

4 비추정과 회귀추정

4.1 비추정(ratio estimation)

예. 총 쌀생산량(τ_y)

$$\hat{\tau}_y = \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n x_i} \tau_x = b\tau_x \quad (\text{평수: } x_i, \text{ 생산량: } y_i, \text{ 총평수: } \tau_x, \text{ 평당생산량: } b)$$

(y_i, x_i) 를 동시에 측정

$$\text{모함: } \tau_y = \sum_{i=1}^N y_i, \tau_x = \sum_{i=1}^N x_i \Rightarrow \text{표본함: } t_y = \sum_{i=1}^n y_i, t_x = \sum_{i=1}^n x_i$$

$$\text{모평균: } \mu_y = \tau_y/N, \mu_x = \tau_x/N \Rightarrow \text{표본평균: } \bar{y} = t_y/n, \bar{x} = t_x/n$$

$$\text{모비: } \beta = \tau_y/\tau_x = \mu_y/\mu_x \Rightarrow \text{표본비: } b = t_y/t_x = \bar{y}/\bar{x}$$

$$\text{모상관계수: } \rho = \frac{\sum_{i=1}^N (x_i - \mu_x)(y_i - \mu_y)/N}{\sigma_x \sigma_y} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

$$\Rightarrow \text{표본상관계수: } \gamma = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})/(n-1)}{s_y s_x} = \frac{s_{xy}}{s_x s_y}$$

$$\Rightarrow \hat{\tau}_y = b\tau_x, \hat{\mu}_y = b\mu_x$$

■ b 의 기대값과 분산(n 이 클 때)

$$E(b) \approx \beta, \text{Var}(b) \approx \frac{1}{\mu_x^2} \frac{\sigma_\epsilon^2}{n} \frac{N-n}{N}, \sigma_\epsilon^2 = \frac{\sum_{i=1}^N \epsilon_i^2}{N} = \frac{\sum_{i=1}^N (y_i - \beta x_i)^2}{N}$$

$$\widehat{\text{Var}}(b) \approx \frac{1}{\mu_x^2} \frac{s_\epsilon^2}{n} \frac{N-n}{N} \left(\approx \frac{1}{\bar{x}^2} \frac{s_\epsilon^2}{n} \frac{N-n}{N} \right), s_\epsilon^2 = \frac{\sum_{i=1}^n \epsilon_i^2}{n-1} = \frac{\sum_{i=1}^n (y_i - bx_i)^2}{n-1} (= s_y^2 + b^2 s_x^2 - 2bs_{xy})$$

$$95\% \text{ 신뢰구간: } b \pm 2s\hat{e}(b)$$

■ 모함 τ_y 에 대한 추정

$$\hat{\tau}_y = b\tau_x, \widehat{\text{Var}}(\hat{\tau}_y) = \tau_x^2 \widehat{\text{Var}}(b) \approx \tau_x^2 \frac{1}{\mu_x^2} \frac{s_\epsilon^2}{n} \frac{N-n}{N} = N^2 \frac{s_\epsilon^2}{n} \frac{N-n}{N} (\tau_x = N\mu_x)$$

$$95\% \text{ 신뢰구간: } \hat{\tau}_y \pm 2s\hat{e}(\hat{\tau}_y)$$

■ 모평균 μ_y 에 대한 추정

$$\hat{\mu}_y = b\mu_x, \widehat{\text{Var}}(\hat{\mu}_y) = \mu_x^2 \widehat{\text{Var}}(b) \approx \frac{s_\epsilon^2}{n} \frac{N-n}{N}$$

$$95\% \text{ 신뢰구간: } \hat{\mu}_y \pm 2s\hat{e}(\hat{\mu}_y)$$

■ 비율 p 에 대한 추정

$$\hat{p} = \bar{y}/\bar{x}$$

예 4.1 전나무들의 평균 나이를 추정(나무의 나이는 나무의 지름에 비례)

수목원에 있는 나무 N=1132 전나무 모두에 대해서 지름을 측정하였고, 지름의 모평균 $\mu_x = 10.3$. n=20 그루의 나무를 임의표집.

나무20	나이	지름	잔차 $e_i = y_i - bx_i$
합	$t_y = 2148$	$t_x = 188.1$	
평균	$\bar{y} = 107.4$	$\bar{x} = 9.405$	
표준편차	$s_y = 28.67$	$s_x = 1.829$	$s_e = 17.94$
비	$b = 11.42$		

$$\hat{\mu}_y = b\mu_x = 11.42(10.3) = 117.6$$

$$\widehat{Var}(\hat{\mu}_y) \approx \frac{s_e^2}{n} \frac{N-n}{N} = 15.81$$

예 4.2 1998년도 대비 2000년의 소나무 생존률을 95% 신뢰구간 추정

지역 10	00묘목수	98묘목수	잔차 $e_i = y_i - bx_i$
합	$t_y = 10753$	$t_x = 15720$	
평균	$\bar{y} = 1075.3$	$\bar{x} = 1572$	
표준편차	$s_y = 987.07$	$s_x = 1452.3$	$s_e = 45.038$
비	$b = 0.684$		

$$\widehat{Var}(b) \approx \frac{1}{\bar{x}^2} \frac{s_e^2}{n} \frac{N-n}{N} = 0.0041 \text{ (유한모집단수정을 무시)}$$

4.2 표본크기

$z_{\alpha/2} \sqrt{\text{Var}(\hat{\theta})} = B$:오차한계

$$z_{\alpha/2} \sqrt{\frac{1}{\mu_x^2} \frac{\sigma_\epsilon^2}{n} \frac{N-n}{N}} = B \Rightarrow n = \frac{N\sigma_\epsilon^2}{NB^2\mu_x^2/z_{\alpha/2}^2 + \sigma_\epsilon^2} = \frac{N\sigma_\epsilon^2}{ND + \sigma_\epsilon^2}, D = \left(\frac{B\mu_x}{z_{\alpha/2}}\right)^2$$

■ μ_x 와 σ_ϵ^2 의 추정량

$$\hat{\mu}_x = \frac{\sum_{i=1}^n x_i}{n}, \hat{\sigma}_\epsilon^2 = \frac{\sum_{i=1}^n (y_i - bx_i)^2}{n-1}$$

■ 모함 τ_y 에 대한 추정

$$n = \frac{N\sigma_\epsilon^2}{ND + \sigma_\epsilon^2}, D = \left(\frac{B}{z_{\alpha/2}N}\right)^2$$

■ 모평균 μ_y 에 대한 추정

$$n = \frac{N\sigma_\epsilon^2}{ND + \sigma_\epsilon^2}, D = \left(\frac{B}{z_{\alpha/2}}\right)^2$$

예 4.3 1,000 채의 아파트를 보유한 아파트 단지 관리인이 현재 아파트 1채당 평균 거주인이 몇 명인지 알고 싶어 표본조사를 계획.(오차한계 $B = \pm 0.5$)

1) 특이 값 존재

아파트10	올해	작년	전차 $e_i = y_i - bx_i$
합	$t_y = 43$	$t_x = 42$	
평균	$\bar{y} = 4.3$	$\bar{x} = 4.2$	
표준편차	$s_y = 1.64$	$s_x = 1.55$	$s_e = 1.612$
비	$b = 1.02 (r = 0.50)$		

$$n = \frac{N\sigma_\epsilon^2}{ND + \sigma_\epsilon^2} = 39.9, D = \left(\frac{B}{z_{\alpha/2}}\right)^2 = 1/16$$

2) 특이 값 제거

아파트9	올해	작년	전차 $e_i = y_i - bx_i$
합	$t_y = 36$	$t_x = 39$	
평균	$\bar{y} = 4.0$	$\bar{x} = 4.33$	
표준편차	$s_y = 1.41$	$s_x = 1.58$	$s_e = 0.818$
비	$b = 0.923 (r = 0.84)$		

$$n = \frac{N\sigma_\epsilon^2}{ND + \sigma_\epsilon^2} = 10.6, D = \left(\frac{B}{z_{\alpha/2}}\right)^2 = 1/16$$

4.3 층화표집에서 비추정

분리비추정(separate ratio estimation) : 각 층별로 독립적으로 비추정을 하여 층합을 추정한 후 그 층별결과를 합하는 방법

병합비추정(combined ratio estimation, 결합비추정) : 모든 층을 통틀어 단 하나의 비를 추정해서 모합을 추정하는 방법

어느 방법이든 층마다 표본크기가 커야 쓸모가 있다.

1) 분리비추정

- 모합 τ_y 에 대한 추정

$$b_h = \frac{t_{yh}}{t_{xh}} = \frac{\bar{y}_h}{\bar{x}_h}, \quad \hat{\tau}_{yh} = b_h \tau_{xh}, \quad \hat{\tau}_{yRs} = \sum_{h=1}^H \hat{\tau}_{yh} = \sum_{h=1}^H b_h \tau_{xh}$$

$$\begin{aligned} \widehat{Var}(\hat{\tau}_{yRs}) &= \widehat{Var}\left(\sum_{h=1}^H b_h \tau_{xh}\right) & s_{eh}^2 &= \frac{\sum_{i=1}^{n_h} e_{hi}^2}{n_h - 1} = \frac{\sum_{i=1}^{n_h} (y_{hi} - b_h x_{hi})^2}{n_h - 1} \\ &= \sum_{h=1}^H \tau_{xh}^2 \widehat{Var}(b_h) \\ &\approx \sum_{h=1}^H \tau_{xh}^2 \frac{1}{\mu_{xh}^2} \frac{s_{eh}^2}{n_h} \frac{N_h - n_h}{N_h} \end{aligned}$$

- 모평균 μ_y 에 대한 추정

$$\hat{\mu}_{yRs} = \sum_{h=1}^H \left(\frac{N_h}{N}\right) b_h \mu_{xh}, \quad \widehat{Var}(\hat{\mu}_{yRs}) = \sum_{h=1}^H \left(\frac{N_h}{N}\right)^2 \frac{s_{eh}^2}{n_h} \frac{N_h - n_h}{N_h}$$

2) 병합비추정

$$\begin{aligned} \hat{\tau}_y &= \sum_{h=1}^H N_h \bar{y}_h, \quad \hat{\tau}_x = \sum_{h=1}^H N_h \bar{x}_h \\ \hat{\mu}_y &= \bar{y}_{st} = \hat{\tau}_y / N, \quad \hat{\mu}_x = \bar{x}_{st} = \hat{\tau}_x / N \end{aligned}$$

- 모합 τ_y 에 대한 추정

$$b_c = \frac{\hat{\tau}_y}{\hat{\tau}_x} = \frac{\bar{y}_{st}}{\bar{x}_{st}}, \quad \hat{\tau}_{yRc} = b_c \tau_x = \sum_{h=1}^H b_c \tau_{xh}$$

$$\begin{aligned} \widehat{Var}(\hat{\tau}_{yRc}) &= \widehat{Var}(b_c \tau_x) = \sum_{h=1}^H \tau_{xh}^2 \widehat{Var}(b_c), \quad s_{eh}^2 = \frac{\sum_{i=1}^{n_h} e_{hi}^2}{n_h - 1} = \frac{\sum_{i=1}^{n_h} (y_{hi} - b_c x_{hi})^2}{n_h - 1} \\ &\approx \sum_{h=1}^H \tau_{xh}^2 \frac{1}{\mu_{xh}^2} \frac{s_{eh}^2}{n_h} \frac{N_h - n_h}{N_h} \\ &= \sum_{h=1}^H N_h^2 \frac{s_{eh}^2}{n_h} \frac{N_h - n_h}{N_h} \end{aligned}$$

- 모평균 μ_y 에 대한 추정

$$\hat{\mu}_{yRc} = b_c \mu_x, \quad \widehat{Var}(\hat{\mu}_{yRc}) = \sum_{h=1}^H \left(\frac{N_h}{N}\right)^2 \frac{s_{eh}^2}{n_h} \frac{N_h - n_h}{N_h}$$

※ 층별 표본크기가 큰 경우에는 분리비추정이 유리하고, 작은 경우는 병합비추정이 유리하다.

예4.4 농가별 벼를 재배한 평균면적 추정

y : 농가별 벼를 재배한 경지의 면적

x : 농가별 총 경작 면적

층	경지면적	N_h	\bar{y}_h	\bar{x}_h	σ_{yh}^2	σ_{xyh}	σ_{xh}^2	β_h	n_h
1	≤ 160	1580	19.40	82.56	312	494	2055	0.2350	70
2	> 160	430	51.63	244.85	922	858	7357	0.2109	30
모집단		2010	26.30	117.28	620	1453	7619	0.2242	100

1) 단순임의표집

$$\nu_1 = \frac{\sigma_y^2}{n} = 6.20$$

2) 단순임의표집에서 비추정

$$\nu_2 = \frac{\sigma_\epsilon^2}{n} = 3.51$$

3) 층화임의표집

$$\nu_3 = \sum_{h=1}^2 \left(\frac{N_h}{N} \right)^2 \frac{\sigma_{yh}^2}{n_h} = 4.16$$

4) 층화임의표집에서 분리비 추정

$$\nu_4 = \sum_{h=1}^2 \left(\frac{N_h}{N} \right)^2 \frac{\sigma_{\epsilon h}^2}{n_h} = 3.06$$

5) 층화임의표집에서 병합비 추정

$$\nu_5 = \sum_{h=1}^2 \left(\frac{N_h}{N} \right)^2 \frac{\sigma_{\epsilon c}^2}{n_h} = 3.10$$

4.4 회귀추정

$$y = \alpha + \beta x + \epsilon \Rightarrow \hat{y} = a + bx, \hat{\beta} = b = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{s_{xy}}{s_x^2} = \frac{\gamma s_y}{s_x}, \hat{\alpha} = a = \bar{y} - b\bar{x}$$

- 모평균 μ_y 에 대한 추정

$$\hat{\mu}_y = a + b\mu_x = \bar{y} + b(\mu_x - \bar{x}) \quad (\nabla \text{ 항상 } (\bar{x}, \bar{y}) \text{를 통과})$$

$$E(\hat{\mu}_y - \mu_y) = -\text{cov}(\bar{x}, b), \widehat{\text{Var}}(\hat{\mu}_y) \approx \frac{s_e^2}{n} \frac{N-n}{N} \approx \frac{MSE}{n} \frac{N-n}{N}$$

$$s_e^2 \approx \frac{\sum_{i=1}^n e_i^2}{n-1} = \frac{\sum_{i=1}^n [(y_i - \bar{y}) - b(x_i - \bar{x})]^2}{n-1}$$

- 모합 τ_y 에 대한 추정

$$\hat{\tau}_y = N \hat{\mu}_y, \widehat{\text{Var}}(\hat{\tau}_y) \approx N^2 \frac{s_e^2}{n} \frac{N-n}{N}$$

예 4.5 산림지역에서 죽은 나무의 숫자를 추정

x : 항공사진만으로 죽은 나무의 수

y : 실사를 통한 죽은 나무의 수

$$\hat{y} = 4.6581 + 0.6629x, r = 0.874, s_e^2 = 129.6904 / (24 - 1) = 5.6387$$

$$\hat{\mu}_y = a + b\mu_x = 4.6581 + 0.6629(12.2) = 12.75$$

1) 비추정

$$\widehat{\text{se}}(\hat{\mu}_y) = \sqrt{\frac{100-24}{100} \frac{5.6387}{24}} = 0.423$$

$$95\% \text{ 신뢰구간 } \hat{\mu}_y \pm t_{(n-2):\alpha/2} \widehat{\text{se}}(\hat{\mu}_y) = (11.87, 13.63)$$

2) 단순임의표집

$$\widehat{\text{se}}(\bar{y}) = \sqrt{\frac{100-24}{100} \frac{4.8812^2}{24}} = 0.869$$

3) 죽은 나무 총수

$$\hat{\tau}_y = N \hat{\mu}_y = 100(12.75) = 1275$$

$$\widehat{\text{se}}(\hat{\tau}_y) = N \widehat{\text{se}}(\hat{\mu}_y) = 100(0.423) = 42.3$$

4.5 차이추정

$\beta = 1$ 임을 아는 경우

$$y = \alpha + 1x + \epsilon \quad \Leftrightarrow \quad \hat{y} = a + x, \quad \hat{\alpha} = a = \bar{y} - \bar{x}$$

- 모평균 μ_y 에 대한 추정

$$\hat{\mu}_y = (\bar{y} - \bar{x}) + \mu_x, \quad \widehat{Var}(\hat{\mu}_y) = \widehat{Var}(\bar{e} + \mu_x) = \widehat{Var}(\bar{e}) \approx \left(\frac{N-n}{N}\right) \frac{s_e^2}{n}, \quad s_e^2 = \frac{\sum_{i=1}^n (e_i - \bar{e})^2}{n-1}$$

- 모합 τ_y 에 대한 추정

$$\hat{\tau}_y = N \hat{\mu}_y = \tau_x + N(\bar{y} - \bar{x}) = \tau_x + N\bar{e}$$

$$\widehat{Var}(\hat{\tau}_y) = N^2 \widehat{Var}(\bar{e}) \approx N^2 \left(\frac{N-n}{N}\right) \frac{s_e^2}{n}$$

예 4.6

$$\hat{\mu}_y = (\bar{y} - \bar{x}) + \mu_x = 1.1 + 100 = 101.1$$

$$\widehat{Var}(\hat{\mu}_y) \approx \left(\frac{N-n}{N}\right) \frac{s_e^2}{n} = \frac{300-10}{300} \frac{2.77}{10} = 0.268$$

4.6 층과 비-블록(block)과 공변수(covariate)