

제 6장: 집락표집

6.1 집락표집의 개요

- 집락표집은 계통표집과 마찬가지로 조사의 편리성을 고려한 조사 방법
- 예를 들어, 서울시에서 각 가구별로 보유한 컴퓨터의 총 수를 추정하고자 하는 경우 서울시에 존재하는 모든 가구의 목록이 필요하나 이런 목록을 얻기가 어렵고 모든 가구의 목록을 얻는다 하더라도 단순임의추출을 시행하는 경우 조사 개체들이 서울시 전체에 산포되게 되어 조사비용이 커지게 된다. 이 경우 서울시에 존재하는 모든 동의 목록을 마련하고 이들 중 50개의 동을 임의로 선택하게 되고 선택된 동에서 5개의 통씩 임의로 선택하고 또 선택된 통 내의 모든 가구들을 조사하는 방식을 선택하게 된다.
- 이런 표집방법을 이단집락표집(two-stage cluster sampling)이라 하고
일차표집단위 (primary sampling unit, psu)는 동 들이
이차표집단위 (secondary sampling unit, ssu)는 통 들이 된다.
- 집락표집을 이용하여 모수를 추정하는 경우 동일 크기의 단순임의표집에 비하여 추정의 정밀도가 낮아지는 경향이 있을 수 있다. 이는 집락표집을 시행할 때 집락내의 구성원들의 성향이 비슷함으로서 나타나는 현상이다. 집락표본에서 최대 정밀도를 얻기 위하여는 집락내의 원소는 서로 이질적이고 집락간은 서로 동질적이어야 한다.
- 간단한 예로 크기가 $N(= Km)$ 인 모집단에서 $n = 2m$ 개의 표본을 얻는 경우 크기가 m 인 집락을 2개 뽑았다고 하면

$$\hat{\mu}_{cl} = \frac{1}{2m} \sum_{i=1}^2 \sum_{j=1}^m Y_{ij} = \frac{1}{2}(\bar{Y}_1 + \bar{Y}_2)$$

이다. 이의 분산은 집락간 독립을 가정하면

$$\text{var}(\hat{\mu}_{cl}) = \frac{1}{4}(\text{var}(\bar{Y}_1) + \text{var}(\bar{Y}_2))$$

이고

$$\text{var}(\bar{Y}_1) = \text{var}(\bar{Y}_2) = \frac{1}{m^2} \left\{ \sum_{j=1}^m \text{var}(Y_{1j}) + \sum_{j1 \neq j2} \text{cov}(Y_{1,j1}, Y_{2,j2}) \right\}.$$

- 만일 집락내가 동질적이라면 $\text{var}(\bar{Y}_1) \geq \frac{\sigma^2}{m}$ 이 되어 $\text{var}(\hat{\mu}_{cl})$ 이 $\text{var}(\hat{\mu}_{sys})$ 보다 커지게 된다.
- 집락표집에서는 표집단위 (sampling unit)이 집락이므로 집락표집의 표본크기는 선택된 집락의 개수를 의미하고 조사한 개체들의 수가 아니다.

6.2 집락표집에 대한 표기

y_{ij} : i 번째 psu의 j 번째 ssu에 해당하는 관측값

모수	표본통계량
N =모집단에서 psu들의 총수, 집락 수	n =표본으로 선택된 psu들의 총수, 표집된 집락수
M_i = i 번째 집락의 크기: i 번째 psu에 속한 ssu들의 총수	m_i : i 번째 psu에서 추출하는 표본의 크기
$K = \sum_{i=1}^N M_i$: 모집단에서 ssu들의 총수	$k = \sum_{i=1}^n m_i$: 총 표본크기
$\bar{M} = K/N = \sum_{i=1}^N M_i/N$: psu당 평균 크기	$\bar{m} = \sum_{i=1}^n m_i/n$: 추출된 psu당 평균 표본 크기
$\tau_i = \sum_{j=1}^{M_i} y_{ij}$: i 번째 psu의 모합	$t_i = \sum_{j=1}^{m_i} y_{ij}$: i 번째 psu의 표본합 $\hat{\tau}_i = \frac{M_i}{m_i} t_i$: i 번째 psu의 모합에 대한 추정량
$\tau = \sum_{i=1}^N \tau_i = \sum_{i=1}^N \sum_{j=1}^{M_i} y_{ij}$: 모합	$\hat{\tau} = \frac{N}{n} \sum_{i=1}^n \hat{\tau}_i$: 모합에 대한 추정량
$\bar{\tau} = \tau/N$: psu당 평균합	$\bar{t} = \hat{\tau}/N$: psu당 평균합에 대한 추정량
$\sigma_b^2 = \frac{1}{N-1} \sum_{i=1}^N (\tau_i - \bar{\tau})^2$: psu 합들에 대한 모분산	$s_b^2 = \frac{1}{n-1} \sum_{i=1}^n (\hat{\tau}_i - \bar{t})^2$
$\mu = \tau/K$: (ssu당)모평균	$\hat{\mu} = \hat{\tau}/K$
$\mu_i = \tau_i/M_i$: i 번째 psu의 모평균	$\hat{\mu}_i = \hat{\tau}_i/M_i$
$\sigma^2 = \frac{1}{K-1} \sum_{i=1}^N \sum_{j=1}^{M_i} (y_{ij} - \mu)^2$: (ssu당)모분산	$s^2 = \frac{1}{k-1} \sum_{i=1}^n \sum_{j=1}^{m_i} (y_{ij} - \hat{\mu})^2$
$\sigma_i^2 = \frac{1}{M_i-1} \sum_{j=1}^{M_i} (y_{ij} - \mu_i)^2$: i 번째 psu의 모분산	$s_i^2 = \frac{1}{m_i-1} \sum_{j=1}^{m_i} (y_{ij} - \hat{\mu}_i)^2$

6.3 일단집락표집

일단집락표집은 선택된 psu에 속한 ssu들을 전부 조사한다. 일단집락표집의 예로는 고등학생들에 대한 설문조사의 경우 psu는 학급이고 ssu는 학생, 폐렴으로 입원한 환자에 대한 조사의 경우 병원을 psu로 폐렴환자를 ssu로 조사한다.

6.3.1 동일한 집락의 크기

- 편의상 $M_1 = M_2 = \dots = M_N = M$ 을 가정하고 이 때 전체 개체의 수 $K = N \cdot M$ 이다.
- 관측값은 $\{\tau_1, \tau_2, \dots, \tau_n\}$ 이 되고 여기서 τ_i 는 i 번째 선택된 집락에서 모든 원소들에 대한 합이다.

- 전체 모함에 대한 추정량은

$$\hat{\tau} = \frac{N}{n} \sum_{i=1}^n \tau_i$$

가 된다.

- 추정량의 분산: $\{\tau_1, \tau_2, \dots, \tau_N\}$ 의 모집단에서 n 개를 SRS 했다고 생각하여 SRS의 분산공식을 이용하여 계산할 수 있다.

$$\text{var}(\hat{\tau}) = N^2 \frac{N-n}{N} \frac{\sigma_b^2}{n}$$

이고 여기서

$$\sigma_b^2 = \frac{1}{N-1} \sum_{i=1}^N (\tau_i - \mu)^2.$$

이다.

- 분산의 추정량:

$$\widehat{\text{var}}(\hat{\tau}) = N^2 \frac{N-n}{N} \frac{s_b^2}{n}$$

이고

$$s_b^2 = \frac{1}{n-1} \sum_{i=1}^n (\tau_i - \bar{t})^2, \quad \bar{t} = \frac{\hat{\tau}}{N}$$

이다.

- 모평균 μ 의 추정

6.3.2 집락표집의 효율

- 집락의 크기가 같음을 가정하고 표본의 크기가 n 인 집락표집과 표본의 크기가 nM 인 단순임의추출의 효율을 비교한다. $K = NM$ 이라 하자.
- 각각의 분산이

$$\text{var}(\hat{\tau}_{\text{cl}}) = N^2 \frac{N-n}{N} \frac{1}{n} \sigma_b^2$$

과

$$\text{var}(\hat{\tau}_{\text{srs}}) = K^2 \frac{K-nM}{K} \frac{\sigma^2}{nM} = N^2 \left(\frac{N-n}{N} \right) \frac{M}{n} \sigma^2$$

이다. 따라서 σ_b^2 과 $M\sigma^2$ 을 비교한다. (σ_b^2 과 σ^2 의 정의에서 분모를 $N-1$ 과 $NK-1$ 을 사용하였다.)

• 계산:

No.

Date.

$$\sum_{i=1}^N \sum_{j=1}^M (y_{ij} - \mu)^2 = \sum_{i=1}^N \sum_{j=1}^M (\mu_i - \mu)^2 + \sum_{i=1}^N \sum_{j=1}^M (y_{ij} - \mu_i)^2$$

SST

$NM - 1$

SSB

$N - 1$

SSW

$N(M - 1)$

$$\sigma_b^2 = \frac{1}{N-1} \sum_{i=1}^N (\bar{y}_i - \bar{y})^2 = \frac{1}{N-1} \sum_{i=1}^N (M\mu_i - M\mu)^2$$

$$= \frac{M}{N-1} SSB = M \cdot MSB$$

군내 상관관계수

$$\rho_w = \frac{E((y_{ij} - \mu)(y_{ij'} - \mu))}{E(y_{ij} - \mu)^2} = 1 - \frac{M}{M-1} \frac{SSW}{SST}$$

임의의 교재의 p158, 160을 통하여 확인한다.

$$MSB = \frac{1}{N-1} \{SST - SSW\}$$

$$= \frac{1}{N-1} \left\{ 1 - \frac{M-1}{M} (1 - \rho_w) \right\}$$

$$= SST \frac{1}{M} \frac{1}{N-1} \{M - M + 1 + \rho_w(M-1)\}$$

$$= \sigma^2 \frac{NM-1}{M(N-1)} \{1 + (M-1)\rho_w\}$$

$$\approx \sigma^2 \frac{N}{N-1} \{1 + (M-1)\rho_w\}$$

$$1 - \rho_w = \frac{M}{M-1} \frac{SSW}{SST}$$

$$SSW = SST \cdot \frac{M-1}{M} (1 - \rho_w)$$

$$\frac{1}{NM} SST = \sigma^2$$

- 따라서

$$\begin{aligned}\frac{\text{var}(\hat{\tau}_{\text{cl}})}{\text{var}(\hat{\tau}_{\text{srs}})} &= \frac{\sigma_b^2}{M\sigma^2} = \frac{M\text{MSB}}{M\sigma^2} = \frac{\text{MSB}}{\sigma^2} \\ &= \frac{N}{N-1} \left\{ 1 + (M-1)\rho_w \right\}.\end{aligned}$$

- 급내상관계수의 범위는

$$-\frac{1}{M-1} \leq \rho_w \leq 1$$

이고 ρ_w 가 음수이면 - 즉, 집락내의 원소의 성질이 이질적이면 - 집락추출이 단순임의추출보다 효율적이다.

6.3.3 크기가 다른 집합들

일반적으로 집합의 크기는 다르고 i 번째 집합의 크기를 M_i 라 하자

(1) 비편향 추정 (unbiased estimator)

• 추정량

$$\hat{T}_w = \frac{N}{n} \sum_{i \in A} Z_i T_i \quad Z_i = \begin{cases} 1 & \text{만일 } i \text{ 번째 집합 추출} \\ 0 & \text{o.w} \end{cases}$$

• 평균

$$E(\hat{T}_w) = \frac{N}{n} \sum_{i=1}^N \frac{n}{N} T_i = T$$

• 추정량의 분산

$$\text{var}(\hat{T}_w) = N^2 \frac{N-n}{N} \frac{\sigma_b^2}{n}$$

$$\sigma_b^2 = \frac{1}{N-1} \sum_{i=1}^N (T_i - \bar{T})^2, \quad \bar{T} = \frac{T}{N} = \text{psu당 평균}$$

• 분산의 추정량

$$\hat{\text{var}}(\hat{T}_w) = N^2 \frac{s_b^2}{n} \frac{N-n}{N}$$

$$s_b^2 = \frac{1}{n-1} \sum_{i=1}^n (T_i - \bar{T})^2, \quad \bar{T} = \frac{1}{N} \hat{T}_w$$

모집단 원소의 총수를 $K = \sum_{i=1}^N M_i$ 라면 모평균의 추정값

$$\hat{\mu}_u = \frac{1}{K} \hat{T}_u$$

$$\text{var}(\hat{\mu}_u) = \frac{1}{K^2} N^2 \frac{s_b^2}{n} \left(\frac{N-n}{N} \right)$$

$$\widehat{\text{var}}(\hat{\mu}_u) = \frac{1}{K^2} N^2 \frac{s_b^2}{n} \left(\frac{N-n}{N} \right)$$

이 된다.

• 일반적으로 N 은 알려져 있으나 K 은 알려져 있지 않음을 기억한다

(2) 비추정

모평균의 다른 추정량으로 비추정량을 생각할 수 있다.

비추정량은 ① $K = \sum_{i=1}^N M_i$ 를 알 필요가 없고

② 편향 (bias) 가 있으나 종종 분산이 작아져 (1)절의
불편추정량보다 효율적이다

$$\hat{\mu}_r = \frac{\sum_{i=1}^n T_i}{\sum_{i=1}^n M_i} \quad ; \quad \text{비추정량}$$

이 추정량은 4장의 비 β 에 대한 추정량 b 와 동일하다.

교재의 P120 을 참고하면

$$\text{var}(b) \approx \frac{1}{M_y^2} \frac{s_y^2}{n} \frac{N-n}{N} \quad \text{로 부터}$$

$$\widehat{\text{var}}(\hat{\mu}_r) \approx \frac{1}{\bar{M}^2} \frac{s_e^2}{n} \frac{N-n}{N} \quad \text{이고}$$

$$\bar{M} = \frac{1}{N} \sum_{i=1}^N M_i = \frac{K}{N}$$

$$\begin{aligned} s_e^2 &= \frac{1}{n-1} \sum_{i=1}^n (\tau_i - \hat{\mu}_r M_i)^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n M_i^2 (\bar{Y}_i - \hat{\mu}_r)^2 \quad \text{이다.} \end{aligned}$$

만일 \bar{M}^2 보다는 $\bar{M}_3 = \frac{1}{n} \sum_{i=1}^n M_i^2$ 사용한다.

위의 비측정량을 이용한 모형의 측정량

$$\hat{\tau}_r = K \hat{\mu}_r \quad \text{or} \quad \hat{\tau}_r = \hat{K} \hat{\mu}_r, \quad \hat{K} = N \bar{M}_3$$

$$\widehat{\text{var}}(\hat{\tau}_r) = N^2 \frac{s_e^2}{n} \frac{N-n}{N}$$

정리

- 비측정량을 편향(bias)가 존재하나 unbiased estimator 보다 분산이 작은 경향이 있다.
- 모형에 대한 측정시 비측정량은 모형변인소의 총수 $K = \sum_{i=1}^N M_i$ 를 알아야 한다.
- 모평균에 대한 측정시 비편향측정량은 모형변인소의 총수 $K = \sum_{i=1}^N M_i^2$ 를 알아야 한다.

만일 $M_i = M$ 이라면

$$\hat{\mu}_r = \frac{1}{nM} \sum_{i=1}^n \tau_i = \frac{1}{K} \frac{N}{n} \sum_{i=1}^n \tau_i = \frac{1}{K} \hat{T}_w$$

$$\hat{\tau}_r = K \cdot \hat{\mu}_r = \frac{N}{n} \sum_{i=1}^n \tau_i = \hat{\tau}_w$$

이다.

(3) 모비율의 추정

2 번째 계층에서 관성특성을 보유한 단위들의 총수를 a_i 라 표기하자.

① $M_1 = M_2 = \dots = M_N = M$ 인 경우

$$p_i = \frac{1}{M} a_i \quad i=1, 2, \dots, N$$

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n p_i$$

$\{p_1, p_2, \dots, p_n\}$ 은 N 개의 단위 $\{P_1, P_2, \dots, P_N\}$ 의

비율은 크기 n 이 비모양인 집단의 비율이다.

$$\text{Var}(\hat{p}) = \frac{N-n}{N} \frac{\sigma_p^2}{n} \quad \sigma_p^2 = \frac{1}{N-1} \sum_{i=1}^N (p_i - \bar{p})^2$$

$$\begin{aligned} \bar{p} &= \frac{1}{K} \sum_{i=1}^N \frac{a_i}{M} \\ &= \frac{1}{N \cdot M} \sum_{i=1}^N a_i \end{aligned}$$

$$\widehat{\text{var}}(\hat{p}) = \frac{N-n}{N} \frac{s_p^2}{n}, \quad s_p^2 = \frac{1}{n-1} \sum_{i=1}^n (p_i - \bar{p})^2$$

$$\bar{p} = \frac{1}{nM} \sum_{j=1}^n a_j$$

② M_j 등이 서로 다른 경우

이 경우 비동질성을

$$\hat{p}_r = \frac{\sum_{j=1}^n a_j}{\sum_{j=1}^n M_j}$$

$$\widehat{\text{var}}(\hat{p}_r) = \frac{N-n}{N} \frac{1}{\bar{M}^2} \frac{s_e^2}{n}$$

$$s_e^2 = \frac{1}{n-1} \sum_{j=1}^n (a_j - \hat{p}_r M_j)^2$$

$$\bar{M} = \frac{K}{N} = \frac{1}{N} \sum_{i=1}^N M_i$$

(4) 표본 크기의 결정

① 모평균 추정에 표본 크기

 $\hat{\mu}_r$ 사용한다

$$\text{var}(\hat{\mu}_r) = \frac{1}{M^2} \frac{\sigma_e^2}{n} \frac{N-n}{N}$$

$$Z_{\frac{\alpha}{2}} \sqrt{\text{var}(\hat{\mu}_r)} = \text{오차한계}(B) \quad (*)$$

 $(*)$ $\frac{Z}{2}$ $\frac{P}{2}$ 면

$$Z_{\frac{\alpha}{2}}^2 \frac{1}{M^2} \frac{\sigma_e^2}{n} \frac{N-n}{N} = B^2$$

$$n = \frac{N \sigma_e^2}{ND + \sigma_e^2}, \quad D = \left(\frac{B \cdot M}{Z_{\frac{\alpha}{2}}} \right)^2 \approx 14.$$

② 모합에 대한 추정에 표본 크기 (비편향 추정량을 사용한다)

$$\text{var}(\hat{\tau}_u) = N^2 \frac{\sigma_b^2}{n} \left(\frac{N-n}{N} \right)$$

$$Z_{\frac{\alpha}{2}} \sqrt{\text{var}(\hat{\tau}_u)} = \text{오차한계}(B) \quad \text{의 방정식을 통하여}$$

$$n = \frac{N \sigma_b^2}{ND + \sigma_b^2}, \quad D = \left(\frac{B}{Z_{\frac{\alpha}{2}} N} \right)^2$$

 $\frac{Z}{2}$ 사용한다.

$$\sigma_b^2 = \frac{1}{N-1} \sum_{i=1}^N (\tau_i - \bar{\tau})^2, \quad \bar{\tau} = \frac{\tau}{N} = \text{psu 당 평균값}$$



6.4. 층화집락 표집 (Stratified cluster sampling)

층화집락 표집은 모집단을 동질적인 층으로 층화한 다음 각 층마다 독립적인 집락 표집을 하는 방법이다.

예를 들어 서울시 교육청에서 고 3생들을 대상으로 조사하는 경우

서울에 총 600 여개의 고등학교가 있다고 한다면

각 고등학교는 지역적 동질성을 지닐 수 있으므로 고등학교를

층으로 한 모든 고등학교에서 3 학년 학생수의 $\frac{1}{3}$ 은 단순방의

선락하여 선락된 학생의 모든 학생을 대상으로 하는 조사를

생각 할 수 있다

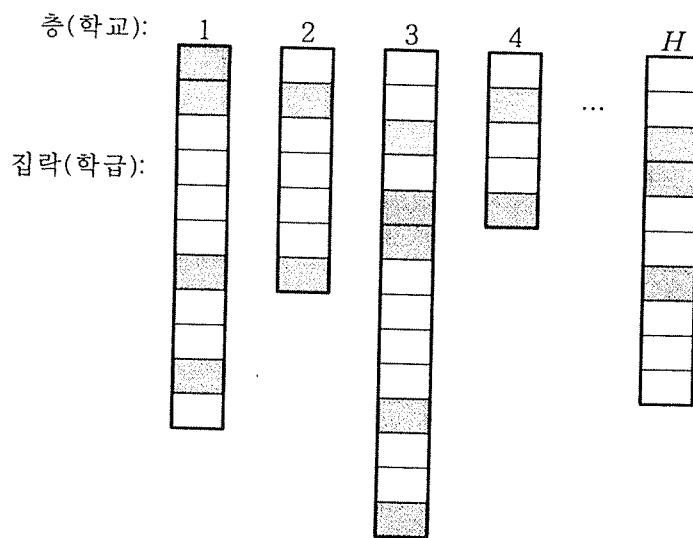


그림 6.6 층화집락표집의 도해

(1) 모함의 추정

- 모집단 전체의 총 집락 수를 N 이라 하고

총 h 의 총 집락 수를 N_h , $h=1, 2, \dots, H$ 라 하자

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

- 각 층에서 n_h 개의 집락을 비복원 단순 임의 추출 하고

해당 집락의 집락합은 $\{T_{h1}, T_{h2}, \dots, T_{hn_h}\}$ 라 하자.

- $\bar{T}_h = \frac{1}{n_h} \sum_{j=1}^{n_h} T_{hj}$ 라 하면

$\hat{\bar{T}}_h$ = h 번째 층의 총합 T_h 이 대한 추정치

$$= \frac{N_h}{n_h} \sum_{j=1}^{n_h} T_{hj} = N_h \cdot \bar{T}_h$$

- 전체 모함 T 의 추정치로는

$$\hat{T} = \sum_{h=1}^H \hat{\bar{T}}_h = \sum_{h=1}^H N_h \bar{T}_h = \sum_{h=1}^H \frac{N_h}{n_h} \sum_{j=1}^{n_h} T_{hj}.$$

∴ 사망하게 되고 보상추정량은

$$\widehat{\text{Var}}(\hat{T}) = \sum_{h=1}^H N_h^2 \frac{N_h - n_h}{N_h} \frac{S_h^2}{n_h}$$

$$S_h^2 = \frac{1}{n_h - 1} \sum_{j=1}^{n_h} (T_{hj} - \bar{T}_h)^2$$



(2) 모평균의 추정

모평균의 추정값으로 \bar{y} 를 표집관내 총 원소수로 나누는 추정값을 사용한다. 표집관 크기 K 는 일반적으로 알려져 있지 않고, K 를 추정하기 위해

$$K = \sum_{h=1}^H \sum_{j=1}^{n_h} M_{hj} \quad \text{이러} \quad \text{let} \quad M_{hj} = \sum_{i=1}^{N_{hj}} M_{hij}$$

따라서

$$\hat{M}_{hj} = N_{hj} \bar{M}_{hj} = \frac{N_{hj}}{n_{hj}} \sum_{i=1}^{n_{hj}} M_{hij}$$

이러

$$\hat{K} = \sum_{h=1}^H \hat{M}_{hj} = \sum_{h=1}^H N_{hj} \bar{M}_{hj} = \sum_{h=1}^H \frac{N_{hj}}{n_{hj}} \sum_{i=1}^{n_{hj}} M_{hij}$$

가 된다

따라서 평균 추정값

$$\hat{M}_G = \frac{\sum_{h=1}^H N_{hj} \bar{M}_{hj}}{\sum_{h=1}^H N_{hj}} = \frac{\sum_{h=1}^H \frac{N_{hj}}{n_{hj}} \sum_{i=1}^{n_{hj}} M_{hij}}{\sum_{h=1}^H \frac{N_{hj}}{n_{hj}} \sum_{i=1}^{n_{hj}} M_{hij}} \quad \text{이러}$$



이는 1차항 비 추정이므로 β 이 대각 추정량 b_c 이
 =항으로 항이고 이 추정량의 일차항

$$\begin{aligned}\widehat{\text{var}}(\hat{\mu}_c - \mu) &= \widehat{\text{var}} \left(\frac{\sum_{h=1}^H \frac{1}{N} \hat{\tau}_h}{\sum_{h=1}^H \frac{1}{N} \hat{\mu}_h} - \mu \right) \\ &= \widehat{\text{var}} \left\{ \frac{1}{\sum_{h=1}^H \frac{1}{N} \hat{\mu}_h} \left(\sum_{h=1}^H \frac{1}{N} (\hat{\tau}_h - \mu \hat{\mu}_h) \right) \right\}\end{aligned}$$

$$\approx \frac{1}{\bar{\mu}^2} \widehat{\text{var}} \left\{ \sum_{h=1}^H \frac{1}{N} (\hat{\tau}_h - \mu \hat{\mu}_h) \right\}$$

$$= \frac{1}{\bar{\mu}^2} \frac{1}{N^2} \sum_{h=1}^H \widehat{\text{var}}(\hat{\tau}_h - \mu \hat{\mu}_h)$$

$$= \frac{1}{\bar{\mu}^2} \frac{1}{N^2} \sum_{h=1}^H N_h^2 \frac{N_h - n_h}{N_h} \frac{S_{e.e.}^2}{n_h} \quad \text{이항}$$

$$\text{여기서 } S_{e.e.}^2 = \frac{1}{n_e - 1} \sum_{j=1}^{n_e} (\tau_{ej} - \hat{\mu}_c \mu_{ej})^2 \quad \text{여기}$$

6.5. 크기비례 확률 표집.

- Sampling with probabilities proportional to size (PPS)

= 복원 부등 확률 표집

- Sampling weight (표집 비중)

= 표집 확률 (sampling probability)의 역수

- 총화 추출 vs 크기비례 확률 표집

- 총화 추출에서는 각 층을 새로운 모집단으로 간주하여 층마다

독립적으로 단순임의 비복원 표집을 하기 때문에 등 확률 표집법이다.

- 층의 크기가 다르므로 표집비중이 같지 않기 때문에

이를 보정하기 위하여 층의 크기를 고려하여 추정값을 계산한다.

$$\hat{\mu}_{st} = \sum_{h=1}^H \sum_{j=1}^{n_h} \frac{N_h}{n_h} y_{hj} = \sum_{h=1}^H \sum_{j=1}^{n_h} w_h \cdot Y_{hj}$$

- 반면 PPS에서는 총화를 먼저 하지 않고

층 크기를 고려하여 ~~부등 확률 표집을 한다.~~
전체 모집단으로부터

- 일반적으로 부등 확률 표집의 경우 selection probabilities를

알지 못하고 이견현상을 selection bias라 한다.

이제 PPS를 살펴 보자.

M_i = i 번째 집락내 원소의 수

K = 모집단 총 원소의 수

$$Z_i = \frac{M_i}{K} = \text{표집 확률}$$

$$\Rightarrow \text{Sampling weight} = \frac{1}{Z_i}$$

N 개의 집락이 있고 n 개의 집락은 PPS sampling 된다면
(복원추출)

$$\hat{\tau}_{PPS} = \frac{1}{n} \sum_{i=1}^n \frac{1}{Z_i} \tau_i$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{K}{M_i} \tau_i = \frac{K}{n} \sum_{i=1}^n \frac{\tau_i}{M_i} = \frac{K}{n} \sum_{i=1}^n \mu_i$$

$$\hat{\mu}_{PPS} = \frac{1}{K} \hat{\tau}_{PPS} = \frac{1}{n} \sum_{i=1}^n \frac{\tau_i}{M_i} = \frac{1}{n} \sum_{i=1}^n \mu_i$$

$$\hat{\text{var}}(\hat{\mu}_{PPS}) = \frac{1}{n} \frac{1}{n-1} \sum_{i=1}^n (\mu_i - \hat{\mu}_{PPS})^2$$

$$\hat{\text{var}}(\hat{\tau}_{PPS}) = \frac{K^2}{n} \frac{1}{n-1} \sum_{i=1}^n (\mu_i - \hat{\mu}_{PPS})^2$$

6.7 이단 집락 표집

이단 집락 표집에서는 선택된 PSU에 속한 모든 SSU들을 조사 하였으나, 집락내 원소들을 모두 조사할 필요가 없거나, SSU에 대한 조사비용이 상대적으로 높으면 이단 집락 표집 (two-stage cluster sampling)을 고려한다. 즉 선택된 PSU에서 부분본 (sub sample)을 얻는다.

이단 집락 표집은

- ① 모집단 N 개의 PSU에서 n 개의 PSU들을 단순임의 추출한다
- ② 선택된 PSU들로 부터 크기가 M_i 인 i 번째 PSU가 선택되었다면 크기 m_i 인 단순임의 표집을 시행한다.

6.7.1. 모수의 추정

(1) 모평균의 추정

① 추정량

이단 집락 표집에서 모평균의 추정은

$$\hat{\bar{Y}}_u = \frac{N}{n} \sum_{i=1}^n \bar{Y}_i \quad \text{이웃으나 이단 집락 표집에서는}$$

집락의 모평균 \bar{Y}_i 가 관측이 되지 않으므로 추정하여야 한다.

$$\hat{\bar{Y}}_i = \frac{M_i}{m_i} \sum_{j=1}^{m_i} y_{ij} = M_i \frac{\bar{y}_i}{m_i} = M_i \bar{\bar{Y}}_i$$

따라서

$$\begin{aligned}\hat{\tau}_u &= \frac{N}{n} \sum_{i=1}^n \hat{\tau}_i = \frac{N}{n} \sum_{i=1}^n \frac{M_i}{m_i} \sum_{j=1}^{m_i} Y_{ij} \\ &= \frac{N}{n} \sum_{i=1}^n \frac{M_i}{m_i} t_i = \frac{N}{n} \sum_{i=1}^n M_i \bar{Y}_i\end{aligned}$$

② $\hat{\tau}_u$ 의 통계적 성질

$$\begin{aligned}\hat{\tau}_u &= \frac{N}{n} \sum_{i=1}^n \frac{M_i}{m_i} t_i \\ &= \frac{N}{n} \sum_{i=1}^n \frac{M_i}{m_i} \sum_{j=1}^{m_i} Y_{ij} \\ &= \frac{N}{n} \sum_{i=1}^N U_i \left\{ \frac{M_i}{m_i} \sum_{j=1}^{M_i} V_{ij} y_{ij} \right\}\end{aligned}$$

$$U_i = \begin{cases} 1 & \text{if } i \text{ 번째 cluster 가 select 된다} \\ 0 & \text{o.w.} \end{cases}$$

$$V_{ij} = \begin{cases} 1 & \text{if } y_{ij} \text{ is selected} \\ 0 & \text{o.w.} \end{cases}$$

평균

$$\begin{aligned}E(\hat{\tau}_u) &= E_u \left\{ E_v [\hat{\tau}_u | u] \right\} \\ &= E_u \left\{ \frac{N}{n} \sum_{i=1}^N U_i \frac{M_i}{m_i} \sum_{j=1}^{m_i} \frac{m_j}{M_i} y_{ij} \right\} \\ &= E_u \left\{ \frac{N}{n} \sum_{i=1}^N U_i \sum_{j=1}^{M_i} y_{ij} \right\}\end{aligned}$$

$$= \frac{N}{n} \sum_{i=1}^N \frac{n}{N} \sum_{\bar{d}=1}^{M_i} y_{i\bar{d}} = \bar{y}$$

계산

$$\begin{aligned} \text{var}(\hat{\tau}_w) &= \text{var}_{\underline{u}} \left\{ E_{\underline{v}} (\hat{\tau}_w | \underline{u}) \right\} \\ &\quad + E_{\underline{u}} \left\{ \text{var}_{\underline{v}} (\hat{\tau}_w | \underline{u}) \right\} \end{aligned}$$

Let.

$$\textcircled{1} = \text{var}_{\underline{u}} \left\{ E_{\underline{v}} (\hat{\tau}_w | \underline{u}) \right\}$$

$$\textcircled{2} = E_{\underline{u}} \left\{ \text{var}_{\underline{v}} (\hat{\tau}_w | \underline{u}) \right\}$$

먼저 ①을 계산한다.

$$\text{var}_{\underline{u}} \left\{ E_{\underline{v}} (\hat{\tau}_w | \underline{u}) \right\}$$

$$= \text{var}_{\underline{u}} \left\{ \frac{N}{n} \sum_{i=1}^N U_i \sum_{\bar{d}=1}^{M_i} y_{i\bar{d}} \right\}$$

$$= \text{var}_{\underline{u}} \left\{ \frac{N}{n} \sum_{i=1}^n \tau_i \right\} = N^2 \frac{N-n}{N-1} \frac{\sigma_b^2}{n}$$

↑
sample 원 cluster 평균에 대한
분산

$$\sigma_b^2 =$$

② $\frac{2}{n}$ 계산한다

$$\textcircled{2} = E_{\underline{u}} \left\{ \text{var}_{\underline{v}} (\hat{\tau}_u | \underline{u}) \right\}$$

$$= E_{\underline{u}} \left\{ \frac{N^2}{n^2} \sum_{i=1}^N M_i^2 U_i \text{var} \left(\frac{1}{m_i} \sum_{j=1}^{m_i} Y_{ij} \right) \right\}$$

$$= E_{\underline{u}} \left\{ \frac{N^2}{n^2} \sum_{i=1}^N M_i^2 U_i \frac{M_i - m_i}{M_i - 1} \frac{\sigma_i^2}{m_i} \right\}$$

$$\sigma_i^2 = \frac{1}{M_i - 1} \sum_{j=1}^{M_i} (y_{ij} - \mu_i)^2, \quad \mu_i = \frac{1}{M_i} \tau_i$$

$$= \frac{N^2}{n^2} \sum_{i=1}^N M_i^2 \frac{n}{N} \frac{M_i - m_i}{M_i - 1} \frac{\sigma_i^2}{m_i}$$

$$= \frac{N}{n} \sum_{i=1}^N M_i^2 \frac{M_i - m_i}{M_i - 1} \frac{\sigma_i^2}{m_i}$$

따라서

$$\text{var}(\hat{\tau}_u) = \textcircled{1} + \textcircled{2}$$

$$= N^2 \left(\frac{N-n}{N} \right) \frac{\sigma_b^2}{n} + \frac{N}{n} \sum_{i=1}^N M_i^2 \left(\frac{M_i - m_i}{M_i} \right) \frac{\sigma_i^2}{m_i}$$

$$\sigma_b^2 = \frac{1}{N-1} \sum_{i=1}^N (\tau_i - \bar{\tau})^2, \quad \left(\bar{\tau} = \frac{1}{N} \tau, \text{ PSU 당 평균} \right)$$

$$\sigma_i^2 = \frac{1}{M_i - 1} \sum_{j=1}^{M_i} (y_{ij} - \mu_i)^2, \quad \left(\mu_i = \frac{1}{M_i} \tau_i, \text{ i 번째 PSU의 평균} \right)$$

이다.

③ 분산의 추정량

분산에 대한 추정량은

$$\widehat{\text{var}}(\hat{t}_u) = N^2 \left(\frac{N-n}{N} \right) \frac{1}{n} s_b^2 + \left(\frac{N}{n} \right) \sum_{i=1}^n M_i^2 \left(\frac{M_i - m_i}{M_i} \right) \frac{s_i^2}{m_i}$$

$$s_b^2 = \frac{1}{n-1} \sum_{i=1}^n (\hat{t}_i - \bar{t})^2$$

$$t_i = \sum_{j=1}^{m_i} y_{ij}$$

$$\bar{t} = \frac{\hat{t}}{N} = \frac{1}{N} \left\{ \frac{N}{n} \sum_{i=1}^n \hat{t}_i \right\}$$

$$\hat{t}_i = \frac{M_i}{m_i} t_i$$

(2) 표평균에 대한 추정

$$\hat{\mu}_u = \frac{1}{K} \hat{t}_u = \frac{1}{K} \frac{N}{n} \sum_{i=1}^n \frac{M_i}{m_i} t_i$$

$$\widehat{\text{var}}(\hat{\mu}_u) = \frac{1}{K^2} \left\{ N^2 \left(\frac{N-n}{N} \right) \frac{s_b^2}{n} + \left(\frac{N}{n} \right) \sum_{i=1}^n M_i^2 \left(\frac{M_i - m_i}{M_i} \right) \frac{s_i^2}{m_i} \right\}$$

(3) 모비율에 대한 추정

a_i = i 번째 psu에서 관심 특성을 보유한 표본의 총 개수 among m_i

A_i = " " " " 개체들의 총 개수 among M_i

$$p_i = \frac{A_i}{M_i}$$

$$\hat{p}_i = \frac{a_i}{m_i}$$

~~$$\bar{p} = \frac{1}{n} \sum_{i=1}^n \hat{p}_i$$~~

이제

$$\rho = \frac{\sum_{i=1}^N A_i}{\sum_{i=1}^N M_i} \quad \text{이 불편추정량은 모평균의 불편추정량을 갖고있어}$$

$$\hat{P}_w = \frac{1}{K} \frac{N}{n} \sum_{i=1}^n \frac{M_i}{m_i} a_i$$

$$= \frac{1}{K} \frac{N}{n} \sum_{i=1}^n M_i \hat{p}_i = \frac{1}{K} \sum_{i=1}^n \frac{M_i}{M} \hat{p}_i$$

$$\bar{M} = \frac{1}{N} \sum_{i=1}^N M_i$$

$$\text{var}(\hat{P}_w) = \frac{1}{K^2} \left\{ N^2 \left(\frac{N-n}{N} \right) \frac{1}{n} \frac{1}{n-1} \sum_{i=1}^n \left(M_i \hat{p}_i - \frac{1}{n} \sum_{i=1}^n M_i \hat{p}_i \right)^2 \right. \\ \left. + \left(\frac{N}{n} \right) \sum_{i=1}^n M_i^2 \left(\frac{M_i - \bar{M}}{M_i} \right) \cdot \frac{1}{m_i - 1} \hat{p}_i (1 - \hat{p}_i) \right\}$$

이다.

(2) 비추정

앞의 평균 추정의 불편추정량은 $K = \sum_{i=1}^N M_i$ 라고 하였을

가정한다. 이를 모르는 경우 K 를 추정하여 앞의 불편추정량이

plug-in 하고 이는 비추정량 형태의 추정량을 제공한다.

K 에 대한 추정량은 $\hat{K} = \frac{N}{n} \sum_{i=1}^n M_i$ 이고

$$\hat{\mu}_r = \frac{1}{\hat{K}} \cdot \hat{\bar{t}}_n = \frac{\sum_{i=1}^n \frac{M_i}{m_i} t_i}{\sum_{i=1}^n M_i} = \frac{\sum_{i=1}^n M_i \bar{Y}_i}{\sum_{i=1}^n M_i} \quad \text{이다.}$$

이는 크기 N 인 모집단

$$\left\{ (M_1, \bar{t}_1), (M_2, \bar{t}_2), \dots, (M_N, \bar{t}_N) \right\}$$

으로부터 크기 n 인 표본

$$\left\{ (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \right\} \text{ 을 얻고}$$

이로부터

$$\mu = \frac{\sum_{i=1}^N \bar{t}_i}{\sum_{i=1}^N M_i} \quad \text{의 비추정량을 생각한다면}$$

$$\hat{\mu}_r = \frac{\bar{Y}}{\bar{X}} = \frac{\frac{1}{n} \sum_{i=1}^n \bar{t}_i}{\frac{1}{n} \sum_{i=1}^n M_i} = \frac{\sum_{i=1}^n M_i \bar{Y}_i}{\sum_{i=1}^n M_i} \quad \text{이다.}$$

$$\text{var}(\hat{\mu}_r) = \text{var} \left(\frac{\sum_{i=1}^n M_i \bar{Y}_i}{\sum_{i=1}^n M_i} \right)$$

$$\approx \frac{1}{M^2} \text{var} \left(\frac{1}{n} \sum_{i=1}^n M_i \bar{Y}_i \right)$$

$$\begin{aligned}
&= \frac{1}{\bar{M}^2} \frac{1}{N^2} \text{var} \left(\frac{N}{h} \sum_{i=1}^n M_i \bar{Y}_i \right) \\
&= \frac{1}{\bar{M}^2} \frac{1}{N^2} \text{var} (\hat{\tau}_u) \\
&= \frac{1}{\bar{M}^2} \frac{1}{N^2} \left\{ N^2 \left(\frac{N-n}{N} \right) \frac{\sigma_b^2}{n} + \frac{N}{n} \sum_{i=1}^n M_i^2 \left(\frac{M_i - m_i}{M_i} \right) \frac{\sigma_i^2}{m_i} \right\} \\
&= \frac{1}{\bar{M}^2} \left\{ \left(\frac{N-n}{N} \right) \frac{\sigma_b^2}{n} + \frac{1}{nN} \sum_{i=1}^n M_i^2 \left(\frac{M_i - m_i}{M_i} \right) \frac{\sigma_i^2}{m_i} \right\}
\end{aligned}$$

분산의 추정량으로

$$\begin{aligned}
\sigma_b^2 &\rightarrow s_b^2 = \frac{1}{n-1} \sum_{i=1}^n (\hat{\tau}_i - M_i \hat{\mu}_r)^2 \\
&= \frac{1}{n-1} \sum_{i=1}^n (M_i \bar{Y}_i - M_i \hat{\mu}_r)^2 = \frac{1}{n-1} \sum_{i=1}^n M_i^2 (\bar{Y}_i - \hat{\mu}_r)^2 \\
\sigma_i^2 &\rightarrow s_i^2 = \frac{1}{m_i-1} \sum_{j=1}^{m_i} (Y_{ij} - \bar{Y}_i)^2, \quad i=1, 2, \dots, n
\end{aligned}$$

모비율의 추정

추정량으로 $\hat{p}_r = \frac{\sum_{i=1}^n M_i \hat{p}_i}{\sum_{i=1}^n M_i}$ 이고 추정량의 분산에 대한 추정량은

예시 2.3

$$\text{var}(\hat{p}_r) = \frac{1}{\bar{M}^2} \left\{ \frac{N-n}{N} \frac{s_b^2}{n} + \frac{1}{nN} \sum_{i=1}^n M_i^2 \frac{M_i - m_i}{M_i} \frac{\hat{p}_i (1 - \hat{p}_i)}{m_i - 1} \right\}$$

여기서 $s_b^2 = \frac{1}{n-1} \sum_{i=1}^n [M_i \hat{p}_i - M_i \hat{p}_r]^2$

019.

MOORE

8.1.3. 동일 크기의 부표집 (sub-sample)

특별히

$$M_{\bar{i}} = M, \quad \bar{i} = 1, 2, 3, \dots, N$$

$$m_{\bar{i}} = m, \quad \bar{i} = 1, 2, 3, \dots, n$$

에 대하여 생각해 본다. 사전에 잘 계획된 실험의 경우
위의 조건이 만족될 수 있다.

이 경우

$$\begin{aligned} \hat{\mu}_u &= \frac{1}{K} \hat{t}_u \\ &= \frac{1}{K} \frac{N}{n} \sum_{\bar{i}=1}^n \left(\frac{M_{\bar{i}}}{m_{\bar{i}}} \right) t_{\bar{i}} = \frac{M}{m} \\ &= \frac{1}{\cancel{N} \cdot \cancel{M}} \frac{\cancel{N}}{n} \frac{\cancel{M}}{m} \sum_{\bar{i}=1}^n t_{\bar{i}} \\ &= \frac{1}{nm} \sum_{\bar{i}=1}^n \sum_{j=1}^m Y_{\bar{i}j} = \frac{1}{n} \sum_{\bar{i}=1}^n \left\{ \frac{1}{m} \sum_{j=1}^m Y_{\bar{i}j} \right\} \end{aligned}$$

$$\hat{\mu}_r = \frac{\sum_{\bar{i}=1}^n \cancel{M_{\bar{i}}} \bar{Y}_{\bar{i}}}{\sum_{\bar{i}=1}^n \cancel{M_{\bar{i}}}} = \frac{1}{n} \sum_{\bar{i}=1}^n \bar{Y}_{\bar{i}} = \hat{\mu}_u \text{ 이다}$$

$$\text{따라서 } \hat{\mu} \triangleq \hat{\mu}_r = \hat{\mu}_u$$

$$\widehat{\text{var}}(\hat{\mu}) = \frac{1}{M^2} \frac{N-n}{N} \frac{S_b^2}{n} + \frac{1}{nN} \frac{M-m}{M} \sum_{i=1}^n \frac{S_i^2}{m_i}$$

$$S_b^2 = \frac{1}{n-1} \sum_{i=1}^n \left(\hat{\bar{t}}_i - \frac{1}{n} \sum_{i=1}^n \hat{\bar{t}}_i \right)^2$$

$$\uparrow = \frac{M^2}{n-1} \sum_{i=1}^n \left(\bar{Y}_i - \hat{\mu} \right)^2$$

$$\hat{\bar{t}}_i = M \bar{Y}_i$$

$$S_i^2 = \frac{1}{m-1} \sum_{j=1}^m (Y_{ij} - \bar{Y}_i)^2$$

여기에서

$$SSB = \sum_{i=1}^n \sum_{j=1}^m \left(\bar{Y}_i - \hat{\mu} \right)^2 \quad df = n-1$$

$$SSW = \sum_{i=1}^n \sum_{j=1}^m \left(Y_{ij} - \bar{Y}_i \right)^2 \quad df = n(m-1)$$

이라고 정의 하면

$$S_b^2 = \frac{M^2}{n-1} \sum_{i=1}^n \left(\bar{Y}_i - \hat{\mu} \right)^2 = \frac{M^2}{n-1} \frac{SSB}{m}$$

$$S_i^2 = \frac{1}{m-1} \sum_{j=1}^m \left(Y_{ij} - \bar{Y}_i \right)^2$$

$$\begin{aligned} \text{따라서} \quad \sum_{i=1}^n S_i^2 &= \frac{1}{m-1} \sum_{i=1}^n \sum_{j=1}^m \left(Y_{ij} - \bar{Y}_i \right)^2 \\ &= \frac{1}{m-1} SSW \end{aligned}$$

따라서

$$\begin{aligned}
 \widehat{\text{var}}(\hat{\mu}) &= \frac{1}{M^2} \frac{N-n}{N} \frac{S_b^2}{n} + \frac{1}{nN} \frac{M-m}{M} \sum_{i=1}^n \frac{S_i^2}{m} \\
 &= \frac{1}{M^2} \frac{N-n}{N} \frac{1}{n} \frac{M^2}{n-1} \frac{SSB}{m} \\
 &\quad + \frac{1}{nN} \frac{M-m}{M} \frac{1}{m} \frac{1}{m-1} SSW \\
 &= \frac{N-n}{N} \cdot \frac{1}{n} \cdot \frac{1}{m} MSB + \frac{M-m}{NM} \frac{1}{m} MSW.
 \end{aligned}$$

만일 N 이 아주 크면

$$\widehat{\text{var}}(\hat{\mu}) \approx \frac{1}{nm} MSB$$

따라서 MSB 가 작은 경우, 집락간 동질적인 경우,

집락 수준이 효율적이 된다.