

## 제 4장: 비추정량과 회귀추정량: 보조변수를 이용한 추정량

### 4.1 개요와 정의

1. 모집단의 보조정보가 알려져 있는 경우. 연구변수와 상관이 높은 경우 보조변수( $X$ )와 연구의 대상인 변수( $Y$ )사이에 어떤 관계를 파악하여 이 관계를  $Y$ 의 모수 추정에 이용하는 방법.

$$Y \approx bX \quad \text{or} \quad Y \approx a + bX$$

2. 비추정법: 보조변수( $X$ )와 연구변수( $Y$ ) 사이에 비례관계가 성립할 때 이 비 값을 이용하여 모집단의 평균과 총계를 추정하는 방법.
3. 회귀추정법: 보조변수( $X$ )와 연구변수( $Y$ ) 사이에 선형관계가 성립할 때 이 선형관계를 이용하여 모집단의 평균과 총계를 추정하는 방법.
4. 비슷한 개념으로 사후 증화를 생각할 수 있다.
5. 예제:

- 교재의 나무의 평균나이
- 아파트단지의 평균 세대원 수
- 소비자물가지수: 두 시점사이의 물건의 구매값들의 비.  $X$  = 기준 시점의 물건의 값,  $Y$  = 현재 물건의 값.
- 병원의 의약품 소비:  $X$  = 병원의 병상수,  $Y$  = 병원에서 소비하는 의약품의 양.

### 4.2 비추정법

1. 모함:  $\tau_x, \tau_y$ ;

모평균:  $\mu_x, \mu_y$ ;

모비:

$$\beta = \frac{\tau_y}{\tau_x} = \frac{\mu_y}{\mu_x}.$$

## 2. 단순임의비복원 추출을 가정한다.

각 모집단위마다  $(x_i, y_i)$ 를 동시에 관측

크기  $N$ 인 모집단으로 부터  $n$ 개의 표본  $\{(X_i, Y_i), i = 1, 2, \dots, n\}$ 을 단순임의비복원 추출로 관측한다.

표본합:  $t_x, t_y$ ;

표본평균:  $\bar{Y}, \bar{X}$ ;

표본비:

$$b = \frac{t_y}{t_x} = \frac{\bar{Y}}{\bar{X}}.$$

따라서  $\tau_y$ 와  $\mu_y$ 에 대한 비 추정치는

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

## 3. 모집단 비 $\beta$ 에 대한 추정

- 추정량:

$$b = \frac{\sum_{i=1}^n Y_i}{\sum_{i=1}^n X_i} = \frac{\bar{Y}}{\bar{X}}.$$

- 추정량의 분산:

$$\text{var}(b) = \frac{N-n}{n(N-1)} \left( \frac{1}{\mu_x^2} \right) \sigma_r^2 \approx \frac{N-n}{nN} \left( \frac{1}{\mu_x^2} \right) \sigma_r^2$$

이고 여기서

$$\sigma_r^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \beta x_i)^2. \quad \left( \sum_{i=1}^N y_i = \beta \sum_{i=1}^N x_i \right).$$

- 분산 추정량:

$$\widehat{\text{var}}(b) = \frac{N-n}{nN} \left( \frac{1}{\mu_x^2} \right) s_r^2$$

이고 여기서

$$s_r^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - bX_i)^2.$$

#### 4. 모총계 $\tau_y$ 에 대한 추정

- 추정량:

$$\hat{\tau}_y = b\tau_x.$$

- 추정량의 분산:

$$\text{var}(\hat{\tau}_y) = \tau_x^2 \text{var}(b) \approx N^2 \mu_x^2 \frac{N-n}{nN} \left( \frac{1}{\mu_x^2} \right) \sigma_r^2 = \frac{N(N-n)}{n} \sigma_r^2.$$

- 분산 추정량:

$$\widehat{\text{var}}(\hat{\tau}_y) = \tau_x^2 \widehat{\text{var}}(b) = \frac{N(N-n)}{n} s_r^2$$

이다.

#### 5. 모평균 $\mu_y$ 에 대한 추정

- 추정량:

$$\hat{\mu}_y = b\mu_x.$$

- 추정량의 분산:

$$\text{var}(\hat{\mu}_y) = \mu_x^2 \text{var}(b) = \mu_x^2 \frac{N-n}{nN} \left( \frac{1}{\mu_x^2} \right) \sigma_r^2 = \frac{N-n}{nN} \sigma_r^2.$$

- 분산 추정량:

$$\widehat{\text{var}}(\hat{\mu}_y) = \mu_x^2 \widehat{\text{var}}(b) = \frac{N-n}{nN} s_r^2$$

이다.

6. 가중선형회귀를 이용한 비추정량의 계산:

7. [예제 4.1](#)

### 4.3 표본의 크기

표본의 크기는 오차한계( $B$ )를 이용하여 결정한다. 즉,

$$Z_{\alpha/2} \sqrt{\text{var}(\hat{\theta})} = B$$

를 만족시키는  $n$ 을 계산한다.

- 모비  $\beta$ 에 대한 표본크기:

$$n = \frac{N\sigma_r^2}{ND + \sigma_r^2} \quad D = \left( \frac{B\mu_x}{Z_{\alpha/2}} \right)^2.$$

- 모합  $\tau_y$ 에 대한 표본크기:

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

- 모합  $\mu_y$ 에 대한 표본크기:

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

- [예제 4.3](#)

## 4.4 SRS와 비추정량의 효율성 비교

1.  $r_i \equiv y_i - \beta x_i$ 를 정의 한다. 즉,  $y_i = \beta x_i + r_i$  이고  $\sum_{i=1}^N r_i = 0$ .

2. 각 추정량의 분산을 살펴보면

$$\text{var}(\hat{\mu}_y^{\text{srs}}) = \frac{\sigma_y^2}{n} \frac{N-n}{N-1} \approx \frac{\sigma_y^2}{n} \frac{N-n}{N}$$

과

$$\text{var}(\hat{\mu}_y^{\text{ratio}}) = \text{var}(\mu_x b) = \frac{\sigma_r^2}{n} \frac{N-n}{N}$$

이다.

3. 여기서

$$\begin{aligned} \sigma_y^2 &= \frac{1}{N} \sum_{i=1}^N (y_i - \mu_y)^2 \\ &= \frac{1}{N} \sum_{i=1}^N (y_i - \beta x_i + \beta x_i - \mu_y)^2 \\ &= \frac{1}{N} \sum_{i=1}^N \{r_i^2 - 2r_i(x_i - \mu_x)\beta + \beta^2(x_i - \mu_x)^2\} \\ &= \sigma_r^2 + \beta^2 \sigma_x^2 - 2\beta \sigma(r, x) \end{aligned}$$

이고

$$\sigma(r, x) = \frac{1}{N} \sum_{i=1}^N r_i(x_i - \mu_x)$$

이다.

4.

$$\text{RE} = \frac{\text{var}(\hat{\mu}_y^{\text{srs}})}{\text{var}(\hat{\mu}_y^{\text{ratio}})} = \frac{\sigma_y^2}{\sigma_r^2} = 1 + \beta^2 \frac{\sigma_x^2}{\sigma_r^2} - 2\beta \frac{\sigma_x}{\sigma_r} \text{corr}(r, x)$$

이고

$$\text{corr}(r, x) = \frac{\sigma(r, x)}{\sigma_r \sigma_x}.$$

5. 따라서 위의 결과로부터

- $r$ 와  $x$ 의 correlation이 약할 수록/ 독립에 가까울 수록 (대체적으로  $x$ 와  $y$ 의 correlation이 강할 수록) ratio estimator의 SRS estimator에 대한 상대 효율성이 좋아진다.
- $r_i$ 와  $x_i$ 의 비상관성/독립성 은 회귀분석에서와 같이  $R_i = Y_i - bX_i$ 와  $X_i$ 의 잔차도를 통하여 확인할 수 있다.

층화표집에서의 비추정량으로는 다음의 두 가지 형태가 있다.

- 분리 비추정량 (separate ratio estimator)
- 병합 비추정량 (combined ratio estimator)

## 4.5 층화표집 - 분리비추정량

1. 모합  $\tau_y$ 에 대한 분리 비추정량:

- 추정량:

$$b_h = \frac{t_{yh}}{t_{xh}} = \frac{\bar{Y}_h}{\bar{X}_h},$$

$$\hat{\tau}_{yh} = b_h \tau_{xh},$$

$$\hat{\tau}_{yRs} = \sum_{h=1}^H \hat{\tau}_{yh} = \sum_{h=1}^H b_h \tau_{xh}.$$

- 추정량  $\hat{\tau}_{yRs}$ 의 분산:

$$\begin{aligned} \text{var}(\hat{\tau}_{yRs}) &= \sum_{h=1}^H \tau_{xh}^2 \text{var}(b_h) \\ &\approx \sum_{h=1}^H \tau_{xh}^2 \frac{1}{\mu_{xh}^2} \frac{\sigma_{rh}^2}{n_h} \frac{N_h - n_h}{N_h - 1}. \end{aligned} \quad (1)$$

- 분산  $\text{var}(\hat{\tau}_{yRs})$ 의 추정량:

$$\widehat{\text{var}}(\hat{\tau}_{yRs}) = \sum_{h=1}^H \tau_{xh}^2 \frac{1}{\mu_{xh}^2} \frac{s_{rh}^2}{n_h} \frac{N_h - n_h}{N_h}$$

이고

$$s_{rh}^2 = \frac{1}{n_h - 1} \sum_{j=1}^{n_h} (Y_{hj} - b_h X_{hj})^2.$$

2. 모평균  $\mu_y$ 에 대한 분리 비추정량:

- 추정량:

$$\hat{\mu}_{yRs} = \frac{1}{N} \hat{\tau}_{yRs} = \sum_{h=1}^H \left( \frac{N_h}{N} \right) \left( \frac{1}{N_h} \right) b_h \tau_{xh} = \sum_{h=1}^H W_h b_h \mu_{xh}$$

- 추정량  $\hat{\mu}_{yRs}$ 의 분산:

$$\begin{aligned} \text{var}(\hat{\mu}_{yRs}) &= \sum_{h=1}^H W_h^2 \mu_{xh}^2 \text{var}(b_h) \\ &\approx \sum_{h=1}^H W_h^2 \mu_{xh}^2 \frac{1}{\mu_{xh}^2} \frac{\sigma_{rh}^2}{n_h} \frac{N_h - n_h}{N_h - 1} \\ &= \sum_{h=1}^H W_h^2 \frac{\sigma_{rh}^2}{n_h} \frac{N_h - n_h}{N_h - 1} \end{aligned}$$

- 분산  $\text{var}(\hat{\mu}_{yRs})$ 의 추정량:

$$\widehat{\text{var}}(\hat{\mu}_{yRs}) = \sum_{h=1}^H W_h^2 \frac{s_{rh}^2}{n_h} \frac{N_h - n_h}{N_h}$$

이고

$$s_{rh}^2 = \frac{1}{n_h - 1} \sum_{j=1}^{n_h} (Y_{hj} - b_h X_{hj})^2.$$

## 4.6 층화표집 - 병합비추정량

$$\beta = \frac{\tau_y}{\tau_x}$$

를 이용

$$b_c = \frac{\hat{\tau}_y}{\hat{\tau}_x} = \frac{\bar{Y}_{st}}{\bar{X}_{st}}$$

의 추정량을 사용한다.

- 모합  $\tau_y$ 에 대한 병합 비추정량:

- 병합추정량:

$$\hat{\tau}_{yRc} = b_c \tau_x.$$



- 추정량  $\hat{\mu}_{yRc}$  의 분산:

$$\text{var}(\hat{\tau}_{yRc}) = \tau_x^2 \widehat{\text{var}}(b_c).$$

이제  $\widehat{\text{var}}(b_c)$  를 계산하여 보면

$$\text{var}(b_c) = \text{var}(b_c - \beta)$$

이고

$$\begin{aligned} b_c - \beta &= \frac{\bar{Y}_{\text{st}}}{\bar{X}_{\text{st}}} - \beta \\ &= \frac{1}{\bar{X}_{\text{st}}} (\bar{Y}_{\text{st}} - \beta \bar{X}_{\text{st}}) \\ &\approx \frac{1}{\mu_x} (\bar{Y}_{\text{st}} - \beta \bar{X}_{\text{st}}) \\ &= \frac{1}{\mu_x} \left\{ \sum_{h=1}^H \frac{N_h}{N} (\bar{Y}_h - \beta \bar{X}_h) \right\}. \end{aligned} \quad (2)$$

따라서

$$\begin{aligned} \text{var}(b_c) &= \text{var}(b_c - \beta) \\ &\approx \frac{1}{\mu_x^2} \sum_{h=1}^H \left( \frac{N_h}{N} \right)^2 \frac{\sigma_{rh}^2}{n_h} \frac{N_h - n_h}{N_h - 1} \\ &\approx \frac{1}{\mu_x^2} \sum_{h=1}^H \left( \frac{N_h}{N} \right)^2 \frac{\sigma_{rh}^2}{n_h} \frac{N_h - n_h}{N_h}. \end{aligned}$$

이고

$$\begin{aligned} \text{var}(\hat{\tau}_{yRc}) &= \tau_x^2 \text{var}(b_c) \\ &\approx \tau_x^2 \cdot \frac{1}{\mu_x^2} \sum_{h=1}^H \left( \frac{N_h}{N} \right)^2 \frac{\sigma_{rh}^2}{n_h} \frac{N_h - n_h}{N_h - 1}. \end{aligned} \quad (3)$$

- 분산  $\text{var}(\hat{\mu}_{yRc})$  의 추정량:

$$\begin{aligned} \widehat{\text{var}}(\hat{\tau}_{yRc}) &= \tau_x^2 \widehat{\text{var}}(b) \\ &= \tau_x^2 \cdot \frac{1}{\mu_x^2} \sum_{h=1}^H \left( \frac{N_h}{N} \right)^2 \frac{s_{rh}^2}{n_h} \frac{N_h - n_h}{N_h}. \end{aligned} \quad (4)$$

2. 모평균  $\mu_y$ 에 대한 병합 비추정량:

- 추정량:

$$\hat{\mu}_{yRc} = b_c \mu_x.$$

- 추정량의 분산

$$\text{var}(\hat{\mu}_{yRc}) \approx \sum_{h=1}^H \left( \frac{N_h}{N} \right)^2 \frac{\sigma_{rh}^2}{\mu_h} \frac{N_h - n_h}{N_h - 1}.$$

- 분산의 추정량:

$$\widehat{\text{var}}(\hat{\mu}_{yRc}) \approx \sum_{h=1}^H \left( \frac{N_h}{N} \right)^2 \frac{s_{rh}^2}{\mu_h} \frac{N_h - n_h}{N_h}.$$

#### 4.7 분리 비추정량과 병합 비 추정량의 비교

	분리비 추정법	병합비 추정법
층별모비	--	일정
$x$ 의 층별모합	알고있음	--
$x$ 의 전체모합	알고있음	알고있음
표본크기	커야함	--

#### 4.7 회귀추정량

1. 모집단  $\{(x_i, y_i), i = 1, 2, \dots, N\}$ 에서  $n$ 개의 표본  $\{(X_j, Y_j), j = 1, 2, \dots, n\}$ 을 얻는다.

변수  $x_i$ 와  $y_i$ 가 선형의 관계  $y = \alpha + \beta x$ 를 가지고 있음을 가정한다.

2. 회귀계수  $\alpha$ 와  $\beta$ 에 대한 추정량은

$$\begin{aligned} \hat{\beta} = b &= \frac{\sum_{j=1}^n (X_j - \bar{X})(Y_j - \bar{Y})}{\sum_{j=1}^n (X_j - \bar{X})^2} = \frac{s_{XY}}{s_X^2} = \frac{Rs_Y}{s_X}, \\ \hat{\alpha} = a &= \bar{Y} - b\bar{X}. \end{aligned}$$

3. 따라서  $y$ 의 평균  $\mu_y$ 의 회귀 추정량은

$$\hat{\mu}_y^{\text{reg}} = a + b \cdot \mu_x = \bar{Y} + b \cdot (\mu_x - \bar{X}).$$

4. 편(bias):

$$\begin{aligned} \hat{\mu}_y^{\text{reg}} - \mu_y &= \bar{Y} + b(\mu_x - \bar{X}) - \mu_y \\ &= \bar{Y} - \mu_y + (b - \beta)(\mu_x - \bar{X}) + \beta(\bar{X} - \mu_x). \end{aligned}$$

따라서

$$\begin{aligned} E(\hat{\mu}_y - \mu_y) &= E[(b - \beta)(\mu_x - \bar{X})] = -\text{cov}(b, \bar{X}) \\ &= -E \left[ \frac{\sum_{i=1}^n R_i (X_i - \bar{X})(\bar{X} - \mu_x)}{\sum_{i=1}^n (X_i - \bar{X})^2} \right] \\ &\approx -\frac{E \left[ \sum_{i=1}^n R_i (X_i - \bar{X})(\bar{X} - \mu_x) \right]}{n\sigma_x^2} \end{aligned}$$

가 된다. 단,  $R_i = Y_i - \alpha - \beta X_i$ .

5. 추정량의 분산:

$$\begin{aligned} \begin{pmatrix} a \\ b \end{pmatrix} &= \begin{pmatrix} n & \sum_{i=1}^n X_i \\ \sum_{i=1}^n X_i & \sum_{i=1}^n X_i^2 \end{pmatrix}^{-1} \begin{pmatrix} \sum_{i=1}^n Y_i \\ \sum_{i=1}^n X_i Y_i \end{pmatrix} \\ &= \left( \sum_{i \in A} \tilde{x}_i \tilde{x}_i^\top \right)^{-1} \left( \sum_{i \in A} \tilde{x}_i y_i \right). \end{aligned}$$

$\tilde{X}_i = (1, X_i)^\top$ 로 정의되고  $A$ 는 표본의 index set.

이제  $\mu_y$  의 회귀추정량은

$$\begin{aligned}
\hat{\mu}_y^{\text{reg}} &= (1, \mu_x) \left( \sum_{i \in A} \tilde{x}_i \tilde{x}_i^\top \right)^{-1} \left( \sum_{i \in A} \tilde{x}_i y_i \right) \\
&= \mu_y + (1, \mu_x) \left\{ \left( \sum_{i \in A} \tilde{x}_i \tilde{x}_i^\top \right)^{-1} \sum_{i \in A} \tilde{x}_i y_i - \begin{pmatrix} \alpha \\ \beta \end{pmatrix} \right\} \\
&= \mu_y + (1, \mu_x) \left\{ \left( \sum_{i \in A} \tilde{x}_i \tilde{x}_i^\top \right)^{-1} \sum_{i \in A} \tilde{x}_i (y_i - \alpha - \beta x_i) \right\} \\
&= \mu_y + (1, \mu_x) \left( \sum_{i \in A} \tilde{x}_i \tilde{x}_i^\top \right)^{-1} \sum_{i \in A} \tilde{x}_i (y_i - \alpha - \beta x_i) \\
&\approx \mu_y + (1, \mu_x) \begin{pmatrix} n & n\mu_x \\ n\mu_x & n(\mu_x^2 + \sigma_x^2) \end{pmatrix}^{-1} \sum_{i \in A} \tilde{x}_i (y_i - \alpha - \beta x_i) \\
&= \mu_y + (1, \mu_x) \begin{pmatrix} n & n\mu_x \\ n\mu_x & n(\mu_x^2 + \sigma_x^2) \end{pmatrix}^{-1} \begin{pmatrix} \sum_{i \in A} (y_i - \alpha - \beta x_i) \\ \sum_{i \in A} x_i (y_i - \alpha - \beta x_i) \end{pmatrix} \\
&= \mu_y + \frac{1}{n} \sum_{i \in A} (y_i - \alpha - \beta x_i). \tag{5}
\end{aligned}$$

가 된다. 따라서

$$\text{var}(\hat{\mu}^{\text{reg}}) \approx \frac{\sigma_e^2}{n} \frac{N-n}{N-1}, \quad \sigma_e^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \alpha - \beta x_i)^2$$

을 얻을 수 있다.

## 4.8 차이검정