

NOTE 1 Census 人口普查

Target population: we wish to study. Sample: A subset of a population. Sampled population: we can study

Sample unit: household. Observation unit: individuals living in the household.

Sampling frame: (like telephone survey) a list of all numbers.

④ transcription from records

Method of data collection: ① physical observations ② personal interview ③ mail. ④ web-based. ⑤ registration

Two basic sampling designs: probability sampling. nonprobability sampling

probability sampling = known probability. can judge reliability and validity.

① simple random sampling: known equal chance. 所有人名单拿出来随便抽

② systematic random sampling: ordered. 每5个抽1个.

③ stratified random sampling: 之间差很多. 内部差不多 [male female] ← strata

④ cluster sampling: 小或 cluster. 之间差不多. 里面内容不一样. (如每个小学一个cluster), 不需要所有 cluster

nonprobability sampling: unknown. selection is based on some intuitive knowledge

① convenience samples: (pilot study) 先抽几个问问. Assumption: target population is homogeneous. similar to overall

② judgement sampling: 找专家 ③ Quota sampling ④ snowball sampling: 我问你→问你朋友→...

NOTE 2. Basic concepts of statistics.

Mean: $\bar{y} = \sum y_i P(y_i)$. Variance: $E(y-\bar{y})^2 = \sum (y_i - \bar{y})^2 P(y_i)$. S.d. $\sigma^2 = E(y-\bar{y})^2 = E(y^2) - \bar{y}^2 \Rightarrow E(y^2) = \sigma^2 + \bar{y}^2$.

Infinite population: Sample Mean: $\bar{y} = \frac{1}{n} \sum y_i$ Sample Variance: $S^2 = \frac{\sum (y_i - \bar{y})^2}{n-1}$

$$E(\bar{y}) = \bar{y}, V(\bar{y}) = \frac{\sigma^2}{n}, E(\bar{y}^2) = \frac{\sigma^2}{n} + \bar{y}^2.$$

Finite population: Population total $T = n_1 + n_2 + \dots + n_N$ mean: $\bar{n} = \frac{T}{N}$ Variance: $\sigma^2 = \frac{1}{N} \sum (n_i - \bar{n})^2$.

⇒ with replacement. Population $\{n_1, \dots, n_N\}$. Selection prob. $\{p_1, \dots, p_N\}$. Sample $\{y_1, \dots, y_n\}$. $T = n_1 + \dots + n_N$

$$\hat{T} = \frac{1}{n} \sum_{i=1}^n y_i \stackrel{(n_i \text{ is } n \text{ times in sample})}{=} \frac{1}{n} \sum_{i=1}^n a_i \frac{n_i}{p_i} \quad (a_i \text{ is } n_i \text{出现次数}) \quad E(\hat{T}) = E\left(\frac{1}{n} \sum_{i=1}^n a_i \frac{n_i}{p_i}\right) = \frac{1}{n} \sum_{i=1}^n n_i \cdot \frac{n_i}{p_i} = \sum_{i=1}^n n_i = T. \text{ (unbiased)}$$

$$V(\hat{T}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n \frac{y_i}{p_i}\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}\left(\frac{y_i}{p_i}\right) = \frac{1}{n^2} \sum_{i=1}^n \left(\frac{n_i}{p_i} - T\right)^2 \cdot p_i \quad \hat{V}(\hat{T}) = \frac{1}{n} \cdot \frac{1}{n-1} \sum_{i=1}^n \left(\frac{y_i}{p_i} - \hat{T}\right)^2 \text{ (also unbiased)}$$

Ex: {1, 2, 3, 4} S: {0.1, 0.1, 0.4, 0.4} for $n=2$.

$$\text{Sample } \{1, 2\}. \text{ Prob} = 0.1 \times 0.1 \times 2 = 0.02, \hat{T} = \frac{1}{2} \sum \frac{y_i}{p_i} = \frac{1}{2} \left(\frac{1}{0.1} + \frac{2}{0.1} \right) = 15 \quad \hat{V}(\hat{T}) = \frac{1}{2} \cdot \frac{1}{2-1} \sum \left(\frac{y_i}{p_i} - \hat{T} \right)^2 = \frac{1}{2} \times \left[\left(\frac{1}{0.1} - 15 \right)^2 + \left(\frac{2}{0.1} - 15 \right)^2 \right] = 25$$

$$E(\hat{T}) = \sum \hat{T} P = (15 \times 0.02 + \dots) = 10 \quad V(\hat{T}) = \sum (\hat{T} - E(\hat{T}))^2 P = (15 - 10)^2 \times 0.02 + \dots = 6.250.$$

⇒ without replacement $\pi_{ij}: P(\text{the } i\text{-th element in the population, } n_i \text{ is selected in the sample})$

$$\frac{\pi_{ij}}{n} = \text{average prob. the } n_i \text{ is selected across the } n \text{ draws that will occur in a sample}$$

$$\text{Ex: } \{1, 2, 3\}. \quad T=6. \quad n=2. \quad \text{then } \pi_1 = \pi_2 = \pi_3 = \frac{2}{3}. \quad \hat{T} = \frac{1}{2} \sum \frac{y_i}{\pi_i} = \frac{1}{2} \sum \frac{y_i}{\frac{2}{3}} \quad (\pi_{ij} = \sum_{j \neq i} P\{j\})$$

$$\text{Sample } \{1, 2\}. \text{ prob} = \frac{1}{3}. \quad \hat{T} = \frac{1}{2} \cdot \frac{1}{2} = 4.5 \quad E(\hat{T}) = 4.5 \times \frac{1}{3} + 6 \times \frac{1}{3} + 7.5 \times \frac{1}{3} = 6. \quad V(\hat{T}) = 4.5^2 \times \frac{1}{3} + 6^2 \times \frac{1}{3} + 7.5^2 \times \frac{1}{3} - 6^2 = 1.5 \quad (E(\hat{T}^2) - E(\hat{T})^2)$$

NOTE 2 (2) Supp- Categorical

Test Statistic for Goodness-of-Fit Tests: H_0 : the distribution is the same today as it was in 2000.

$$\chi^2_{(k, F-1)} = \sum \frac{(O_i - E_i)^2}{E_i} \text{ for } O_i: \text{observed. } E_i: \text{expected } (i=1 \dots k) \quad \frac{(252-228)^2}{228} + \dots$$

Answer: since χ^2_0 is less/larger than $\chi^2_{0.05}$. we (fail to) reject H_0 .

There's (in)sufficient evidence to conclude that the distribution for ... at the $\alpha=0.05$ level of significance

Test for Independence (chi-square test) H_0 : the variables are independent for $r=2, c=3$ degree: $(2-1)(3-1)=2$.

$$\text{expected frequency} = \frac{(\text{row total})(\text{column total})}{\text{table total}} \quad \chi^2_0 = \frac{(82-63.5808)^2}{63.5808} + \dots \approx 36.82$$

Test for Homogeneity of Proportions $H_0: p_1 = p_2 = p_3$ $\chi^2_0 = \frac{(418-475.554)^2}{475.554} + \dots$

Number of Years	Observed Counts	Expected Counts
<1	252	228
1-2	255	239
3-4	162	176
≥5	331	357

	1992	2002	2008
Yes	418 (475.554)	479 (475.554)	525 (470.892)
No	602 (544.446)	541 (544.446)	485 (539.108)

128 1020 608
(128 × 883 + 1776)

1020 1020 1010
 $1020 \times 1420 + 3050 = 475.554$

NOTE 3. Simple Random Sampling

⇒ without replacement. possible samples: $\binom{N}{n} = C_n^N$. prob of select each is $\frac{1}{C_n^N}$

$$P(y_i = u_i) = \frac{1}{C_n^N}, P(y_i = u_i, y_j = u_j) = \frac{1}{C_{n-1}^{N-1}}$$

Estimation of a population mean: Population $\{n_1, n_2, \dots, n_N\}$, mean $\bar{n} = \frac{n_1 + \dots + n_N}{N} = \frac{1}{N} \sum n_i$, variance $\sigma^2 = \frac{1}{N} \sum (n_i - \bar{n})^2$

$$\text{Sample } \{y_1, \dots, y_n\}. \quad \text{mean } \bar{y} = \frac{1}{n} \sum y_i$$

$$S^2 = \frac{1}{n-1} \sum (y_i - \bar{y})^2$$

$$E(y_i) = \sum_j \frac{1}{N} u_j = u_i. \quad \text{Var}(y_i) = \sum_j \frac{1}{N} (u_j - u_i)^2 = 6^2. \quad \text{Cov}(y_i, y_j) = E(y_i y_j) - u^2 = -\frac{6^2}{N-1} \quad \text{① negative. ② } N \uparrow \text{ then Var} \rightarrow 0.$$

$$E(\bar{y}) = u_i. \quad \text{Var}(\bar{y}) = \frac{6^2}{n} \frac{N-n}{N-1} \quad \text{③ 不放回会小些. ④ } n \rightarrow \infty, \text{Var}(\bar{y}) \rightarrow 0. \quad \text{⑤ } N \rightarrow \infty, \text{Var}(\bar{y}) \rightarrow \frac{6^2}{n}.$$

Ex. precision. $\frac{1}{n} \frac{N-n}{N-1} 6^2$, like An SRS of size 400 from population 4000. is $\frac{1}{400} \frac{4000-400}{4000-1} 6^2 = 0.0225 6^2$. S 找自己知量 (Z* = 1.96)

(Estimation of population variance). $E(S^2) = \frac{N}{N-1} 6^2$. Unbiased estimator for $\sigma^2 = \frac{N-1}{N} S^2$, for $\text{Var}(\bar{y}) = (1-f) \frac{S^2}{n}$. $f = \frac{n}{N}$ (called finite population correction)

if $n, N, N-n$ sufficiently large. $(\bar{y} - \bar{u}) / \sqrt{\text{Var}(\bar{y})} \sim N(0, 1)$ and $100(1-\alpha)\%$ C.I. for u is $B = \sum_{i=1}^n \sqrt{\text{Var}(y_i)} = \sum_{i=1}^n (1-f) \frac{S^2}{n}$

• to find minimum required sample size (n). $D = \frac{B^2}{Z_{\alpha/2}^2} = \frac{6^2}{n} \frac{N-n}{N-1} \Rightarrow n \approx \frac{N 6^2}{(N-1) D + 6^2}$ (B 会缩, $G = \frac{\text{range}}{4}$) $6^2 = E(X^2) - E(X)^2$

Estimation of a population total: Population $\{u_1, \dots, u_N\}$. total $T = u_1 + \dots + u_N = Nu$.

Sample: $\{y_1, \dots, y_n\}$. Estimator of T : $\hat{T} = N\bar{y}$. $E(\hat{T}) = E(N\bar{y}) = NE(y) = Nu = T$. $\text{Var}(\hat{T}) = \text{Var}(N\bar{y}) = N^2 \frac{6^2}{n} \frac{N-n}{N-1}$

$\text{Var}(\hat{T}) = N^2 (1-f) \frac{S^2}{n}$ $100(1-\alpha)\%$ C.I. for T is $\hat{T} \pm Z_{\alpha/2} \sqrt{\text{Var}(\hat{T})} = N\bar{y} \pm Z_{\alpha/2} \frac{6}{\sqrt{n}} \sqrt{1-f}$ (if G known. $B = N \frac{6}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}$)

• minimum required sample size (n). $D = \frac{B^2}{Z_{\alpha/2}^2} = \frac{6^2}{n} \left(\frac{N-n}{N-1} \right)$ also $n \approx \frac{N 6^2}{(N-1) D + 6^2}$ (n 要取整) Compared to mean.

Estimation of a population proportion:

for population $\{u_1, \dots, u_N\}$. $u_i = 1$ if i th element has a specified characteristic 0 if not. ($u_i^2 = u_i$)

let p be the population proportion: $p = \frac{\sum u_i}{N}$, population mean $\bar{u} = \frac{1}{N} \sum u_i = p$. $\frac{6^2}{n} = \frac{1}{N} \sum (u_i - \bar{u})^2 = p(1-p)$

Sample $\{y_1, \dots, y_n\}$. proportion: $\hat{p} = \frac{\sum y_i}{n} = \bar{y}$ Estimator of $p = \hat{p}$, $E(\hat{p}) = p$. $S^2 = \frac{1}{n-1} \sum (y_i - \bar{y})^2 = \frac{n}{n-1} \hat{p}(1-\hat{p})$

$E(S^2) = \frac{n}{n-1} p(1-p)$. $\text{Var}(\hat{p}) = \frac{p(1-p)}{n-1}$, $\text{Var}(\hat{p}) = (1-f) \frac{\hat{p}(1-\hat{p})}{n-1}$ $100(1-\alpha)\%$ C.I. for p is $\hat{p} \pm Z_{\alpha/2} \sqrt{\text{Var}(\hat{p})} = \hat{p} \pm Z_{\alpha/2} \frac{\sqrt{\hat{p}(1-\hat{p})}}{\sqrt{n-1}}$

Minimum required sample size n . $D = \frac{B^2}{Z_{\alpha/2}^2}$, $n \approx \frac{N p(1-p)}{(N-1) D + p(1-p)}$ (n 要取整 or max)

Comparing estimates: $E(\bar{x}-\bar{y}) = u_x - u_y$, $\text{Var}(\bar{x}-\bar{y}) = \text{Var}(\bar{x}) + \text{Var}(\bar{y}) - 2\text{Cov}(\bar{x}, \bar{y})$ (if samples are indep. $\text{Cov}(\bar{x}, \bar{y})=0$)

Ex. difference in mean in 95% CI is $(\bar{y}_1 - \bar{y}_2) \pm Z_{\alpha/2} \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}$ ($S^2 = \frac{N}{N-1} 6^2 = \frac{N}{N-1} \cdot \frac{1}{N} (u_1 - \bar{u})^2 = \frac{1}{N-1} (u_1 - \bar{u})^2$ contains 0.)

• $100(1-\alpha)\%$ C.I. for $p_1 - p_2$ is $(\hat{p}_1 - \hat{p}_2) \pm Z_{\alpha/2} \sqrt{\text{Var}(\hat{p}_1 - \hat{p}_2)} = (\hat{p}_1 - \hat{p}_2) \pm Z_{\alpha/2} \left(\frac{\hat{p}_1}{n_1} + \frac{\hat{p}_2}{n_2} + \frac{2\hat{p}_1\hat{p}_2}{n_1+n_2} \right)$ (if ignore, then underestimate)

Ex. multinomial sampling. X = num of responses in the first category (P_1). then $E(X) = np_1$, $E(Y) = np_2$, $E(Z) = np_3$

$\text{Var}(X) = np_1 q_1$, $\text{Var}(Y) = np_2 q_2$, $\text{Cov}(X, Y) = E(XY) - n^2 p_1 p_2 = -np_1 p_2$.

$E(\hat{p}_1) = p_1$, $\text{Var}(\hat{p}_1) = \frac{p_1 q_1}{n}$. $\text{Cov}(\hat{p}_1, \hat{p}_2) = -\frac{p_1 p_2}{n}$ $\text{Var}(\hat{p}_1 - \hat{p}_2) = \frac{p_1 q_1}{n} + \frac{p_2 q_2}{n} + 2 \frac{p_1 p_2}{n}$ Population stratum like male/female

NOTE 6. Stratified Random Sampling

separating the population element into nonoverlapping groups

population total = number of sampling units \times population mean ($T_i = N_i u_i$), $T = \sum_i N_i u_i$. $n = \frac{1}{N} \sum_i N_i u_i$ \uparrow 有 CI

Estimation of population mean: $\hat{u} = \bar{y}_{st} = \frac{1}{N} \sum_i N_i \bar{y}_i$. $E(\bar{y}_{st}) = E\left[\frac{1}{N} \sum_i N_i \bar{y}_i\right] = \frac{1}{N} T = u$. $\text{Var}(\bar{y}_{st}) = \frac{1}{N} \sum_i N_i^2 \frac{S_i^2}{n_i} (1 - \frac{n_i}{N})$

Estimation of total: $\text{Var}(N \bar{y}_{st}) = \text{Var}(\hat{T}_{st}) = \sum_i N_i^2 \frac{S_i^2}{n_i} (1 - \frac{n_i}{N})$ (for $S_i^2 = \frac{N_i}{N_i-1} 6^2$) $\bar{y}_{st} \pm Z_{\alpha/2} \sqrt{\text{Var}(\bar{y}_{st})}$ ($S_i^2 = \frac{N_i}{N_i-1} 6^2$)

Sample size determination (means) $D = \frac{B^2}{Z_{\alpha/2}^2}$ and $n = (\sum_i N_i^2 \frac{S_i^2}{n_i}) / (N^2 D + \sum_i N_i \frac{S_i^2}{n_i})$ ($n_i = n a_i$) 通过 a_i 例分母

Sample size determination (total) $D = B^2 / (N^2 Z_{\alpha/2}^2)$ and $n = (\sum_i N_i^2 \frac{S_i^2}{n_i}) / (N^2 D + \sum_i N_i \frac{S_i^2}{n_i})$ (题目会告诉你 estimate total?)

Optimal allocation of the sample : to gain the most information for the least cost. $C = c_0 + \sum_i c_i n_i$

cost = overhead cost (operating expenses, fixed) + variable cost C depending on number of observed units)

I. Min $\text{Var}(\bar{y}_{st})$, for a given total cost $C = c_0 + \sum_i c_i n_i$ (min var for fix cost)

$$n_i = \left(\frac{c - c_0}{\sum_{i=1}^L N_i \sigma_i \sqrt{c_i}} \right) \frac{N_i \sigma_i}{\sqrt{c_i}}$$

$$n = \left(\frac{c - c_0}{\sum_{i=1}^L N_i \sigma_i \sqrt{c_i}} \right) \sum_{i=1}^L \frac{N_i \sigma_i}{\sqrt{c_i}}$$

$$\frac{n_i}{n} = a_i = (N_i \sigma_i / \sqrt{c_i}) / (\sum_{i=1}^L N_i \sigma_i / \sqrt{c_i}) \quad B = \sum_{i=1}^L \sqrt{c_i} \quad (\text{给 B 就是给 V})$$

$$\text{II. Min total cost } C, \text{ given } \text{Var}(\bar{y}_{st}) = V \quad \text{就是 } D.$$

$$n_i = \frac{N_i \sigma_i}{N^2 V + \sum_{i=1}^L N_i \sigma_i^2} \quad n_i = \frac{(N_i \sigma_i) \sum N_i \sigma_i}{N^2 D + \sum N_i \sigma_i^2}$$

Proportional allocation (ignore both cost and variability), Allocation scheme $\frac{n_i}{n} = \frac{N_i}{N} = a_i$

$$n = \left(\frac{1}{N} \sum_i N_i \sigma_i^2 \right) / (ND + \frac{1}{N} \sum_i N_i \sigma_i^2 / N) = \left(\frac{1}{N} \sum_i N_i \sigma_i^2 / a_i \right) / (N^2 D + \frac{1}{N} \sum_i N_i \sigma_i^2 / a_i), \text{ where } D = \frac{B^2}{Z_{\alpha/2}^2}, \text{ then } V(\bar{y}_{st}) = \sum_i N_i \frac{6^2}{n_i} \frac{N_i - n_i}{N_i - 1}$$

Estimation of population proportion: $S_i^2 = \frac{n_i}{n_i-1} \hat{p}_i q_i$, unbiased estimator of p is $\hat{p}_{st} = \frac{1}{N} \sum_i N_i \hat{p}_i$ where $\hat{p}_i = \frac{n_i}{N_i}$, $100(1-\alpha)\%$ C.I. of p is $\hat{p}_{st} \pm Z_{\alpha/2} \sqrt{V(\hat{p}_{st})}$

→ sample size: $n = \left(\frac{1}{N} \sum_i N_i p_i q_i / a_i \right) / (N^2 D + \frac{1}{N} \sum_i N_i p_i q_i)$ where $D = \frac{B^2}{Z_{\alpha/2}^2}$ (estimate proportion with error B and level 1- α)

Optimal allocation: $n_i = n \left(\frac{N_i \sqrt{p_i q_i} / \sqrt{c_i}}{\sum_{i=1}^L N_i \sqrt{p_i q_i} / \sqrt{c_i}} \right)$ Neyman allocation: $n_i = n \left(\frac{N_i \sqrt{p_i q_i}}{\sum_{i=1}^L N_i \sqrt{p_i q_i}} \right)$ Proportional allocation: $n_i = n \left(\frac{N_i}{N} \right)$ $\sqrt{p_i q_i} = b_i$

• if we want to find CI for $y_1 - y_2$: $y_1 - y_2 \pm Z_{\alpha/2} \frac{\sqrt{N_1 - n_1} \frac{S_1^2}{n_1} + \sqrt{N_2 - n_2} \frac{S_2^2}{n_2}}{\sqrt{N_1 - n_1} + \sqrt{N_2 - n_2}}$

① $\text{Var}(\bar{y}_{st}) \leq \text{Var}(\text{opt}(\bar{y}_{st}))$ ② $\text{Var}(\text{Neyman}(\bar{y}_{st})) \leq \text{Var}(\text{prop}(\bar{y}_{st}))$

③ if $N_i, i=1 \dots L$ are large, then $\text{Var}(\text{Neyman}(\bar{y}_{st})) \leq \text{Var}(\text{prop}(\bar{y}_{st})) \leq \text{Var}(\text{opt}(\bar{y}_{st}))$

④ $\frac{\text{Var}(\bar{y}_{st})}{\text{Var}(\bar{y}_{st})} = \frac{\sum_i N_i \sigma_i^2 + \sum_i N_i (u_i - \bar{u})^2}{\sum_i N_i \sigma_i^2} = 1 + \frac{\sum_i N_i (u_i - \bar{u})^2}{\sum_i N_i \sigma_i^2}$

(要先找 Allocation 才能回去算 n.)

NOTE 7. Ratio, Regression and Difference Estimation

Population N pairs $(U_1, V_1) \dots (U_N, V_N)$ when $U_x = \frac{1}{N} \sum_i U_i$. $T_x = \frac{1}{N} \sum_i V_i$. $S_{rx}^2 = \frac{1}{N} \sum_i (U_i - U_x)^2$, 把 U_i 换成 V_i 就是 $V_y T_y$ by

Sample n pairs $(X_1, Y_1) \dots (X_n, Y_n)$, X_i is called auxiliary variable/subsidiary variable. (to be used for forecasting)

Estimation of population ratio: $R = \frac{\sum_i V_i}{\sum_i U_i} = \frac{T_y}{T_x} = \frac{V_y}{U_x}$ $r = \frac{V_y}{U_x}$. $E(r) - R \approx \frac{N-n}{N-1} \left(\frac{1}{n n_x^2} \right) (R S_{rx}^2 - P_b S_{rx} S_{by})$ ($E(r) = R$ 如果不说)

$\text{Var}(r) \approx \frac{1}{n} \frac{S_{rx}^2}{U_x^2} \left(\frac{N-n}{N-1} \right)$ where $S_{rx}^2 = \frac{1}{N} \sum_i (U_i - U_x)^2 = S_x^2 + R^2 S_{rx}^2 - 2 R P_b S_{rx} S_{by}$

$\widehat{\text{Var}}(r) = \frac{1}{U_x^2} \frac{S_{rx}^2}{n} (1-f)$ if U_x is known, $\widehat{\text{Var}}(r) = \frac{1}{U_x^2} \frac{S_{rx}^2}{n} (1-f)$ if U_x is unknown

then $100(1-\alpha)\%$ C.I. of R is $R \pm Z_{\alpha/2} \sqrt{\widehat{\text{Var}}(r)}$, $\text{Var}(r) = \frac{1}{U_x^2} \frac{S_{rx}^2}{n} \left(\frac{N-n}{N-1} \right)$, where $S_{rx}^2 = S_y^2 + r^2 S_x^2 - 2 r \hat{P}_S S_x S_y = \frac{1}{n-1} \sum_i (Y_i - r X_i)^2$

用 ratio 的条件: linear and pass zero. the required sample size is $n = \frac{ND}{N S_{rx}^2}$ where $D = \frac{B^2 n_x^2}{S_{rx}^2}$

Estimation of population total: $T_y = R T_x$, $\hat{T}_y = r T_x$. $\text{Var}(\hat{T}_y) = (T_x)^2 \text{Var}(r)$, $\widehat{\text{Var}}(\hat{T}_y) = (T_x)^2 \widehat{\text{Var}}(r) = (T_x)^2 \frac{1}{U_x^2} \frac{S_{rx}^2}{n} (1-f)$

$100(1-\alpha)\%$ C.I. for T_y is $\hat{T}_y \pm Z_{\alpha/2} \sqrt{\widehat{\text{Var}}(\hat{T}_y)}$. Sample size: $n = \frac{NS_{rx}^2}{ND + S_{rx}^2}$ where $D = \frac{B^2}{S_{rx}^2}$

Estimation of population mean: $U_y = R U_x$, $\hat{U}_y = r U_x$, $\text{Var}(U_y) = (U_x)^2 \text{Var}(r) = \left(\frac{N-n}{N-1} \right) \frac{S_y^2}{n}$, $\widehat{\text{Var}}(\hat{U}_y) = (U_x)^2 \widehat{\text{Var}}(r) = \frac{S_y^2}{n} (1-f)$

$100(1-\alpha)\%$ C.I. for U_y is $\hat{U}_y \pm Z_{\alpha/2} \sqrt{\widehat{\text{Var}}(\hat{U}_y)}$, sample size $n = \frac{NS_{rx}^2}{ND + S_{rx}^2}$ where $D = \frac{B^2}{S_{rx}^2}$

Regression Estimation (x_i, y_i) be $y = \alpha + \beta x$. the least squares for β : $b = \frac{S_{xy}}{S_{xx}}$, α : $a = \bar{y} - b \bar{x}$, best fitted line $\hat{y} = a + b x = \bar{y} + b(x - \bar{x})$

Regression estimator of population mean U_y $\hat{U}_{yl} = \bar{y} + b(U_x - \bar{x})$, $\widehat{\text{Var}}(\hat{U}_{yl}) = \frac{(1-f)MSE}{n}$, $MSE = \frac{\sum (y_i - \hat{y})^2}{n-2}$, $\text{Var}(U_{yl}) = \left(\frac{N-n}{N-1} \right) \frac{1}{n} B^2$

Difference Estimation $U_y = U_x + (U_y - U_x)$ estimate difference the difference $D = U_y - U_x$ $S_E^2 = S_y^2 + B^2 S_x^2$

$\hat{U}_{yd} = U_x + (\bar{y} - \bar{x}) = U_x + \bar{d}$, $\bar{d} = \bar{y} - \bar{x}$, $d_i = y_i - x_i$, $\text{Var}(\hat{U}_{yd}) = \frac{(1-f)S_d^2}{n} = \left(\frac{N-n}{N-1} \right) \frac{S_d^2}{n}$

and $S_d^2 = \frac{1}{n-1} \sum (d_i - \bar{d})^2 = \frac{1}{n-1} \sum [(y_i - \bar{y}) - (x_i - \bar{x})]^2 = S_y^2 + S_x^2 - 2 \hat{P}_S S_x S_y$

这种是不可以的

NOTE 8. Systematic Sampling /-in-K systematic sampling (抽第1, 4, 7, 10...个), 只抽一个group. 确保 Var between group small

Estimation of population mean: Sample $\{Y_1, Y_2, \dots, Y_n\}$. $\bar{Y}_{sy} = \frac{1}{n} \sum_i Y_i$, variance $S_{sy}^2 = \frac{1}{n-1} \sum (Y_i - \bar{Y}_{sy})^2$, $\text{Var}(\bar{Y}_{sy}) = \frac{1}{K} \sum_i (U_i - \bar{U})^2$

$\widehat{\text{Var}}(\bar{Y}_{sy}) = (1-f) \frac{S_{sy}^2}{n}$, $100(1-\alpha)\%$ C.I. for U : $\bar{Y}_{sy} \pm Z_{\alpha/2} \frac{(1-f) S_{sy}^2}{n}$

$SST = \sum_{i=1}^k \sum_{j=1}^n (U_{ij} - \bar{U})^2$, $SSW = \sum_{i=1}^k \sum_{j=1}^n (Y_{ij} - \bar{U}_i)^2$, $SSB = \sum_{i=1}^k n (U_i - \bar{U})^2$, $SST = SSW + SSB$, $MST = \frac{SST}{N-1}$, $MSW = \frac{SSW}{k(n-1)} = \frac{SSW}{N-k}$, $MSB = \frac{SSB}{k-1}$

then $\text{Var}(\bar{Y}_{sy}) = \frac{S_{sy}^2}{n} (1 + (n-1)P_w)$ where $P_w = \frac{E(Y_{ij} - \bar{U}_i)(Y_{il} - \bar{U}_l)}{E(Y_{ij} - \bar{U}_i)^2} = \frac{\sum_{i=1}^k \sum_{j=1}^n (U_{ij} - \bar{U}_i)(U_{il} - \bar{U}_l)}{N K (n-1) \cdot S_{sy}^2}$

(Compare) if $\text{Var}(\bar{Y}_{sy}) < \text{Var}(\bar{Y}_{SPS}) \Leftrightarrow B^2 - \frac{N-k}{N} MSW < \frac{S_{sy}^2}{n} \Leftrightarrow MSW > B^2 \frac{N}{N-1}$

Estimation of population total $\hat{T} = N \bar{Y}_{sy} = \frac{N}{n} \sum_i Y_i$, $\widehat{\text{Var}}(\hat{T}_{sy}) = N^2 (1-f) \frac{S_{sy}^2}{n}$, $100(1-\alpha)\%$ C.I. for T is $\hat{T}_{sy} \pm Z_{\alpha/2} \sqrt{N^2 (1-f) \frac{S_{sy}^2}{n}}$

Estimation of population proportion $\hat{P}_{sy} = \bar{Y}_{sy} = \frac{1}{n} \sum_i Y_i$, $\widehat{\text{Var}}(\hat{P}_{sy}) = (1-f) \frac{P_{sy} Q_{sy}}{n-1}$, $100(1-\alpha)\%$ C.I. for p : $\hat{P}_{sy} \pm Z_{\alpha/2} \sqrt{(1-f) \frac{\hat{P}_{sy} \hat{Q}_{sy}}{n-1}}$

NOTE 9&10 Cluster Sampling externally homogeneous but internally heterogeneous 每个cluster都大致能代表总体性质

(Clusters of equal sizes) Primary sampling units (PSU): Sampling units chosen in the first stage of selection

Secondary sampling units (SSU) Sampling units within the PSUs that are chosen in the second stage of selection.

One-stage cluster sampling: 这一个cluster里面所有都要 Two-stage cluster sampling: 在选定的cluster里 SRS 之类

Let $M_1 = \dots = M_N = M$, total population size = NM , Total sample size = nM

Estimation of population mean: Select n clusters in N by SRS.

$$\hat{Y} = \bar{Y}_{cl} = \frac{1}{n} \sum_i^N \bar{y}_i \text{ (average of cluster mean)}, \hat{u} = \bar{y}_c \text{ unbiased of } u$$

$$\text{Var}(\bar{Y}_{cl}) = \frac{N-n}{NM} S_b^2 \text{ where } S_b^2 = \frac{1}{N-1} \sum_i^N (u_i - \bar{u})^2 \text{ mean sum of square between cluster means}$$

$$\text{Var}(\hat{Y}_{cl}) = \frac{N-n}{NM} S_b^2 \text{ where } S_b^2 = \frac{1}{n-1} \sum_i^N (\bar{y}_i - \bar{Y}_{cl})^2$$

(compare to SRS) define $S^2 = \frac{1}{NM-1} \sum_{i=1}^{N-1} \sum_{j=1}^m (u_{ij} - \bar{u})^2$ and \bar{y} sample mean of SRS

$$\text{relative efficiency RE} = \frac{\text{Var}(\bar{y})}{\text{Var}(\bar{y}_{cl})} = \frac{1}{NM-1} \left[\frac{N(N-1)}{M} \frac{S_w^2}{S_b^2} + (N-1) \right] = \frac{S_w^2}{M S_b^2} \frac{1}{N} \sum_i^N S_i^2 \text{ within}$$

$$RE = \frac{M-1}{M} \frac{S_w^2}{MS_b^2} + \frac{1}{M}$$

cluster sampling is good when S_w^2 is large and S_b^2 is small

Estimation of population proportion p : $\hat{P}_{cl} = \frac{1}{n} \sum_i^N p_i$ unbiased of p

$$\text{Var}(\hat{P}_{cl}) = \frac{N-n}{NM} \frac{1}{N-1} (NPq - \frac{1}{n} P_{cl} q_{cl}) \quad \text{Var}(\hat{P}_{cl}) = \frac{N-n}{NM} \frac{1}{N-1} \sum_i^N (p_i - \hat{P}_{cl})^2$$

(Clusters of unequal sizes) 好处是: m_0 不知道时, 抽几个 cluster 去做就好.

Estimation of population mean:

$$(\text{unbiased}) \hat{Y} = \bar{Y}_n = \frac{1}{nm} \sum_i^N M_i \bar{y}_i \text{ 但不知道就无法估计.}$$

$$E(\bar{Y}_n) = u, \text{ Var}(\bar{Y}_n) = \frac{N-n}{NM} S_b^2 \text{ where } S_b^2 = \frac{1}{n-1} \sum_i^N \left(\frac{M_i}{m} \bar{y}_i - \bar{Y}_n \right)^2$$

$$(\text{biased}) \hat{Y} = \bar{Y} = \frac{\sum_i^N \bar{y}_i}{\sum_i^N m_i}, E(\bar{Y}) = u, \text{ Var}(\bar{Y}) = (1-\frac{n}{N}) \frac{S_r^2}{NM} \text{ where } S_r^2 = \frac{\sum_i^N (y_i - \bar{y}_m)^2}{n-1} \text{ ratio observation}$$

Estimation of population total (when M_0 Known) $\hat{T} = M_0 \bar{Y}$. $\text{Var}(\hat{T}) = N^2 (1-\frac{n}{N}) \frac{S_r^2}{NM}$

$$(\text{M}_0 \text{ unknown}) \hat{T} = N \bar{Y}_T = N \left(\frac{1}{n} \sum_i^N y_i \right) = N \bar{M} \bar{Y}_n, \text{ Var}(\hat{T}) = \text{Var}(N \bar{M} \bar{Y}_n) = N^2 \bar{M}^2 \frac{N-n}{NM} \frac{1}{n-1} \sum_i^N \left(\frac{M_i}{m} \bar{y}_i - \bar{Y}_n \right)^2 = N^2 \frac{N-n}{NM} S_T^2 \text{ where } S_T^2 = \frac{1}{n-1} \sum_i^N (y_i - \bar{Y}_T)^2$$

$$\bar{Y}_T = \bar{M} \bar{Y}_n = \frac{1}{n} \sum_i^N \bar{y}_i, \text{ Var}(\bar{Y}_T) = \frac{N-n}{NM} \frac{1}{n-1} \sum_i^N \left(\frac{M_i}{m} \bar{y}_i - \bar{Y}_n \right)^2, B = 1.96 \sqrt{\text{Var}(\hat{T})} = 1.96 \sqrt{\frac{N(N-n)}{NM}} S_T$$

$$\text{Estimation of population proportion } p: \hat{P} = \frac{\sum_i^N a_i}{\sum_i^N m_i} \text{ (unbiased of } p), \text{ Var}(\hat{P}) = (1-\frac{n}{N}) \frac{S_p^2}{NM} \text{ where } S_p^2 = \frac{\sum_i^N (a_i - \hat{P} M_0)^2}{n-1}$$

Notation

Population quantities (refer to population data on p.8)

- N = number of clusters (psus) in the population
- M_i = number of elements (ssus) in cluster i
- $M_0 = \sum_{i=1}^N M_i$ = total number of elements (ssus) in the population
- $\bar{M} = \frac{M_0}{N}$ = Average cluster size for the population
- $t_i = \sum_{j=1}^{M_i} u_{ij}$ = total in cluster i
- $\tau = \sum_{i=1}^N t_i = \sum_{i=1}^N \sum_{j=1}^{M_i} u_{ij}$ = population total
- $\mu = \frac{\tau}{M_0}$ = population mean
- $\mu_i = \frac{t_i}{M_i}$ = population mean in cluster i
- $\sigma^2 = \frac{1}{M_0-1} \sum_{i=1}^N \sum_{j=1}^{M_i} (u_{ij} - \mu)^2$ = population variance
- $\sigma_i^2 = \frac{1}{M_i-1} \sum_{j=1}^{M_i} (u_{ij} - \mu_i)^2$ = population variance in cluster i

Sample quantities

- n = number of clusters (psus) in the sample
- m_i = number of elements (ssus) in cluster i in the sample
- $\sum_{i=1}^n m_i$ = total number of elements (ssus) in the sample
- $\bar{m} = \frac{1}{n} \sum_{i=1}^n m_i$ = Average cluster size for the sample
- y_i = total of all observations in cluster i
- $\bar{y} = \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n m_i}$ = sample mean