

# STAT 344 MT Review

1. Target population : "All"
2. sampling :  $\bar{y}$  in sampling method giving it  $\bar{y}$  / total

{ simple random sample (SRS)  
 stratified sampling

3. observational unit : individual "i"

4. variable of interest : Response variable
  - Binary variable -  $\hat{p} \rightarrow p$
  - numeric variable -  $\bar{x}_s \rightarrow \bar{x}_p$   
 $t_s \rightarrow t_p$

5. confidence interval

$$(\text{estimator}) \pm \text{critical value} \times \text{S.E.}(\text{estimator})$$

$$Z_{1-\frac{\alpha}{2}}$$

For mean,  $\bar{y}_s \pm$   
 $\uparrow$   
 sample mean

critical value  
 $Z_{1-\frac{\alpha}{2}}$

$$\frac{s_s^2}{n} \rightarrow \begin{array}{l} \text{sample variance} \\ \text{sample size} \end{array}$$

$$SE(\bar{y}_s)$$

ME (Margin of Error)

For proportion,  $\hat{p}_s \pm$   
 $\uparrow$   
 sample proportion

critical value  
 $Z_{1-\frac{\alpha}{2}}$

$$\frac{\hat{p}_s(1-\hat{p}_s)}{n}$$

$$SE(\hat{p}_s)$$

ME

i) confidence interval width / length =  $2ME$

ii) fixed confidence level  $c\% (1-\alpha\%)$ ,  $Z_{1-\frac{\alpha}{2}}$  fixed

→ sample size  $n \uparrow$ , S.E. (estimator) ↓, ME ↓, width ↓

iii) fixed sample size  $n$ ,

→ confidence level  $c\% \uparrow (90\% \rightarrow 95\%)$

→  $Z_{1-\frac{\alpha}{2}} \uparrow$ , ME ↑, width ↑

iv) ME or width → find / determine  $n$

↓  
within-error      population SD (population guessed SD)

For mean,  $ME = Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \leq ME$

$$\therefore n \geq \frac{(Z_{1-\frac{\alpha}{2}} \sigma)^2}{ME^2}$$

round up      e.g.  $n \geq 22.3$   
 $n \geq 23$

For proportion,  $Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{p}_s(1-\hat{p}_s)}{n}} \leq ME$

$\hat{p}_s$ : guessed proportion

If no info given

worst case:  $\hat{p}_s = 0.5$

$\hat{p}(1-\hat{p})$  值  $\hat{p}^2$  大

$$n \geq \frac{(Z_{1-\frac{\alpha}{2}})^2 \hat{p}_s (1-\hat{p}_s)}{ME^2}$$

round up

b.	population		sample	
	Numeric	Binary	Numeric	Binary
size	$N$	$N$	$n$	$n$
mean	$\mu(\bar{y}_p)$	$P_p$	$\bar{y}_s$	$\hat{P}_s$
sd	$\sigma$	$P_p(1-P_p)$	$s$	$\hat{P}_s(1-\hat{P}_s)$

$$E(\bar{y}_s) = \bar{y}_p$$

$$\text{Var}(\bar{y}_s) = \frac{\sigma^2}{n} \quad SD(\bar{y}_s) = \frac{\sigma}{\sqrt{n}}$$

if  $\sigma^2$  unknown,  $s_s^2$  replace  
sample variance

$$\text{Var}(\bar{y}_s) = \frac{s_s^2}{n} \quad \text{SD}(\bar{y}_s) = \frac{s_s}{\sqrt{n}}$$

$$E(\hat{P}_s) = P_p$$

$$\text{Var}(\hat{P}_s) = \frac{P_p(1-P_p)}{n} \quad \text{SD}(\hat{P}_s) = \sqrt{\frac{P_p(1-P_p)}{n}}$$

$$P_p \text{ unknown, } \hat{P}_s \text{ replace: } SE(\hat{P}_s) = \sqrt{\frac{\hat{P}_s(1-\hat{P}_s)}{n}}$$

$$\text{Total: } t_p = \sum_{i=1}^N y_i \quad \hat{t}_p = N \bar{y}_s$$

$$\text{Var}(\hat{t}_p) = N^2 \text{Var}(\bar{y}_s) = N^2 \left( \frac{s_s^2}{n} \right)$$

$$SE(\hat{t}_p) = N \sqrt{\frac{s_s^2}{n}}$$

## Part I : Simple Random Sample (SRS)

↳ each unit in population selected as a sample with

$$\underline{\text{same}} \text{ probability} = \frac{1}{\binom{N}{n}}$$

↳ total ways to collect sample =  $\binom{N}{n}$

### ① Estimate $\bar{y}_p$ [population mean]

$\bar{y}_s$  estimate  $\bar{y}_p$

$$\textcircled{i} \quad \bar{y}_s = \frac{\sum_{i=1}^n y_i}{n}$$

$$\textcircled{ii} \quad \text{se}(\bar{y}_s) = \sqrt{\frac{s_s^2}{n}} \quad \text{sample variance} \quad s_s^2 = \frac{\sum_{i=1}^n (y_i - \bar{y}_s)^2}{n-1}$$

### ② Estimate $p_p$ (population proportion)

$$\textcircled{i} \quad \hat{p}_s = \frac{\text{特征数} \frac{n}{N}}{n} \quad \text{known}$$

$$\textcircled{ii} \quad \text{se}(\hat{p}_s) = \sqrt{\frac{\hat{p}_s(1-\hat{p}_s)}{n}}$$

### ③ Finite population correction (FPC)

$$N = \text{finite} \quad (\text{Not very large}) \quad \left. \right\} \text{need FPC}$$

$$\frac{n}{N} > 0.05$$

$$N \text{ very large / infinite} \quad \left. \right\} \text{No need FPC.}$$

$$\frac{n}{N} \leq 0.05$$

tip: in exam,  $\frac{n}{N} = 0.05$  建议 FPC.

By FPC

$$\text{Var}(\bar{y}_s) = \left(1 - \frac{n}{N}\right) \frac{s_s^2}{n}$$

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

$$\text{Var}(\hat{t}_s) = N^2 \left(1 - \frac{n}{N}\right) \frac{s_s^2}{n}$$

$$se(\bar{y}_s) = \sqrt{\left(1 - \frac{n}{N}\right) \frac{s_s^2}{n}}$$

$$se(\hat{p}_s) = \sqrt{\left(1 - \frac{n}{N}\right) \frac{\hat{p}_s(1-\hat{p}_s)}{n}}$$

$$se(\hat{t}_s) = N \sqrt{\left(1 - \frac{n}{N}\right) \frac{s_s^2}{n}}$$

Part II: study planning  $\rightarrow$  determine sample size  $n$

ME  $\rightarrow$  determine sample size

But FPC

Input: how narrow of a confidence interval  $\rightarrow$  width  
within — error  $\rightarrow$  ME

Output: sample size  $n$  ( $\hat{p}_s$  or  $\hat{t}_s$ )

$$ME = \delta = Z_{1-\frac{\alpha}{2}} \sqrt{\left(1 - \frac{n}{N}\right) \frac{s_{\text{guess}}^2}{n}} \quad (\text{FPC})$$

$$\delta = Z_{1-\frac{\alpha}{2}} \sqrt{\left(1 - \frac{n}{N}\right) \frac{\hat{p}_{\text{guess}}(1-\hat{p}_{\text{guess}})}{n}} \quad (\text{FPC})$$

without FPC

$$ME = \delta = Z_{1-\frac{\alpha}{2}} \sqrt{\frac{s_{\text{guess}}^2}{n'}}$$

$$n' = \frac{(Z_{1-\frac{\alpha}{2}})^2 s_{\text{guess}}^2}{\delta^2}$$

since  $ME = \bar{ME}$

$$\cancel{Z_{1-\frac{\alpha}{2}}} \sqrt{\left(1 - \frac{n}{N}\right) \frac{s_{\text{guess}}^2}{n}} = \cancel{Z_{1-\frac{\alpha}{2}}} \sqrt{\frac{s_{\text{guess}}^2}{n'}}$$

$$\left(1 - \frac{n}{N}\right) \frac{s_{\text{guess}}^2}{n} = \frac{s_{\text{guess}}^2}{n'}$$

$$\therefore \frac{1 - \hat{p}}{n} = \frac{1}{n'}$$

$$\therefore \frac{1}{n} - \frac{1}{N} = \frac{1}{n'}$$

$$\therefore \frac{N-n}{nN} = \frac{1}{n'}$$

$$\therefore n'(N-n) = nN$$

$$\therefore n'N - n'n - nN = 0$$

$$\therefore n'n + nN = n'N$$

$$\therefore n = \frac{n'N}{n' + N}$$

Tip: ME  $\rightarrow$   $n$  with FPC

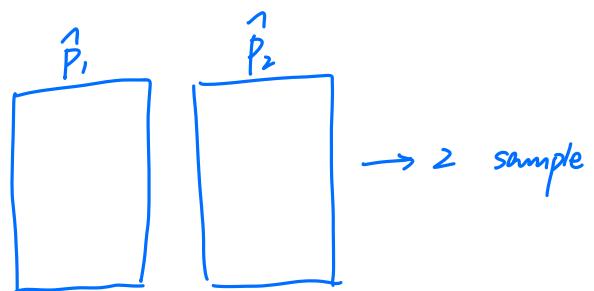
先求  $n'$  with No FPC

$$\Rightarrow n = \frac{n'N}{n' + N} \Rightarrow n \text{ with FPC}$$

Part III: proportion change "Δ" panel data

a. 2 independent sample

$$SE(\hat{p}_1 - \hat{p}_2) = \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$



b. - 2 time period

- same observational units in sample

e.g. 同一病人在 7 月和 8 月  
分别做两次检查

$\therefore$  samples are dependent

		Current month	month
		NO	YES
previous month	NO	$P_{NN}$	$P_{NY}$
	YES	$P_{YN}$	$P_{YY}$

proportion change :  $P(\text{Yes, current}) - P(\text{Yes, previous})$   
 $= (P_{NY} + P_{YY}) - (P_{YN} + P_{NN})$

$\Delta = P_{NY} - P_{YN}$

$$\begin{aligned} \text{se}(\Delta) &= \text{se}(P_{NY} - P_{YN}) \\ &= \sqrt{\frac{\hat{P}_{NY}(1-\hat{P}_{NY})}{n} + \frac{\hat{P}_{YN}(1-\hat{P}_{YN})}{N} - 2 \frac{\hat{P}_{NY}\hat{P}_{YN}}{n}} \\ &\quad \text{Var}(\hat{P}_{NY}) \quad \text{Var}(\hat{P}_{YN}) \quad -2 \text{cov}(\hat{P}_{NY}, \hat{P}_{YN}) \\ &\quad \downarrow \\ &\quad \text{cov negative} \end{aligned}$$

Part IV instrument (随机抽样方法)

	unit	$y$	$Z_i$	
population	1	$y_1$	1	$Z_i = \begin{cases} 0 & i \notin \text{sample} \\ 1 & i \in \text{sample} \end{cases}$
Data	2	$y_2$	0	被选中1, 没选中0
	:	:	:	
	$N$	$y_N$	:	$Z_i \sim \text{Bernoulli}$

$$P(Z_i = 1) = \frac{\binom{N-1}{n-1}}{\binom{N}{n}} = \frac{n}{N} = P(Z_j = 1)$$

$$P(z_i z_j = 1) = \frac{\binom{N-2}{n-2}}{\binom{N}{n}} = \frac{n(n-1)}{N(N-1)}$$

$$z_i \quad 0 \quad 1$$

$$E(z_i) = 0 \cdot (1 - \frac{n}{N}) + 1 \cdot \frac{n}{N} = \frac{n}{N} \quad P(z_i) = 1 - \frac{n}{N} = \frac{n}{N}$$

$$\text{Var}(z_i) = P(1-P) = \frac{n}{N} \left(1 - \frac{n}{N}\right)$$

$$\text{cov}(z_i, z_j) = E(z_i z_j) - E(z_i)E(z_j)$$

$$= 0 \times 0 \times P(z_i=0, z_j=0) + 0 \times 1 \times P(z_i=0, z_j=1) + 1 \times 0 \times P(z_i=1, z_j=0) \\ + 1 \times 1 \times P(z_i=1, z_j=1) - E(z_i)E(z_j)$$

$$= P(z_i z_j = 1) - E(z_i)E(z_j) \\ = \frac{n(n-1)}{N(N-1)} - \left(\frac{n}{N}\right)^2$$

$$\bar{y}_s = \frac{\sum_{i=1}^n y_i z_i}{n}$$

$$\text{Var}(\bar{y}_s) = \text{Var}\left(\frac{\sum_{i=1}^n y_i z_i}{n}\right) = \left(\frac{1}{n}\right)^2 \text{Var}\left(\sum_{i=1}^n y_i z_i\right) \\ = \left(\frac{1}{n}\right)^2 \text{Var}(y_1 z_1 + \dots + y_N z_N) \\ = \left(\frac{1}{n}\right)^2 \left( \text{Var}(y_1 z_1) + \dots + \text{Var}(y_N z_N) + 2 \sum \text{cov}(y_i z_i, y_j z_j) \right) \\ = \left(\frac{1}{n}\right)^2 (n s_p^2) + 2 \left( \sum \text{cov}(y_i z_i, y_j z_j) \right)$$

$y_i$  independent  
 $\text{cov}(z_i, z_j)$

$$\text{Instrument } z_i = \left(1 - \frac{n}{N}\right) \frac{s_p^2}{n}$$

$\Rightarrow$  估计 FPC Variance

更准确

X

## Part V : Estimation $\rightarrow$ under SRS

[ mean  $\bar{y}_p$  estimate ]

### a. Vanilla

$$\bar{y}_{s, \text{vanilla}} = \frac{\sum_{i=1}^n y_i}{n} = \bar{y}_s$$

$$se(\bar{y}_{s, \text{vanilla}}) = \sqrt{\frac{s_s^2}{n}} \quad \text{or} \quad \sqrt{(1 - \frac{n}{N}) \frac{s_s^2}{n}}$$

No FPC FPC

$$\hat{t}_s = N \bar{y}_{s, \text{vanilla}}$$

$$se(\hat{t}_s) = N \sqrt{\frac{s_s^2}{n}} \quad \text{or} \quad N \sqrt{(1 - \frac{n}{N}) \frac{s_s^2}{n}}$$

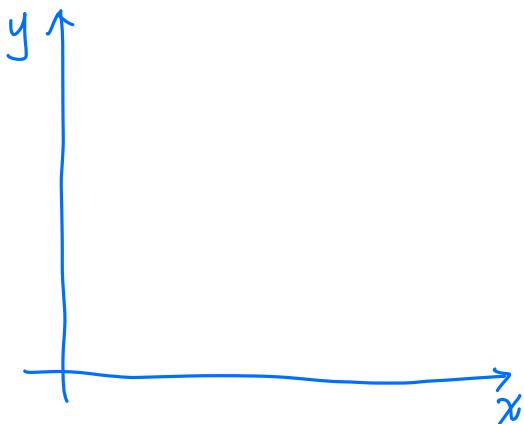
No FPC FPC

### b. Ratio Estimation ( $x$ and $y$ correlated)

$$\bar{y}_{s, \text{ratio}} = \left( \frac{\bar{y}_s}{\bar{x}_s} \right) \bar{x}_p$$

$B = \frac{\bar{y}_p}{\bar{x}_p}$

$\hookrightarrow$  ratio estimate  $\hat{B}$



$$\hat{y}_i = \left( \frac{\bar{y}_s}{\bar{x}_s} \right) x_i = \hat{B} x_i$$

$$e_i = y_i - \hat{y}_i = y - \hat{B} x_i$$

$$\hookrightarrow \text{Var}(e_i) = \frac{\sum_{i=1}^n (e_i - \bar{e})^2}{n-1} \quad \text{where } \bar{e} = \frac{\sum_{i=1}^n e_i}{n}$$

$$\text{Var}(\hat{B}) = \frac{1}{\bar{x}_p^2} \text{Var}(\bar{e}) = \frac{1}{\bar{x}_p^2} \frac{s_e^2}{n} \quad \text{No FPC}$$

$$= \frac{1}{\bar{x}_p^2} (1 - \frac{n}{N}) \frac{s_e^2}{n} \quad \text{FPC}$$

$$\begin{aligned}\text{Var}(\bar{y}_{s,\text{ratio}}) &= \text{Var}(\hat{\beta} \bar{x}_p) = \bar{x}_p^2 \text{Var}(\hat{\beta}) \\ &= \frac{s_e^2}{n} \quad \text{No FPC} \\ &= \left(1 - \frac{n}{N}\right) \frac{s_e^2}{n} \quad \text{FPC}\end{aligned}$$

$$se(\bar{y}_{s,\text{ratio}}) = \sqrt{\frac{s_e^2}{n}} \quad \text{No FPC}$$

$$= \sqrt{\left(1 - \frac{n}{N}\right) \frac{s_e^2}{n}} \quad \text{FPC}$$

$$\hat{t}_{s,\text{ratio}} = N \bar{y}_{s,\text{ratio}}$$

$$\begin{aligned}se(\hat{t}_{s,\text{ratio}}) &= N \sqrt{\frac{s_e^2}{n}} \quad \text{No FPC} \\ &= N \sqrt{\left(1 - \frac{n}{N}\right) \frac{s_e^2}{n}} \quad \text{FPC}\end{aligned}$$



### c. Regression estimate

$$\begin{array}{lll}\hat{y} = b_0 + b_1 x & b_0 = \bar{y} - b_1 \bar{x} & (\bar{x}, \bar{y}) \text{ on the regression line} \\ \uparrow \text{estimate} & b_1 = r \frac{s_y}{s_x} & \end{array}$$

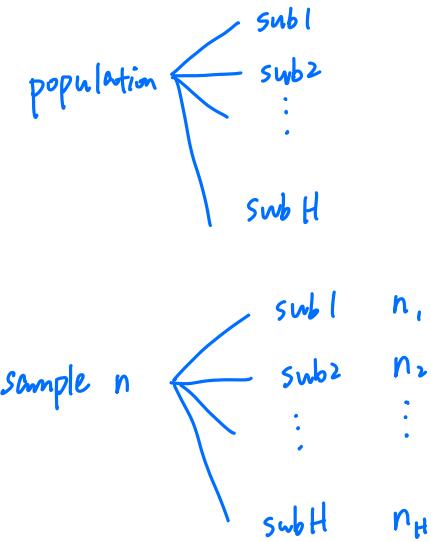
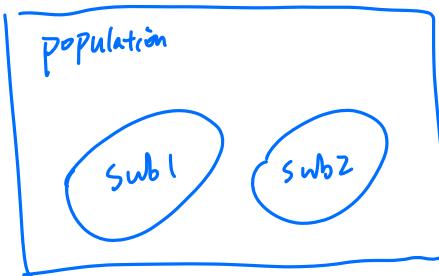
$$\text{residual: } e_i = y_i - \hat{y}_i$$

$$\text{Var}(b_1) = \frac{\sum_{i=1}^n e_i^2}{(n-2) \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\bar{y}_{s,\text{regression}} = b_0 + b_1 \bar{x}_p \quad ???$$

## Part VI: Estimation in Domain (subpopulation)

sample method 依 SRS  
但是計算各 subpopulation



First,  $\bar{y}_p$  population mean

Second,  $\bar{y}_p$  in subpopulation  $\bar{y}_{p1}, \bar{y}_{p2}, \dots, \bar{y}_{pH}$

$$\bar{y}_{s, \text{domain}} = \frac{1}{n_d} \sum_{i=1}^{n_d} y_i \quad \text{estimate} \quad \hat{t}_{s, \text{domain}} = \frac{1}{N_d} \sum_{i=1}^{N_d} y_i$$

$$\hat{t}_{s, \text{domain}} = N_d \bar{y}_{s, \text{domain}}$$

a. known domain size  $N_d$

$$\hat{t}_{s, \text{domain}} = N_d \bar{y}_{s, \text{domain}}$$

$$SE(\hat{t}_{s, \text{domain}}) = N_d \sqrt{\frac{s_{sd}^2}{n_d}} = N_d \sqrt{\left(1 - \frac{n_d}{N_d}\right) \frac{s_{sd}^2}{n_d}}$$

No FPC                          FPC.

b. unknown domain size  $N_d$ , only know population size  $N$

$$u_i = \begin{cases} y_i, & i \in \text{sample} \\ 0, & i \notin \text{sample} \end{cases}$$

$$\hat{t}_{s, \text{domain}} = N \bar{u}_s$$

$$SE(\hat{t}_{s, \text{domain}}) = N \sqrt{\frac{s_u^2}{n}} \quad \text{No FPC}$$

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

C. ratio estimate

$$u_i = \begin{cases} y_i, & i \in \text{sample} \\ 0, & i \notin \text{sample} \end{cases}$$

$$x_i = \begin{cases} 1, & i \in P_d \\ 0, & i \notin P_d \end{cases}$$

$$\bar{y}_{s, \text{domain}} = \frac{\bar{u}_s}{\bar{x}_s}$$

$$\hat{t}_{s, \text{domain}} = N_d \bar{y}_{s, \text{domain}}$$

 Part VII: Stratified Sampling  $\rightarrow$  each group randomly as sample

$$\text{population size } N = N_1 + N_2 + \dots + N_H$$

$$N_1 \dots N_H \text{ subpopulation size}$$

$$H: \text{the number of subpopulations}$$

$$\text{sample size } n = n_1 + n_2 + \dots + n_H$$

$$n_1, n_2 \dots n_H \text{ subpopulation size}$$

$$\bar{y}_p = \frac{1}{N} \sum_{i=1}^N y_i = \frac{\sum_{h=1}^H n_h \bar{y}_{ph}}{N} = \sum_{h=1}^H \left( \frac{n_h}{N} \right) \bar{y}_{ph}$$

subpopulation 1 2 3 ... H

$$\text{population size } N_1, N_2, N_3 \dots N_H$$

$$\text{sample size } n_1, n_2, n_3 \dots n_H$$

$$\text{sample mean } \bar{y}_{s1}, \bar{y}_{s2}, \dots, \bar{y}_{sH}$$

$$\text{sample variance } s_{s1}^2, s_{s2}^2, \dots, s_{sH}^2$$

$$\hat{t}_{str} = N \bar{y}_{str}$$

$$\hat{P}_{str} = \sum_{h=1}^H \left( \frac{n_h}{N} \right) \hat{P}_{sh}$$

h<sup>th</sup> subpopulation  
sample proportion

$$- \text{Var}(\bar{y}_{\text{str}}) = \text{Var}\left(\sum_{h=1}^H \left(\frac{N_h}{N}\right) \bar{y}_{sh}\right) = \sum_{h=1}^H \left(\frac{N_h}{N}\right)^2 \text{Var}(\bar{y}_{sh})$$

↑  
constant

$$SE(\bar{y}_{\text{str}}) = \sqrt{\sum_{h=1}^H \left(\frac{N_h}{N}\right)^2 \frac{s_{sh}^2}{n_h}} \quad \text{No FPC}$$

$$= \sqrt{\sum_{h=1}^H \left(\frac{N_h}{N}\right)^2 \left(1 - \frac{n_h}{N_h}\right) \frac{s_{sh}^2}{n_h}} \quad \text{FPC}$$

$$SE(\hat{p}_{\text{str}}) = \sqrt{\sum_{h=1}^H \left(\frac{N_h}{N}\right)^2 \frac{\hat{p}_{sh}(1-\hat{p}_{sh})}{n_h}} \quad \text{No FPC}$$

$$= \sqrt{\sum_{h=1}^H \left(\frac{N_h}{N}\right)^2 \left(1 - \frac{n_h}{N_h}\right) \frac{\hat{p}_{sh}(1-\hat{p}_{sh})}{n_h}} \quad \text{FPC.}$$

$$SE(\hat{t}_{\text{str}}) = N \sqrt{\sum_{h=1}^H \left(\frac{N_h}{N}\right)^2 \frac{s_{sh}^2}{n_h}} \quad \text{No FPC}$$

$$= N \sqrt{\sum_{h=1}^H \left(\frac{N_h}{N}\right)^2 \left(1 - \frac{n_h}{N_h}\right) \frac{s_{sh}^2}{n_h}} \quad \text{FPC.}$$

### study planning

Allocation: ① fixed  $n$

② equal variance

③ large population (No FPC)

$$s_{1,\text{guess}}^2 = s_{2,\text{guess}}^2 = \dots = s_{H,\text{guess}}^2$$

$$SE(\bar{y}_{\text{str}}) = \sqrt{\sum_{h=1}^H \left(\frac{N_h}{N}\right)^2 \frac{s_{h,\text{guess}}^2}{n_h}}$$

1. proportional allocation:

$$\frac{N_h}{N} = \frac{n_h}{n} \quad \text{OR} \quad \frac{n}{N} = \frac{n_h}{N_h}$$

2. optimize allocation:

↓  
max or min

- ~~risk pref~~,  $\text{var}(\bar{y}_{\text{str}})$  minimize

$$SE(\bar{y}_{\text{str}})$$

money constraint  $\rightarrow$  optimize

## Part I: Study planning

proportional allocation:  $\frac{N_h}{N} = \frac{n_h}{n}$  stratified sampling

無 FPC:  $\text{Var}(\bar{y}_{\text{str}}) = \sum_{h=1}^H \left( \frac{N_h}{N} \right)^2 \frac{s_{sh}^2}{n}$

FPC:  $\text{Var}(\bar{y}_{\text{str}}) = \sum_{h=1}^H \left( \frac{N_h}{N} \right)^2 \left( 1 - \frac{n_h}{n} \right) \frac{s_{sh}^2}{n_h}$

$S_p^2 = \sum_{i=1}^N (y_i - \bar{y}_p)^2$  total variance (sum of square total)

$$S_p^2 = S_{p,w}^2 + S_{p,B}^2$$

$$(S_{p,w}^2) \quad (S_{p,B}^2)$$

$$= \underbrace{\sum_{h=1}^H \left( \frac{N_h}{N} \right)^2 S_h^2}_{\substack{\text{within group} \\ \text{sum of square}}} + \underbrace{\sum_{h=1}^H \left( \frac{N_h}{N} \right) (\bar{y}_{ph} - \bar{y}_p)^2}_{\substack{\text{between group} \\ \text{sum of square}}}$$

$\text{Var}(\bar{y}_{\text{str}}) = \left( 1 - \frac{n}{N} \right) \frac{S_{p,w}^2}{n}$  FPC stratified sampling

$\text{Var}(\bar{y}_s) = \left( 1 - \frac{n}{N} \right) \frac{S_p^2}{n}$  FPC SRS

stratified / SRS = confident Ratio in webwork

$$= \frac{\text{Var}(\bar{y}_{\text{str}})}{\text{Var}(\bar{y}_s)}$$

$$= \frac{\left( 1 - \frac{n}{N} \right) \frac{S_{p,w}^2}{n}}{\left( 1 - \frac{n}{N} \right) \frac{S_p^2}{n}} = \frac{S_{p,w}^2}{S_p^2} = \frac{S_{p,w}^2}{S_{p,w}^2 + S_{p,B}^2} = 1 - \frac{S_{p,B}^2}{S_p^2}$$

Reduction in percentage =  $\left[ 1 - \frac{S_{p,w}^2}{S_p^2} \right] \times 100\%$  in webwork

② optimize allocation:

↳ Var( $\bar{Y}_{str}$ ) minimized  $\rightarrow$  by sample size  $n$

↳ constraint for money - "cost constraint"

$$n_1 C_1 + n_2 C_2 + \dots + n_h C_h = C_{\text{total}}$$

↑  
cost per unit out subpopulation 1

$$\frac{n_h}{n} = \frac{N_h \left( \frac{S_{h, \text{guess}}}{C_h} \right)}{\sum_{k=1}^h N_k \frac{S_{k, \text{guess}}}{C_k}} \quad \text{at } h \text{ subpopulation}$$

$$\Rightarrow n_h \propto \frac{N_h S_{h, \text{guess}}}{\sqrt{C_h}} \quad \text{to } h \text{ subpopulation is } \frac{1}{\sqrt{C_h}} \text{ of sample size}$$

Example:  $N_1 = N_2 = N_3 = N_4$        $S_{1, \text{guess}} = S_{2, \text{guess}} = S_{3, \text{guess}} = S_{4, \text{guess}} = 2S_{1, \text{guess}}$

$$C_1 = C_2 = 4 \quad C_3 = C_4 = 9 \quad C = \$1560$$

$$n_1 \propto \frac{N_1 S_{1, \text{guess}}}{\sqrt{C_1}} = \frac{NS}{\sqrt{4}} = \frac{NS}{2}$$

$$n_2 \propto \frac{N_2 S_{2, \text{guess}}}{\sqrt{C_2}} = \frac{NS}{\sqrt{4}} = \frac{NS}{2}$$

$$n_3 \propto \frac{N_3 S_{3, \text{guess}}}{\sqrt{C_3}} = \frac{NS}{\sqrt{9}} = \frac{NS}{3}$$

$$n_4 \propto \frac{N_4 S_{4, \text{guess}}}{\sqrt{C_4}} = \frac{2NS}{\sqrt{9}} = \frac{2NS}{3}$$

$$n = n_1 + n_2 + n_3 + n_4 = \frac{NS}{2} + \frac{NS}{2} + \frac{NS}{3} + \frac{2NS}{3} = 2NS$$

$$\frac{n_1}{n} = \frac{\frac{NS}{2}}{2NS} = \frac{1}{4} \quad \frac{n_2}{n} = \frac{\frac{2}{2}NS}{2NS} = \frac{1}{4}$$

$$\frac{n_3}{n} = \frac{\frac{NS}{3}}{2NS} = \frac{1}{6} \quad \frac{n_4}{n} = \frac{\frac{2}{3}NS}{2NS} = \frac{1}{3}$$

$$n_1, n_2, n_3, n_4 \propto (\frac{1}{4}, \frac{1}{4}, \frac{1}{6}, \frac{1}{3})$$

$$n_1 c_1 + n_2 c_2 + n_3 c_3 + n_4 c_4 = 1560$$

$$\frac{1}{4}n(4) + \frac{1}{4}n(4) + \frac{1}{6}n(9) + \frac{1}{3}n(9) = 1560$$

$$\Rightarrow n = \frac{1560}{6.5} = 240$$

Part 2: weighted → optimal allocation by minimizing  $\text{Var}(\bar{y}_s)$

$$\bar{y}_s = w\bar{y}_{s1} + (1-w)\bar{y}_{s2} \quad w: \text{the sample size weighted for group 1}$$

$$\boxed{\text{Var}(\bar{y}_s)} = \text{Var}(w\bar{y}_{s1} + (1-w)\bar{y}_{s2}) \quad S_1 \text{ 和 } S_2 \text{ independent}$$

$$\text{minimized} = \text{Var}(w\bar{y}_{s1}) + \text{Var}((1-w)\bar{y}_{s2}) = w^2 \text{Var}(\bar{y}_{s1}) + (1-w)^2 \text{Var}(\bar{y}_{s2})$$

why?

min Var → smaller SE → narrower CI → more accurate  $\bar{y}_s$

$$\text{to minimize: } \frac{\partial \text{Var}(\bar{y}_s)}{\partial w} = 2w \text{Var}(\bar{y}_{s1}) + 2(1-w)(-1) \text{Var}(\bar{y}_{s2}) \\ = 2w \text{Var}(\bar{y}_{s1}) - 2 \text{Var}(\bar{y}_{s2}) + 2w \text{Var}(\bar{y}_{s2})$$

$$= 0$$

$$\Rightarrow w = \frac{\text{Var}(\bar{y}_{s2})}{\text{Var}(\bar{y}_{s1}) + \text{Var}(\bar{y}_{s2})}$$

↑  
group 1 sample size weighted

weighted  $\bar{y}_{s2}$  的方差是  $\text{Var}(\bar{y}_{s2})$   
而  $\bar{y}_{s1}$  的方差是  $\text{Var}(\bar{y}_{s1})$

Part 3: polls of polls  $\rightarrow$  proportion; elections

Inputs:  $\hat{p}_i \text{ SE}_i \quad i = 1, 2, \dots k$

$$\hat{P}_{\text{agr}} = \frac{\sum_{i=1}^k \left( \frac{1}{\text{SE}_i^2} \right) \hat{p}_i}{\sum_{i=1}^k \left( \frac{1}{\text{SE}_i^2} \right)}$$

$$SE(\hat{P}_{\text{agr}}) = \sqrt{\frac{1}{\sum_{i=1}^k \left( \frac{1}{\text{SE}_i^2} \right)}}$$

Random effect:

$$\hat{P}_{\text{agr}} = \frac{\sum_{i=1}^k \left( \frac{1}{\text{SE}_i^2 + \tau^2} \right) \hat{p}_i}{\sum_{i=1}^k \left( \frac{1}{\text{SE}_i^2 + \tau^2} \right)}$$

$$SE(\hat{P}_{\text{agr}}) = \sqrt{\frac{1}{\sum_{i=1}^k \left( \frac{1}{\text{SE}_i^2 + \tau^2} \right)}}$$