$ -based estimated  $x$ .  
if this is small, the  $d_i$  is regarded as already calibrated.
- $\vec{T}$  weighted cross product
- $\vec{T} \vec{x}_i$  can shrink or expand  $d_i$  depending on the sign of  $\vec{x} - \hat{\vec{x}}$

## Point VI.

The discussion above are under a design-based perspective. Each unit  $i$  represents  $w_i$  units in the population.  
Can also think of it as a model alternatively

First, in the above discussion  $\bar{Y}$  is fixed — a population characteristic that the sample tries to capture.

But can also conceptualize  $\bar{Y}$  as a random variable — Bayesian

$$\bar{Y} = \sum_{h=1}^H P_h \cdot \bar{Y}_h = \sum_{h=1}^H P_h \left( \frac{n_h}{N_h} \bar{y}_h + \frac{N_h - n_h}{N_h} \cdot \bar{Y}_{h,ns} \right)$$

↑                          ↑  
 Sample                      non-sampled.  
 data/observation         need to work out a posterior distribution

$y_{hi} | \mu_h, \sigma_h^2, z_{i=h} \sim N(\mu_h, \sigma_h^2)$  allow each stratum to have its own mean and variance  
 $P(\mu_h, \log \sigma_h^2) \propto 1$  flat prior

Bayes rule  
is used here  
to flip the probability  
to become the probability  
of the underlying conditions

⇒ posterior distribution  $\mu_h | y, z \sim N(\bar{y}_h, \frac{\sigma_h^2}{n_h})$

(deduction not shown here)  $\sigma_h^2 | y, z \sim \text{Inverse } \chi^2(n_h-1, S_h)$

on the observed. posterior distributions of  $\mu_h$  and  $\sigma_h^2$  are not enough. Since we are dealing with a finite population situation here, need to leverage  $\mu_h$  and  $\sigma_h^2$  and further work out a distribution of  $\bar{Y}_{h,ns}$  conditioning on data.

we know that  $\bar{Y}_{h,ns} | \mu_h, \sigma_h^2 \sim N(\mu_h, \frac{\sigma_h^2}{N_{h,ns}})$ . observing data doesn't change this distribution  $\bar{Y}_{h,ns} | \mu_h, \sigma_h^2, y, z \sim N(\mu_h, \frac{\sigma_h^2}{N_{h,ns}})$

$$\begin{aligned}
 E(\bar{Y}_{h,ns} | \sigma_h^2, y, z) &= E(E(\bar{Y}_{h,ns} | \sigma_h^2, y, z, \mu_h) | \sigma_h^2, y, z) \\
 &= E(\mu_h | \sigma_h^2, y, z) = \bar{y}_h \\
 V(\bar{Y}_{h,ns} | \sigma_h^2, y, z) &= E(V(\bar{Y}_{h,ns} | \sigma_h^2, y, z, \mu_h) | \sigma_h^2, y, z) + V(E(\bar{Y}_{h,ns} | \sigma_h^2, y, z, \mu_h) | \sigma_h^2, y, z) \\
 &= E\left(\frac{\sigma_h^2}{N_{h,ns}} | \sigma_h^2, y, z\right) + V(\mu_h | \sigma_h^2, y, z) \\
 &= \frac{\sigma_h^2}{N_{h,ns}} + \frac{\sigma_h^2}{n_h} = \frac{\sigma_h^2}{(1 - \frac{n_h}{N_h}) n_h}
 \end{aligned}$$

Bayes rule is no longer applied in this part.

$\mu_h$  and  $\sigma_h^2$  are like a helper bridge.

with their help, work out a distribution of  $\bar{Y}_{h,ns}$  continue to the next page the rest of the finite population conditioned on the sample

Remove  $\sigma_h^2$  from condition

$$E(\bar{Y}_{h,ns} | y, z) = E\left(E(\bar{Y}_{h,ns} | y, z, \sigma_h^2) | y, z\right) = E(\bar{y}_h | y, z) = \bar{y}_h$$

$$V(\bar{Y}_{h,ns} | y, z) = E(V(\bar{Y}_{h,ns} | y, z, \sigma_h^2) | y, z) + V(E(\bar{Y}_{h,ns} | y, z, \sigma_h^2) | y, z)$$

$$= E\left(\frac{\sigma_h^2}{(1 - \frac{n_h}{N_h})n_h} | y, z\right) + V(\bar{y}_h | y, z)$$

$$= \frac{n_h - 1}{n_h - 3} \frac{s_h^2}{(1 - \frac{n_h}{N_h})n_h} \quad \text{sample variance}$$

$\bar{y}_h$  doesn't have a variance conditioning on data.

$$E(\bar{Y}_h | y, z) = E\left(\frac{n_h}{N_h} \bar{y}_h + \frac{N_h - n_h}{N_h} \bar{Y}_{h,ns} | y, z\right) = \bar{y}_h$$

$V(\bar{Y}_h | y, z)$  only comes from  $\bar{Y}_{h,ns}$

$$= V\left(\frac{n_h}{N_h} \bar{y}_h + \frac{N_h - n_h}{N_h} \bar{Y}_{h,ns} | y, z\right)$$

$$= \left(\frac{N_h - n_h}{N_h}\right)^2 \cdot \frac{n_h - 1}{n_h - 3} \frac{s_h^2}{(1 - \frac{n_h}{N_h})n_h}$$

$$= \frac{N_h - n_h}{N_h} \frac{n_h - 1}{n_h - 3} \frac{s_h^2}{n_h}$$

$$E(\bar{Y} | y, z) = E\left(\sum_h p_h \bar{Y}_h | y, z\right) = \sum_h p_h \bar{y}_h$$

$$V(\bar{Y} | y, z) = V\left(\sum_h p_h \bar{Y}_h | y, z\right) = \sum_h p_h^2 \frac{n_h - 1}{n_h - 3} \frac{N_h - n_h}{N_h} \frac{s_h^2}{n_h}$$

+ covariance

0

Extend the above discussion:

$$y_{hi} | \mu_h, \sigma_h^2, z_i = h \sim N(\mu_h, \sigma_h^2)$$

If the prior of  $\mu_h$  is not flat  $\mu_h \sim N(\lambda, \tau^2)$   
still  $p(\log \sigma_h^2) \propto 1$

$$E(\bar{Y}_h | y, z, \sigma_h^2, \tau^2)$$

# Module 10

## Theoretical differentiation

2 Situations :

- 1) model as a mean to an end. The goal is to estimate descriptive statistics of the population. Models have two usages.
  - a. One is model-assisted inference:  
models are used to improve efficiency of design-based approach
  - b. The other is model-based inference:  
take survey design into account. When use of model is required for inference we need to borrow information across different elements of the data.
- 2) model is the end. The model itself is the focus of inference.  
Now the question is whether to incorporate design into the inference about the model.

## Model-assisted inference Point I.

Where is model?

Recall GREG. Reweight the original weights  $d_i$  to  $w_i$ :

$$w_i = d_i [1 + (\vec{X} - \hat{\vec{X}})' \vec{T}^{-1} \vec{x}_i]$$

There is a model behind this formula — why the calculation is this way  
Consider using  $w_i$  for estimating population total.

$$\begin{aligned}\hat{T}_w &= \sum^n w_i y_i = \sum^n d_i [1 + (\vec{X} - \hat{\vec{X}})' \vec{T}^{-1} \vec{x}_i] y_i = \sum^n d_i y_i + \sum^n d_i (\vec{X} - \hat{\vec{X}})' \vec{T}^{-1} \vec{x}_i \cdot y_i \\ &= \sum^n d_i y_i + \sum^n d_i \vec{x}' \vec{T}^{-1} \vec{x}_i y_i - \sum^n d_i \vec{x}' \vec{T}^{-1} \vec{x}_i y_i \\ &= \underbrace{\sum^n d_i y_i}_{\text{standard design-based estimator.}} + \underbrace{\vec{x}' (\sum^n d_i \vec{x}_i \vec{x}_i') (\sum^n d_i \vec{x}_i y_i)}_{\text{by definition } \hat{\beta}_w} + \underbrace{\vec{x}' (\sum^n d_i \vec{x}_i \vec{x}_i')^{-1} (\sum^n d_i \vec{x}_i y_i)}_{\text{model}} \\ &= \underbrace{\sum^n d_i y_i}_{\text{standard design-based estimator.}} + (\vec{X} - \hat{\vec{X}})' \hat{\beta}_w\end{aligned}$$

∴ By changing  $d_i$  to  $w_i$ , GREG essentially converts the original estimator based on  $d_i$  by adding an extra part  $(\vec{X} - \hat{\vec{X}})' \hat{\beta}_w$ . The model is between  $y$  (variable of interest) and  $\vec{X}$  (auxiliary variables). The model part makes correction if the original weights  $d_i$  are not sufficient to bring  $\vec{X}$  to the population and there is a relationship between  $\vec{X}$  and  $y$ .

Technically, we only need aggregate  $\vec{X}$  and the sample to make the above GREG correction. But assuming that we have information of  $\vec{X}$  on the entire population (which is not completely impossible), then we can rewrite the above population total estimate.

$$\begin{aligned}\hat{T}_w &= \sum^n d_i y_i + (\vec{X} - \hat{\vec{X}})' \hat{\beta}_w \\ &= \sum^n d_i y_i + \left( \sum^N \vec{x}_i - \sum^n d_i \vec{x}_i \right) \hat{\beta}_w = \sum^n d_i y_i + \sum^N \vec{x}_i \hat{\beta}_w - \sum^n d_i \vec{x}_i \hat{\beta}_w \\ &= \sum^N \vec{x}_i \hat{\beta}_w + \sum^n d_i (y_i - \vec{x}_i \hat{\beta}_w) = \sum^N y_i^w + \underbrace{\sum^n d_i (y_i - \hat{y}_i^w)}_{=0} \\ &= \sum^N y_i^w \quad \text{model}\end{aligned}$$

$\hat{y}_i^w$  is based on the weighted linear regression  
by definition of regression.

$$\therefore \hat{T}_w = \sum^n d_i y_i + (\vec{X} - \hat{\vec{X}})' \hat{\beta}_w = \sum^N \vec{x}_i \hat{\beta}_w$$

Two ways of expressing the GREG-weighted estimated total.

Either use the model ( $\hat{\beta}_w$ ) to adjust the discrepancy on  $\vec{X}$  ( $\vec{X}' - \hat{\vec{X}}'$ )

Or apply the model ( $\hat{\beta}_w$ ) directly on  $\vec{x}_i$  of the entire population units

Either way  $\hat{T}_w$  can be viewed as a regression estimator (based on  $\vec{X}$  and  $y$ )

A general form of regression estimator is  $\hat{T}_R = \sum^N \hat{y}_i + \sum^n d_i (y_i - \hat{y}_i)$  Comparing to above, weights are only used to adjust the difference between  $y_i$  and  $\hat{y}_i$ , but not used for constructing  $\hat{y}_i$  in the first place. Don't equal to 0 as above bc the model is unweighted so  $\sum (y_i - \hat{y}_i) = 0$  and  $\sum d_i (y_i - \hat{y}_i) \neq 0$

Advantages comparing to  $\hat{T}_w$ : ①  $\hat{y}_i$  is more stable than  $\hat{y}_i^w \rightarrow$  reduce variance

② Can still correct for model misspecification  $y_i - \hat{y}_i \rightarrow$  but increase variance

$$T = \sum^N (y_i + u_i) = \sum^N \hat{y}_i + \sum^N u_i$$

Is  $\hat{T}_R$  biased (deviate from  $T$ )?  $\leftarrow \sum u_i \neq \sum d_i (y_i - \hat{y}_i)$

## Model-assisted inference continue

The previous discussion talks about how GREG can be framed as a regression estimator—An estimator that is based on the relationship between auxiliary variable's  $\vec{x}$  and variable of interest  $y$  ( $\vec{y}_w$ ). Recall that in GREG:  $w_i = d_i [1 + (\vec{x} - \vec{\bar{x}})^T T^{-1} \vec{x}_i]^2$ . One example to interpret  $w_i$  and  $d_i$  is to view  $d_i$  as selection weights and  $w_i$  as adjustment on top of the selection weight.

It seems straightforward to see  $w_i$ -based estimator as a model between  $y$  and  $x$  because  $w_i$  is achieved by using  $x_i$ .

However, not only can  $w_i$ -based estimator be framed as a model, the  $d_i$ -based estimator can also be framed as a model, present below:

How to frame HT-estimator (di-based estimator) as a regression model:

with  $d_i = \frac{1}{\pi_i}$  reverse of selection probability, HT Estimator =  $\sum_{i=1}^n d_i y_i = \sum_{i=1}^n \frac{y_i}{\pi_i}$   
Since  $\pi_i$  is only information, consider a regression that models the relationship between  $y_i$  and  $\pi_i$ :  $E(Y_i) = \beta \pi_i$   $V(Y_i) = \sigma^2 \pi_i^2$   $\rightarrow$  A weighted least square regression.  
(no intercept)

That is:  $Y_i \sim N(\beta \pi_i, \sigma^2 \pi_i^2)$

Equivalent to:  $Y_i = \beta \cdot \pi_i + \epsilon_i$   $\epsilon_i \sim N(0, \pi_i^2 \sigma^2)$

Assuming a helper variable  $Z_i = Y_i / \pi_i$   $\therefore E(Z_i) = \beta$   $V(Z_i) = \sigma^2$   
given a sample of realization of  $Z_i$ ,

$$\frac{n}{\sum n} \frac{Z_i}{\pi_i} = \frac{n}{\sum \pi_i} \frac{Y_i}{\pi_i} = \frac{1}{n} \cdot \text{mean of this sample} = \hat{\beta}$$

$\therefore$  HT-estimator =  $n \cdot \hat{\beta}$   $\leftarrow$   $\hat{\beta}$  is based on the model between  $y_i$  and  $\pi_i$

Recall that if a regression model is used to assist the estimate, then the regression estimate (of population total) can be written as

$$\hat{T}_R = \sum_{i=1}^N \hat{y}_i + \sum_{i=1}^n d_i (y_i - \hat{y}_i)$$

$$\begin{aligned} \hat{T}_R &= \sum_{i=1}^N \beta \cdot \pi_i + \sum_{i=1}^n \frac{1}{\pi_i} (y_i - \beta \pi_i) = \sum_{i=1}^N \beta \cdot \pi_i + \sum_{i=1}^n \frac{y_i}{\pi_i} - \sum_{i=1}^n \beta = \sum_{i=1}^n \frac{y_i}{\pi_i} + \beta \left( \sum_{i=1}^n \pi_i - n \right) \\ &= \sum_{i=1}^n \frac{y_i}{\pi_i} = \text{HT-estimator} \end{aligned}$$

This demonstrates how HT-estimator (the di-based estimator) can be written as a regression estimator:  $Y_i \sim \pi_i$  regression estimator  $\leftrightarrow$  HT estimator

Thus, if the model captures the data well, the HT estimator works.

If the model doesn't match with the data at all, the regression estimator would be a bad estimator, so will HT estimator.

## Model-based inference. Point II

The discussion above, whether it's GRER or regression estimator in general, or the model for  $H$  estimator, mainly deals with weights.

Recall that weights are survey design are not the same thing! The same weights can result from many different designs.

Model-based inference is about incorporating Survey-design (rather than just weights) into population inference.

Take stratification as an example: In the above discussion, the estimator  $\hat{T}_w$  or  $\hat{T}_k$  is an estimator for the population total  $T$ . There is only one overall  $T$  for the entire population. Following the perspective of stratification, we could allow  $T_h$  for each stratum and then sum up different strata to get to the population estimate. (Rather than having one random variable, have  $H$  random variables)

2 Stratification — different strata have different mean and variance, and thus different models  
 (In the previous module) Under a design-based approach

$$\bar{Y} = \sum_{h=1}^H \frac{N_h}{N} \cdot \bar{Y}_h \leftarrow \bar{Y}_h / N_h \text{ is presented as a true value, which } \bar{y}_h \text{ tries to capture.}$$

$$\hat{\bar{Y}} = \sum_{h=1}^H \frac{N_h}{N} \hat{y}_h$$

Under a Bayesian model-based approach:  $\bar{Y}_h / N_h$  is framed as a random variable (btw this is a new element comparing to GRER etc. discussion above. Model doesn't have to be Bayesian, but this is how it is discussed here)

$$y_{hi} | \mu_h, \sigma_h^2 \sim N(\mu_h, \sigma_h^2) + \text{flat prior } p(\mu_h, \sigma_h^2) \propto 1$$

posterior distribution  $\mu_h | y_h \sim t_{n_h-1}(\bar{y}_h, \frac{s_h^2}{n_h})$  where  $s_h^2$  is the sample variance

[Where are "models"?

Recall that "model" essentially means assumption of distributions. So all these distributions, from prior, to  $y_{hi}$ , to posterior, are models.]

We already saw this in the previous module.

Utilizing the posterior distribution of  $\mu_h | y_h$ , can now calculate  $E(\bar{Y} | y, z)$

1) given the posterior distribution  $\mu_h | y_h \rightarrow \bar{Y}_{h, ns} \sim t_{n_h-1}(\bar{y}_h, (1-f_h) \frac{s_h^2}{n_h})$  \*

Draw  $\hat{\bar{Y}}_{h, ns}$  values

2) compute  $\hat{\bar{Y}}_h = \frac{n_h}{N_h} \bar{y}_h + \frac{N_h - n_h}{N_h} \hat{\bar{Y}}_{h, ns}$

3) repeat 1000 times  $\Rightarrow E(\hat{\bar{Y}}_h)$

4)  $\bar{\bar{Y}} = \sum_{h=1}^H P_h \cdot \bar{Y}_h \approx \sum_{h=1}^H P_h E(\hat{\bar{Y}}_h)$

Essentially one model for each stratum. That's incorporating design (strata) into model-based inference.

## 2 Clustering — random effect model

$$y_{ki} | \mu, \sigma_k^2, b_k \sim N(\mu + b_k, \sigma_k^2)$$

$$b_k \sim N(0, \tau^2)$$

Assuming that  $\sigma_k^2$  and  $\tau^2$  are known

- the posterior distribution of  $\bar{Y}_{k,ns} \sim N(\tilde{y}_k, v_k^2)$

where for the sampled clusters :

$$\tilde{y}_k = \frac{n_k \cdot \tau^2}{n_k \tau^2 + \sigma_k^2} \cdot \bar{y}_k + \left(1 - \frac{n_k \tau^2}{n_k \tau^2 + \sigma_k^2}\right) \cdot \boxed{\frac{\sum_{k=1}^K \frac{n_k}{n_k \tau^2 + \sigma_k^2} \bar{y}_k}{\sum_{k=1}^K \frac{n_k}{n_k \tau^2 + \sigma_k^2}}} \quad \tilde{y}$$

$$v_k = \left(1 - \frac{n_k}{N_k}\right) \left(\sum_{k=1}^K \frac{n_k}{n_k \tau^2 + \sigma_k^2}\right)^{-1}$$

and for the non-sampled clusters :

$$\tilde{y}_k = \tilde{y} \quad \text{and} \quad v_k = \left(\sum_{k=1}^K \frac{n_k}{n_k \tau^2 + \sigma_k^2}\right)^{-1}$$

can draw  $\bar{Y}_{k,ns}$  from  $N(\tilde{y}_k, v_k^2)$

- compute  $\hat{Y}$  from < for sampled cluster:  $n_k$  sampled, mean  $\bar{y}_k$  +  $(N_k - n_k)$  non-sampled, mean  $\bar{Y}_{k,ns}$   
for non-sampled cluster: all units  $N_k$  non-sampled  $\bar{Y}_{k,ns}$

$$\hat{Y} = \frac{1}{N} \left[ \sum_{k=1}^K [n_k \cdot \bar{y}_k + (N_k - n_k) \bar{Y}_{k,ns}] + \sum_{k=k+1}^M N_k \cdot \bar{Y}_{k,ns} \right]$$

- Repeat  $\Rightarrow E(\hat{Y})$

Recall multilevel analysis :

$$y_{ki} = \beta_0 + \epsilon_{ki} \quad \epsilon_{ki} \sim N(0, \sigma^2)$$

$$\beta_{0k} = \beta_0 + u_{0k} \quad u_{0k} \sim N(0, \sigma_0^2)$$

As for why  
it's like  
this, not  
shown here

## 2 Unequal probability of selection.

### a. If categories of selection probabilities.

This can be viewed as a stratification problem: Each stratum refer to the elements of the population with equal probability of selection.

Intuitive example: if men and women have different selection probabilities, then

$$\frac{n_g}{n} \neq \frac{N_g}{N} \text{ and } \frac{n_d}{n} \neq \frac{N_d}{N}.$$

A way out of this would be to calculate  $\bar{y}_g$  and  $\bar{y}_d$  separately, then combine the two means proportional to their population portion  $\Rightarrow$  stratification.

$$\begin{aligned}\bar{E}(\bar{Y}|y, z) &= \bar{E}\left(\sum_{h=1}^H p_h \cdot \bar{Y}_h | y, z\right) \\ &= \sum_{h=1}^H p_h \underbrace{\bar{E}(\bar{Y}_h | y, z)}_{\left(\text{since } \bar{E}(\bar{Y}_h | y, z) = \bar{y}_h\right)} \\ &= \sum_{h=1}^H p_h \cdot \bar{y}_h \\ &= \sum_{h=1}^H \frac{n_h}{N} \cdot \frac{\sum_i y_i}{n_h} \\ &= \frac{1}{N} \sum_{h=1}^H \frac{n_h}{n_h} \sum_i y_i\end{aligned}$$

What's "model" about this?

$\bar{E}(\bar{Y}_h | y, z)$  exists bc  $\bar{Y}_h | y, z$  follows a (assumed) distribution. This is model!

Especially if taken a Bayesian perspective, to know  $\bar{Y}_h | y, z$ , we need  $\bar{Y}_{h, ns} | y, z$ , which depends on  $\mu_h | y, z$  (posterior distribution)

### b. If rather than categories of selection probabilities, each individual has an unique selection probability:

Recall that HT-estimator can be expressed as a regression estimator based on the regression model between  $\pi_i$  and  $y_i$

$$y_i \sim N(\pi_i \cdot \beta, \pi_i^2 \sigma^2)$$

It's possible to generalize this simple linear model between  $\pi_i$  and  $y_i$  to, e.g.,  $\pi_i \beta_1 + \beta_0$  or  $\pi_i \beta_1 + \pi_i^2 \beta_2 + \beta_0$  (of course, the estimator would not be HT estimator)

$$\text{In general: } y_i \sim N(g(\pi_i), \pi_i^k \sigma^2)$$

How to predict  $\bar{Y}$ :

1) Given a sample of data on  $y_i$  and  $\pi_i$ , model relationship between  $y_i$  and  $\pi_i$ , for example,  $y_i = \beta_0 + \pi_i \beta_1 + \pi_i^2 \beta_2$

$$\begin{aligned}\downarrow \text{Based on the formula of linear regression } y_i &\sim N(\vec{x}_i \cdot \vec{\beta}, \sigma^2) \\ \hat{\beta} &= (\vec{x}^T \vec{x})^{-1} \vec{x}^T \cdot y \text{ and } SE(\hat{\beta}) = \sigma^2 (\vec{x}^T \vec{x})^{-1}\end{aligned}$$

$$\beta_1 \sim N(\beta_1, SE(\beta_1)^2)$$

same for  $\beta_2$  and  $\beta_0$

2) draw  $\hat{\beta}$  from these estimated distribution of  $\beta$ 's

3)  $\pi_i$  is known for all units in the population.

$$\text{for all units in population } \hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 \cdot \pi_i + \hat{\beta}_2 \cdot \pi_i^2$$

4) Sum up all  $\hat{Y}_i$  for unsampled units and  $y$  sample  $\Rightarrow \bar{Y}$

Selection probability

- Comparing the situation a. and b. above. Notice how the information of  $\pi_i$  is used completely differently. In a.,  $\pi_i$  is used for constructing strata, but  $\pi_i$  is not incorporated in the "models" directly. In contrast, in b.,  $\pi_i$  are at the individual level and  $\pi_i$  is directly used for building the "model" for  $Y$

### Point III.

The above discussions are all about situation 1), in which (descriptive) statistics of the population is the goal. E.g., population mean or total — a scalar quantity. Now, consider situation 2), where the goal of inference is model.

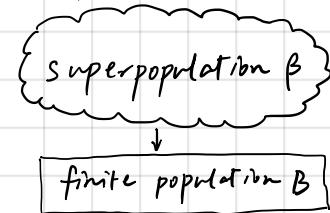
Consider a simple example of no intercept linear model.

$$Y_i = \beta \cdot X_i + \epsilon_i \quad \epsilon_i \sim N(0, \sigma^2)$$

In terms of finite population inference, we can fit the entire population data to the linear model to obtain the population slope.

$$B = \frac{\sum_{i=1}^N Y_i X_i}{\sum_{i=1}^N X_i^2}$$

*why  
Y<sub>i</sub> is a random variable?*



- Note that  $B$  would be an unbiased estimate of  $\beta$  if the modelled relationship between

$Y_i$  and  $X_i$  is the true relationship:

In superpopulation —  $Y_i = \beta \cdot X_i + \epsilon_i \quad \epsilon_i \sim N(0, \sigma^2)$

$$E(Y_i) = \beta X_i$$

$$B = \frac{\sum_{i=1}^N Y_i X_i}{\sum_{i=1}^N X_i^2} \quad E(B) = E\left(\frac{\sum_{i=1}^N Y_i X_i}{\sum_{i=1}^N X_i^2}\right) = \frac{\sum_{i=1}^N E(Y_i) X_i}{\sum_{i=1}^N X_i^2} = \frac{\sum_{i=1}^N X_i^2 \beta}{\sum_{i=1}^N X_i^2} = \beta$$

- However, if the model is misspecified, say, in the superpopulation the true relationship is

$$Y_i = \beta X_i^2 + \epsilon_i \quad \epsilon_i \sim N(0, \sigma^2)$$

$$E(Y_i) = \beta X_i^2$$

$$E(B) = E\left(\frac{\sum_{i=1}^N Y_i X_i}{\sum_{i=1}^N X_i^2}\right) = \frac{\sum_{i=1}^N E(Y_i) X_i}{\sum_{i=1}^N X_i^2} = \frac{\sum_{i=1}^N X_i^3}{\sum_{i=1}^N X_i^2} \beta$$

Then  $B$  no longer estimates the true parameter governing the quadratic relationship. *but the best linear approximation (bc residuals will sum to 0 ← regression's definition)*

Now we consider the realistic situation. Rather than dealing with how  $B$  infers  $\beta$ , we consider how to infer  $B$  from a sample. We have no information of population, but only a sample, but weights are available.

$$\sum_{i=1}^n w_i (y_i - \hat{B}_w X_i) X_i = 0 \quad \therefore \hat{B}_w = \frac{\sum_{i=1}^n w_i y_i X_i}{\sum_{i=1}^n w_i X_i^2}$$

$\hat{B}_w$  is an unbiased estimation of  $B$ :

*if  $w_i$  is independent from  $Y|X$ , then whether  $w_i$  is included here doesn't matter?*

If the model is correct and sampling is independent of the conditional distribution of  $Y|X$ , then we do not necessarily need to consider  $w_i$ . The unweighted estimate will already be consistent for  $B$ :

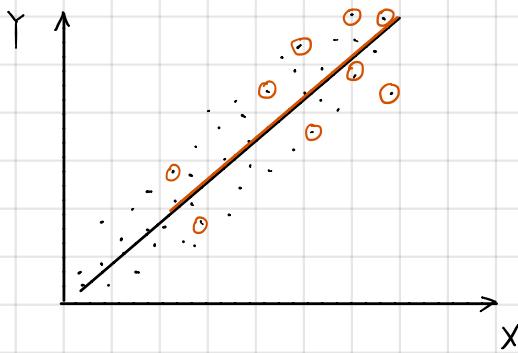
$$\hat{B}_{uw} = \frac{\sum_{i=1}^n y_i X_i}{\sum_{i=1}^n X_i^2} = \frac{\sum_{i=1}^n z_i Y_i X_i}{\sum_{i=1}^n z_i X_i^2}$$

$$E(\hat{B}_{uw}) = \frac{\sum_{i=1}^n z_i E(Y_i) X_i}{\sum_{i=1}^n z_i X_i^2} = \frac{\sum_{i=1}^n z_i X_i B X_i}{\sum_{i=1}^n z_i X_i^2} = B \frac{\sum_{i=1}^n z_i X_i^2}{\sum_{i=1}^n z_i X_i^2}$$

\* why need  $E(\hat{B}_{uw})$ ?

Continuing the previous page.

An intuitive example about when we need to consider  $w_i$  for  $\hat{B}$

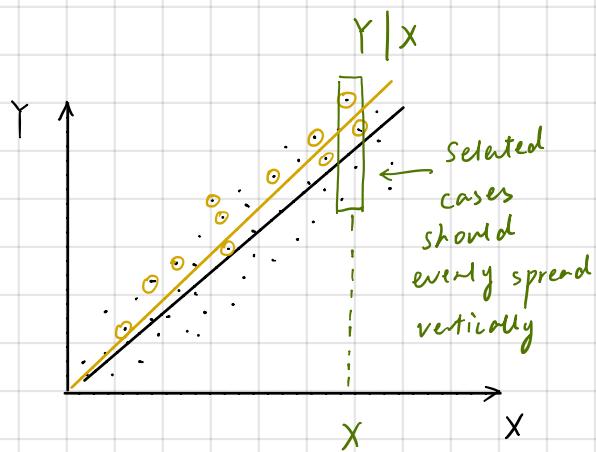


In the population, the true relationship between  $X$  and  $Y$  is like this ( $B$  captures  $\beta$ )

If the selection mechanism  $I$  ( $\rightarrow w_i$ ) is skewed, and that we ended up with these orange cases, that's okay!

→ the estimated  $\hat{B}$  will be very similar to  $B$ .

That's why we can ignore the sampling mechanism (i.e., weights) when estimating coefficient.



The above discussion simplifies the situation a little bit.

Can we unconditionally ignore the sampling mechanism? The answer is no, we can ignore sampling mechanism (weight) only if  $I \perp Y | X$ .

That is to say, it's okay if  $I$  is related to  $Y$ , like the example above, cases with higher valued  $Y$  are selected more. But it's not okay if  $I$  is related to  $Y | X$ , see example in the left:

If  $I$  is related to  $Y | X$  (meaning that the selection doesn't spread evenly around the line in the vertical direction) and we ended up with the yellow cases,

→ the estimated  $\hat{B}$  will be biased.  
(no intercept model)

Thus, 2 conditions need to be satisfied for us to ignore  $w_i$  in estimating  $B$ .

1) the model is correct (See previous page)

2) the sampling mechanism is independent of  $Y | X$  (demonstrated above)