

$$\bar{x} - R\bar{y} = \bar{x} - \frac{\bar{x}}{\bar{y}} \cdot \bar{y} = 0$$

$$\Rightarrow s_d^2 \approx \frac{1}{n-1} \sum_i^n (x_i - R y_i)^2$$

และโดยปกติแล้วเราไม่อาจทราบค่าของ \bar{Y} ดังนั้นจึงแทนที่ \bar{Y} ด้วยค่าประมาณของ \bar{Y} คือ \bar{y} ดังนั้น

$$\hat{V}(R) \approx \frac{1}{n\bar{y}^2} \frac{N-n}{N} \cdot \frac{1}{n-1} \sum_i^n (x_i - R y_i)^2 \quad \text{โดยที่ } R = \bar{x}/\bar{y}$$

สำหรับสูตรดังกล่าวนี้สามารถจัดให้เป็นรูปที่สะดวกต่อการคำนวณ (Computation Formular) ได้ดังนี้

$$\hat{V}(R) \approx \frac{1}{n\bar{y}^2} \frac{N-n}{N} \frac{1}{n-1} \left\{ \sum_i^n x_i^2 - 2R \sum_i^n x_i y_i + R^2 \sum_i^n y_i^2 \right\}$$

ตัวอย่าง 2.3 จากการสำรวจจำนวนนักเรียนในโรงเรียนมัธยมศึกษาในเขตกทม. ซึ่งมีอยู่ทั้งสิ้น 468 โรงเรียน (ข้อมูลสมมุติ) การสำรวจโดยใช้วิธี SRS โดยสุ่มตัวอย่างโรงเรียนมาจากกรอบตัวอย่าง 100 โรงเรียน ปรากฏว่าเป็นโรงเรียนราษฎร์ 46 โรงเรียน รัฐบาล 54 โรงเรียน ผลจากการสำรวจปรากฏดังนี้

	n	จำนวน ครู	จำนวน นักเรียน	Σy^2	Σx^2	Σxy
ร.ร. รัฐบาล	54	2,024	31,281	111,090	29,881,219	1,729,349
ร.ร. ราษฎร์	46	1,075	13,707	33,119	6,366,785	431,041

- ก. ในโรงเรียนแต่ละประเภทจงกะประมาณอัตราส่วนระหว่าง นักเรียน-ครู
 ข. จงกะประมาณอัตราส่วนนักเรียน-ครูในข้อ ก. ด้วยช่วงเชื่อมั่น 95%
 ค. จงกะประมาณอัตราส่วนระหว่างนักเรียน-ครู ในเขต กทม. พร้อมทั้งกะประมาณ
 อัตราส่วนดังกล่าวด้วยช่วงเชื่อมั่น 99%

วิธีทำ

ก. ให้ R_1 = อัตราส่วนระหว่างนักเรียน-ครูในโรงเรียนรัฐบาล

$$= \frac{\text{จำนวนนักเรียนในโรงเรียนรัฐบาล}}{\text{จำนวนครูในโรงเรียนรัฐบาล}} = \frac{\sum_i^N x_i}{\sum_i^N y_i}$$

R_2 = อัตราส่วนระหว่างนักเรียน-ครูในโรงเรียนราษฎร์

$$= \frac{\text{จำนวนนักเรียนในโรงเรียนราษฎร์}}{\text{จำนวนครูในโรงเรียนราษฎร์}} = \frac{\sum_i^N x_i}{\sum_i^N y_i}$$

โดยอาศัยข้อมูลจากกลุ่มตัวอย่างโรงเรียนขนาด $n = 100$ โรงเรียน จะพบว่า

$$\hat{R}_1 = \frac{54}{\sum_i x_i} / \frac{54}{\sum_i y_i} = 31,281/2,024 = 15.455$$

$$\hat{R}_2 = \frac{46}{\sum_i x_i} / \frac{46}{\sum_i y_i} = 13,707/1,075 = 12.751$$

จากผลการสำรวจชี้ให้เห็นดังนี้คือ

1. อัตราส่วนระหว่างนักเรียน-ครูในโรงเรียนรัฐบาล คือครูหนึ่งคนต่อนักเรียน
 ประมาณ 15 คนหรือครู 1000 คนต่อนักเรียน 15,455 คน
2. อัตราส่วนระหว่างนักเรียน-ครูในโรงเรียนราษฎร์คือครูหนึ่งคนต่อนักเรียน
 ประมาณ 13 คน หรือครู 1000 คนต่อนักเรียน 12,751 คน

ข. การกะประมาณค่า R_1 และ R_2 ด้วยช่วงเชื่อมั่น

1. ในโรงเรียนรัฐบาล (R_1)

$$\hat{V}(\hat{R}_1) = \frac{1}{n\bar{y}^2} \cdot \frac{N-n}{N} \cdot \frac{1}{n-1} \left(\sum_i^n x_i^2 - 2\hat{R}_1 \sum_i^n x_i y_i + \hat{R}_1^2 \sum_i^n y_i^2 \right)$$

$$N = 468, n = 54, \hat{R}_1 = 15.455, \hat{R}_1^2 = 238.857$$

$$\bar{y} = \text{อัตราเฉลี่ยจำนวนครูโรงเรียนรัฐบาลในกทม.ต่อ 1 โรงเรียน} = \frac{2024}{54} = 37.48$$

$$\bar{y}^2 = 1,404.86 \quad \Sigma x_i^2 = 29,881,219$$

$$\Sigma y_i^2 = 111,090 \quad \Sigma x_i y_i = 1,729,349$$

$$\text{ดังนั้น } \hat{V}(\hat{R}_1) = \frac{1}{(54)(1404.86)} \cdot \frac{468-54}{468} \cdot \frac{1}{53} \{29881219 - 2(15.455)(1729349) + (238.857)(111090)\}$$

$$= .6516$$

$$s_{\hat{R}_1}^2 = .8072$$

$$\text{ดังนั้น } 13.87 < R_1 < 17.04$$

หรือสามารถเชื่อถือได้ถึง 95% ว่าอัตราส่วนจริงระหว่างนักเรียน-ครูของโรงเรียนรัฐบาลในเขตกทม.จะปรากฏอยู่ระหว่างครู 1 คนต่อนักเรียนประมาณ 14 คนถึง 17 คน

2. ในโรงเรียนราษฎร์ (R_2)

$$\hat{V}(\hat{R}_2) = \frac{1}{n\bar{y}^2} \cdot \frac{N-n}{N} \cdot \frac{1}{n-1} \left(\sum_i^n x_i^2 - 2\hat{R}_2 \sum_i^n x_i y_i + \hat{R}_2^2 \sum_i^n y_i^2 \right)$$

$$N = 468, n = 46, \hat{R}_2 = 12.751, \hat{R}_2^2 = 162.588$$

$$\begin{aligned}\bar{y} &= \text{อัตราเฉลี่ยจำนวนครูโรงเรียนราษฎร์ในกทม.ต่อ 1 โรงเรียน} \\ &= \frac{1075}{46} = 23.37 \quad \bar{y}^2 = 546.137, \quad \sum x_i^2 = 6,366,785\end{aligned}$$

$$\sum_i^n y_i^2 = 33,119, \quad \sum_i^n x_i y_i = 431,041$$

$$\hat{V}(R_2) = \frac{1}{(46)(546.137)} \cdot \frac{468 - 46.1}{468} \cdot \frac{1}{45} \{6366785 - 2(12.751)(431041) + (162.588)(33119)\}$$

$$s = .7781$$

$$\text{ดังนั้น } 11.22 < R_2 < 14.28$$

หรือสามารถเชื่อถือได้ถึง 95% ว่าอัตราส่วนระหว่างนักเรียน-ครูของโรงเรียนราษฎร์ในเขตกทม. จะปรากฏอยู่ในระหว่างครู 1 คนต่อนักเรียนประมาณ 11 ถึง 14 คน

สำหรับข้อ ค. นั้น จะขอเว้นไว้เป็นแบบฝึกหัด การวิเคราะห์และอภิปรายผลสามารถกระทำได้โดยง่ายเช่นเดียวกับข้อ ก. และข้อ ข.

หมายเหตุ สำหรับกรณีของอัตราส่วน W.G. Cochran แนะนำว่า ถ้า $n > 30$ และ C.V. ของทั้ง \bar{x} และ \bar{y} มีค่าไม่เกิน 0.1 กล่าวคือ $\frac{s_x}{\bar{x}} \leq 0.1$ และ $\frac{s_y}{\bar{y}} \leq 0.1$ แล้ว Sampling Distribution ของ \hat{R} จะมีการแจกแจงแบบปกติโดยประมาณคือ $\hat{R} \sim N(R, V(\hat{R}))$ ด้วยเหตุนี้เราจึงสามารถประมาณค่าพารามิเตอร์ R ด้วยช่วงเชื่อมั่นได้ ทั้งยังสามารถทดสอบสมมุติฐานเกี่ยวกับพารามิเตอร์ R ได้อีกด้วย

กรณีการประมาณค่าด้วยอัตราส่วนนี้ นักศึกษาจะพบว่า \hat{R} เป็นตัวประมาณค่าที่มีอคติของ R หรือ $E(\hat{R}) \neq R$ เว้นแต่เมื่อขนาดตัวอย่าง n ใหญ่มาก เราจึงขอเสนอให้ใช้ \hat{R} เป็นตัวประมาณค่าของ R ได้

ปริมาณของความเียงเน (Bias) หรืออคติของ \hat{R} วัดได้ดังนี้

$$\therefore \hat{R} = \frac{\bar{x}}{\bar{y}}$$

$$\hat{R} = \frac{\bar{x}}{\bar{y}} = \frac{\bar{X}}{\bar{Y}} \left(1 + \frac{\bar{x} - \bar{X}}{\bar{X}}\right) \left(1 + \frac{\bar{y} - \bar{Y}}{\bar{Y}}\right)^{-1}$$

โดยอาศัยการกระจายเทอมของเทเลอร์

$$\left(1 + \frac{\bar{y} - \bar{Y}}{\bar{Y}}\right)^{-1} = 1 - \frac{\bar{y} - \bar{Y}}{\bar{Y}} + \frac{(\bar{y} - \bar{Y})^2}{\bar{Y}^2} - \dots$$

$$\text{ดังนั้น } \hat{R} = \frac{\bar{X}}{\bar{Y}} \left(1 + \frac{\bar{x} - \bar{X}}{\bar{X}}\right) \left(1 - \frac{\bar{y} - \bar{Y}}{\bar{Y}} + \frac{(\bar{y} - \bar{Y})^2}{\bar{Y}^2} - \dots\right)$$

$$= R \left(1 + \frac{\bar{x} - \bar{X}}{\bar{X}} - \frac{\bar{y} - \bar{Y}}{\bar{Y}} + \frac{(\bar{y} - \bar{Y})^2}{\bar{Y}^2} - \frac{(\bar{x} - \bar{X})(\bar{y} - \bar{Y})}{\bar{X}\bar{Y}} + \dots\right)$$

$$\therefore E(\hat{R}) = R + R \cdot \frac{1}{\bar{X}} E(\bar{x} - \bar{X}) - R \cdot \frac{1}{\bar{Y}} E(\bar{y} - \bar{Y}) + R \frac{1}{\bar{Y}^2} E(\bar{y} - \bar{Y})^2$$

$$- R \cdot \frac{1}{\bar{X}\bar{Y}} E(\bar{x} - \bar{X})(\bar{y} - \bar{Y})$$

$$= R + 0 - 0 + \frac{R}{\bar{Y}^2} V(\bar{y}) - \frac{R}{\bar{X}\bar{Y}} E(\bar{x} - \bar{X})(\bar{y} - \bar{Y})$$

$$= R + \frac{R}{\bar{Y}^2} \cdot \frac{N-n}{N} \cdot \frac{S_y^2}{n} - \frac{R}{\bar{X}\bar{Y}} \cdot \frac{N-n}{N} \cdot \rho \cdot \frac{S_x S_y}{n}$$

¹ ลองกระจายกลับจะพบว่า $\frac{\bar{X}}{\bar{Y}} \left(1 + \frac{\bar{x} - \bar{X}}{\bar{X}}\right) \left(1 + \frac{\bar{y} - \bar{Y}}{\bar{Y}}\right)^{-1}$

$$= \frac{\bar{X}}{\bar{Y}} \cdot \frac{\bar{X} + \bar{x} - \bar{X}}{\bar{X}} \cdot \frac{\bar{Y}}{\bar{Y} + \bar{y} - \bar{Y}}$$

$$= \frac{\bar{X}}{\bar{Y}} \cdot \frac{\bar{x}}{\bar{X}} \cdot \frac{\bar{Y}}{\bar{y}} = \frac{\bar{x}}{\bar{y}}$$

โดยที่ ρ คือสหสัมพันธ์ระหว่างตัวแปร x และตัวแปร y ¹

$$\text{หรือ } \rho = \rho_{xy} = \rho_{yx} = \frac{\Sigma(x - \bar{X})(y - \bar{Y})}{\sqrt{\Sigma(x - \bar{X})^2 \Sigma(y - \bar{Y})^2}}$$

$$S_x^2 = \frac{1}{N-1} \sum_i^N (x - \bar{X})^2, S_y^2 = \frac{1}{N-1} \sum_i^N (y - \bar{Y})^2$$

$$\begin{aligned} \text{ดังนั้น } E(\hat{R}) &= R + \frac{N-n}{N} \cdot \frac{S_x^2}{n} \frac{R}{\bar{Y}^2} - \frac{N-n}{N} \rho \frac{S_x S_y}{n} \cdot \frac{1}{\bar{X}\bar{Y}} \frac{\bar{X}}{\bar{Y}} \\ &\quad : \text{แทนที่ } R \text{ ด้วย } \frac{\bar{X}}{\bar{Y}} \\ &= R + \frac{N-n}{N} \frac{S_x^2}{n} \frac{R}{\bar{Y}^2} - \frac{N-n}{N} \rho \frac{S_x S_y}{n} \cdot \frac{1}{\bar{Y}^2} \end{aligned}$$

$$\begin{aligned} \text{ดังนั้น Bias คือ } E(\hat{R}) - R &= E(\hat{R} - R) \\ &= \frac{N-n}{N} \cdot \frac{S_x^2}{n} \cdot \frac{R}{\bar{Y}^2} - \frac{N-n}{N} \rho \frac{S_x S_y}{n} \cdot \frac{1}{\bar{Y}^2} \\ \Rightarrow E(\hat{R} - R) &= \frac{N-n}{Nn} \frac{1}{\bar{Y}^2} (RS_x^2 - \rho S_x S_y) \end{aligned}$$

นั่นคือปริมาณของความเียงแตรงเท่ากับ $\frac{N-n}{Nn} \cdot \frac{1}{\bar{Y}^2} (RS_x^2 - \rho S_x S_y)$

¹ การพิสูจน์ว่า $E(\bar{x} - \bar{X})(\bar{y} - \bar{Y}) = \frac{N-n}{nN} \rho S_x S_y$ ได้แสดงไว้ท้ายตอนนีแล้ว

² ถ้า $n \rightarrow \infty$ แล้ว $\frac{N-n}{Nn} \rightarrow 0$ เพราะตัวหารมีค่าสูง แสดงว่า ถ้าใช้กลุ่มตัวอย่างให้มีขนาดใหญ่ แล้วปริมาณความเียงแตรงจะเข้าใกล้ 0 หรือ $E(\hat{R} - R) \rightarrow 0$ หรือ หรือ $E(\hat{R}) \rightarrow R$

ซึ่ง \hat{R} จะเป็นตัวประมาณค่าที่ปราศจากอคติของ R ก็ต่อเมื่อ

$$\frac{N-n}{Nn} \frac{1}{\bar{Y}^2} (RS_y^2 - eS_x S_y) = 0$$

นั่นคือ ปริมาณของความเียงแจะจะมีค่าเท่ากับ $\frac{N-n}{Nn} \frac{1}{\bar{Y}^2} (RS_y^2 - eS_x S_y)$ ซึ่ง \hat{R} จะเป็นตัวประมาณค่าที่ปราศจากอคติของ R ก็ต่อเมื่อ

$$\frac{N-n}{Nn} \frac{1}{\bar{Y}^2} (RS_y^2 - eS_x S_y) = 0$$

$$\text{หรือ } RS_y^2 = eS_x S_y$$

$$\text{หรือ } RS_y = eS_x$$

การอภิปรายผลข้างต้นโดยอาศัยกราฟจะทำให้มองเห็นภาพที่ชัดเจนขึ้น ดังนี้

พิจารณาสมการถดถอย $X_i = \beta_0 + \beta_1 Y_i$ เมื่อ Y เป็นตัวแปรอิสระและ X เป็นตัวแปรตาม¹ ถ้าสมการ ผ่านจุดกำเนิดหรือจุดตัดบนแกน X เท่ากับ 0 แล้ว ความสัมพันธ์นี้จะลดรูปเป็น $X_i = \beta_1 Y_i$ และโดยอาศัยวิธี Least Square ทำให้ทราบว่า

$$\beta_1 = \frac{\sum_{i=1}^N (x_i - \bar{X})(y_i - \bar{Y})}{\sum_{i=1}^N (y_i - \bar{Y})^2}$$

- ¹ ก. โดยปกติเรานิยมเสนอสมการถดถอยในรูป $Y_i = \beta_0 + \beta_1 X_i$ ในที่นี้ใช้เป็น $X_i = \beta_0 + \beta_1 Y_i$ เพื่อให้สอดคล้องกับสูตร $V(\hat{R})$ ที่เสนอมานี้แล้ว ความจริงเราจะนิยมให้ใครเป็นตัวแปรอิสระหรือตัวแปรตามก็ได้
- ข. อัตราส่วน R เป็นเรื่องของการเปรียบเทียบกันระหว่างตัวแปร ความสัมพันธ์ระหว่างตัวแปร เช่น ความถดถอยและสหสัมพันธ์ย่อมมีส่วนสำคัญในการอธิบายความหมายและช่วยทำให้เข้าใจการนำค่า R ไปใช้ประโยชน์ได้กระจ่างขึ้น

อนึ่ง จากสมการ $X_i = \beta_1 Y_i$ ถ้ารวมตลอดถึง N ค่าแล้วหารด้วย N (ใส่ \sum_i^N ทั้งสองด้านแล้วหารตลอดด้วย N) จะพบว่า

$$\begin{aligned} \frac{1}{N} \sum_i^N X_i &= \frac{\beta_1}{N} \sum_i^N Y_i \\ \Rightarrow \beta_1 &= \bar{X}/\bar{Y} \\ \text{แต่ } \frac{\bar{X}}{\bar{Y}} &= R \text{ ดังนั้น } \beta_1 = R \text{ นั่นก็คือ } RS_y = \beta_1 S_y \end{aligned}$$

พิจารณา RS_y จะพบว่า

$$\begin{aligned} RS_y &= \beta_1 S_y = \frac{\sum_i^N (x_i - \bar{X})(y_i - \bar{Y})}{\sum_i^N (y_i - \bar{Y})^2} \cdot S_y \\ &= \frac{\text{Cov}(x,y)}{S_y^2} \cdot S_y = \frac{\text{Cov}(x,y)}{S_y} \\ &= \frac{\text{Cov}(x,y)}{S_x S_y} \cdot S_x \\ &= \rho S_x \end{aligned}$$

แสดงว่า ถ้า $\beta_0 = 0$ หรือสมการถดถอยผ่านจุดกำเนิดแล้วจะมีผลให้ $RS_y = \rho S_x$ หรือนัยหนึ่ง ถ้าสมการความสัมพันธ์ระหว่างตัวแปร x และ y พุ่งผ่านจุดกำเนิดมาแล้ว \hat{R} จะเป็นตัวประมาณค่าที่ปราศจากอคติของ R

ดังนั้น ในขั้นต้นเราควรนำค่าสังเกต (x,y) มาพล็อตกราฟดูแนวของจุดว่าเป็นแนวพุ่งผ่านจุดกำเนิดหรือไม่ ถ้าผ่านก็แสดงว่า ค่า $\hat{R} = \frac{\bar{x}}{\bar{y}}$ จะเป็นตัวประมาณค่าที่ปราศจากอคติของ R

อนึ่งเราสามารถใช้ประโยชน์จากข้อสังเกตเหล่านี้ในการคำนวณ $v(\hat{R})$ ได้อีกด้วย เพราะเราสามารถเสนอสูตรของ $v(\hat{R})$ ได้อีกรูปแบบหนึ่งดังนี้

$$\begin{aligned} \text{จาก } V(\hat{R}) &= \frac{1}{\bar{Y}^2} \frac{N-n}{N} \frac{S_x^2}{n} \\ &= \frac{1}{\bar{Y}^2} \cdot \frac{N-n}{Nn} \cdot \frac{1}{N-1} \sum_i^N (x_i - R y_i)^2 \end{aligned}$$

พิจารณาเทอม $\sum_i^N (x_i - R y_i)^2$ จะพบว่า

$$\begin{aligned} \sum_i^N (x_i - R y_i)^2 &= \sum_i^N (x_i - \bar{X} + \bar{X} - R y_i)^2 \\ &= \sum_i^N (x_i - \bar{X} + \bar{X} \cdot \frac{\bar{Y}}{\bar{Y}} - R y_i)^2 \\ &= \sum_i^N (x_i - \bar{X} + R \bar{Y} - R y_i)^2 \\ &= \sum_i^N \{(x_i - \bar{X}) - R(y_i - \bar{Y})\}^2 \\ &= \sum_i^N (x_i - \bar{X})^2 - 2R \sum_i^N (x_i - \bar{X})(y_i - \bar{Y}) \\ &\quad + R^2 \sum_i^N (y_i - \bar{Y})^2 \\ &= (N-1)S_x^2 - 2R(N-1)\rho S_x S_y + (N-1)R^2 S_y^2 \end{aligned}$$

$$\rho = \frac{E(x_i - \bar{X})(y_i - \bar{Y})}{\sqrt{E(x_i - \bar{X})^2} \sqrt{E(y_i - \bar{Y})^2}} = \frac{\sum_i^N (x_i - \bar{X})(y_i - \bar{Y})}{(N-1)S_x S_y}$$

ดังนั้น $\sum_i^N (x_i - \bar{X})(y_i - \bar{Y}) = (N-1)\rho S_x S_y$

$$= (N-1)(S_x^2 + R^2 S_y^2 - 2\rho R S_x S_y)$$

$$\Rightarrow V(\hat{R}) = \frac{1}{\bar{Y}^2} \cdot \frac{N-n}{Nn} \cdot \frac{1}{N-1} \{(N-1)(S_x^2 + R^2 S_y^2 - 2\rho R S_x S_y)\}$$

$$\text{นั่นคือ } V(\hat{R}) = \frac{1}{\bar{Y}^2} \cdot \frac{N-n}{Nn} \cdot (S_x^2 + R^2 S_y^2 - 2\rho R S_x S_y)$$

$$\text{พิจารณาสูตร } V(\hat{R}) = \frac{1}{\bar{Y}^2} \cdot \frac{N-n}{Nn} (S_x^2 + R^2 S_y^2 - 2R\rho S_x S_y)$$

จะพบว่าเนื่องจาก $-1 \leq \rho \leq 1$ ถ้า $\rho = 1$ แล้ว $S_x^2 + R^2 S_y^2 - 2R\rho S_x S_y$ จะมีค่าต่ำที่สุด ดังนั้นเราจึงได้ข้อสรุปอันเป็นเงื่อนไขของ \hat{R} ดังนี้คือ

1. ถ้าสมการ ความสัมพันธ์ระหว่างตัวแปร x และ y เป็นเส้นตรงที่พุ่งผ่านจุดกำเนิด แล้ว \hat{R} จะเป็นตัวประมาณค่าที่ปราศจากอคติของ R
2. ถ้า x และ y มีค่าสหสัมพันธ์เป็น 1 (perfect positive correlation) ความแปรปรวนของการประมาณค่าจะมีค่าต่ำที่สุด
3. ปริมาณความเียงจะลดลงเมื่อใช้กลุ่มตัวอย่างขนาดใหญ่
4. ตัวประมาณค่าของ $V(\hat{R})$ ในกรณีของสูตรรูปนี้คือ

$$\hat{V}(\hat{R}) = \frac{1}{\bar{y}^2} \cdot \frac{N-n}{Nn} (s_x^2 + \hat{R}^2 s_y^2 - 2\hat{R} r s_x s_y)$$

$$\text{โดยที่ } s_x^2 = \frac{1}{n-1} \sum_i (x_i - \bar{X})^2, s_y^2 = \frac{1}{n-1} \sum_i (y_i - \bar{Y})^2, \hat{R} = \bar{x}/\bar{y}$$

$$r = s_{xy}/s_x s_y \text{ ค่าประมาณของสหสัมพันธ์ระหว่างตัวแปร } x \text{ และ } y$$

$$s_{xy} = \frac{1}{n-1} \sum_i (x_i - \bar{x})(y_i - \bar{y})$$

พิสูจน์ $E(\bar{x} - \bar{X})(\bar{y} - \bar{Y}) = \frac{N-n}{Nn} \rho S_x S_y$?

ให้ $u_i = x_i + y_i$ ดังนั้น $\bar{u} = \bar{x} + \bar{y}$ และ $\bar{U} = \bar{X} + \bar{Y}$

$$\begin{aligned} \text{ดังนั้น } E(\bar{u} - \bar{U})^2 &= E\{(\bar{x} + \bar{y}) - (\bar{X} + \bar{Y})\}^2 = E\{(\bar{x} - \bar{X}) + (\bar{y} - \bar{Y})\}^2 \\ &= E(\bar{x} - \bar{X})^2 + E(\bar{y} - \bar{Y})^2 - 2E(\bar{x} - \bar{X})(\bar{y} - \bar{Y}) \end{aligned}$$

$$\begin{aligned} \text{นั่นคือ } E(\bar{x} - \bar{X})(\bar{y} - \bar{Y}) &= \frac{1}{2} \{E(\bar{u} - \bar{U})^2 - E(\bar{x} - \bar{X})^2 - E(\bar{y} - \bar{Y})^2\} \\ &= \frac{1}{2} V(\bar{u}) - \frac{1}{2} V(\bar{x}) - \frac{1}{2} V(\bar{y}) \\ &= \frac{1}{2} \left(\frac{N-n}{N} \cdot \frac{S_u^2}{n} - \frac{N-n}{N} \cdot \frac{S_x^2}{n} - \frac{N-n}{N} \cdot \frac{S_y^2}{n} \right) \\ &= \frac{N-n}{2Nn} (S_u^2 - S_x^2 - S_y^2) \end{aligned}$$

$$\begin{aligned} \text{โดยที่ } S_u^2 &= \frac{1}{N-1} \sum_i^N (u_i - \bar{U})^2 \\ &= \frac{1}{N-1} \sum_i^N \{(x_i + y_i) - (\bar{X} + \bar{Y})\}^2 \\ &= \frac{1}{N-1} \sum_i^N \{(x_i - \bar{X}) + (y_i - \bar{Y})\}^2 \\ &= \frac{1}{N-1} \sum_i^N (x_i - \bar{X})^2 + \frac{1}{N-1} \sum_i^N (y_i - \bar{Y})^2 + \frac{2}{N-1} \sum_i^N (x_i - \bar{X})(y_i - \bar{Y}) \\ &= S_x^2 + S_y^2 + \frac{2}{N-1} \sum_i^N (x_i - \bar{X})(y_i - \bar{Y}) \end{aligned}$$

$$\begin{aligned} \text{ดังนั้น } E(x-\bar{X})(y-\bar{Y}) &= \frac{N-n}{2Nn} \left\{ (S_x^2 + S_y^2 + \frac{2}{N-1} \sum_i^N (x_i - \bar{X})(y_i - \bar{Y})) - S_x^2 - S_y^2 \right\} \\ &= \frac{N-n}{Nn} \cdot \frac{1}{N-1} \sum_i^N (x_i - \bar{X})(y_i - \bar{Y}) \end{aligned}$$

$$\text{แต่ } \therefore e = \frac{\sum_i^N (x_i - \bar{X})(y_i - \bar{Y})}{(N-1)S_x S_y}$$

$$\sum_i^N (x_i - \bar{X})(y_i - \bar{Y}) = (N-1)eS_x S_y$$

$$\begin{aligned} \text{ดังนั้น } E(\bar{x}-\bar{X})(\bar{y}-\bar{Y}) &= \frac{N-n}{Nn} \cdot \frac{1}{N-1} \sum_i^N (x_i - \bar{X})(y_i - \bar{Y}) \\ &= \frac{N-n}{Nn} eS_x S_y \end{aligned}$$

2.4.3 การประมาณค่าสัดส่วนของกลุ่มประชากร (Estimation of Population Proportion or Population Percentage, p)

ในบางครั้งเรามีความสนใจที่จะทราบค่าประมาณของสัดส่วนหรือร้อยละของคุณลักษณะทางประชากรที่สนใจ หรือสนใจจะทราบค่าประมาณของยอดรวมของจำนวนหน่วยสำรวจที่มีลักษณะสอดคล้องกับคุณลักษณะทางประชากรที่เราสนใจ เช่น ร้อยละของผู้เป็นโรคผิวหนัง หรือยอดรวมของผู้ที่เป็นโรคผิวหนัง ร้อยละของนักเรียนสายตาสั้น หรือยอดรวมของนักเรียนที่สายตาสั้น ฯลฯ ในกรณีเช่นนี้เรามีวิธีประมาณค่าเรียกว่าการประมาณค่าสัดส่วนของกลุ่มประชากร

ในที่นี้เราจะเริ่มศึกษาเฉพาะกรณีที่กลุ่มประชากรถูกจำแนกเป็น 2 ฝ่าย (Class) คือมีคุณลักษณะตรงตามความสนใจของเราฝ่ายหนึ่ง เรียกว่าฝ่าย C (Class C) กับไม่มีคุณลักษณะตรงตามความสนใจของเราอีกฝ่ายหนึ่ง เรียกว่า C' (Class C') กล่าวเช่นนี้อาจทำให้มองภาพไม่ออก ขอยกตัวอย่างประกอบเพื่อให้เข้าใจง่ายดังนี้ สมมุติว่าเราสนใจที่จะ

กะประมาณสัดส่วนของนักศึกษาที่เราสนใจในกิจกรรมทางศาสนาและวัฒนธรรม ในที่นี้จะเห็นได้ว่าเราสามารถจำแนกนักศึกษาทั้งหมดออกได้เป็น 2 พวก หรือ 2 ฝ่ายคือ ฝ่ายหนึ่งเป็นฝ่ายที่มีความสนใจในกิจกรรมของศาสนาและวัฒนธรรม อีกฝ่ายหนึ่งไม่สนใจในกิจกรรมดังกล่าว ดังนั้นเป็นต้น เมื่อเป็นเช่นนี้จึงเห็นได้ว่าถ้าดำเนินการสุ่มตัวอย่างหน่วยสำรวจขึ้นมา n หน่วย หน่วยสำรวจหนึ่ง ๆ จึงอาจตกอยู่ในฝ่าย C หรือ C' ก็ได้ หรือนัยหนึ่งหน่วยสำรวจหนึ่ง ๆ อาจมาจากฝ่าย C หรือ C' ก็ได้

ให้ $x = 1$ ถ้า $x \in C$ และ $x = 0$ ถ้า $x \in C'$ ดังนั้น เมื่อสุ่มตัวอย่างหน่วยสำรวจมา n หน่วยจากกลุ่มประชากรขนาด N เราสามารถนิยามสัดส่วนจากกลุ่มตัวอย่างและสัดส่วนของประชากรได้ดังนี้

นิยาม 2.1 เมื่อตัวแปรสุ่ม มีค่าได้เป็น 0 หรือ 1 อย่างใดอย่างหนึ่งเท่านั้นคือ $x = 0$ เมื่อ $x \in C'$ และ $x = 1$ เมื่อ $x \in C$ ดังนั้นสัดส่วนของประชากร P และสัดส่วนจากกลุ่มตัวอย่างจะมีลักษณะดังนี้คือ

$$P = \sum_i^N \frac{x_i}{N} \quad \text{เมื่อ } x=0,1$$

และ

$$p = \sum_i^n \frac{x_i}{n} \quad \text{เมื่อ } x=0,1$$

ข้อสังเกต

1. ขอให้สังเกตว่าทั้ง P และ p มีสูตรโครงสร้างเช่นเดียวกันกับ \bar{X} และ \bar{x} แตกต่างกันเฉพาะในกรณีของสัดส่วนค่าของ x มีได้เพียง 2 ค่าคือ 0 และ 1 เท่านั้นโดยที่ $x=0$ เมื่อ $x \in C'$ และ $x = 1$ เมื่อ $x \in C$ ส่วนในกรณีของค่าเฉลี่ยค่าของ x จะมีได้ไม่จำกัดกล่าวคือ $-\infty < x < \infty$ ดังนั้นเราจึงสามารถสรุปได้ว่าสัดส่วนก็คือกรณีเฉพาะของค่าเฉลี่ยเมื่อ $x=0,1$

2. ในทางปฏิบัติ หมายเลข 0 และ 1 ที่ใช้เป็นค่าของ x นั้นก็คือรหัสที่กำหนดให้แก่ x นั้นเอง เรายินยอมกำหนดรหัสเช่นนี้ให้แก่ตัวแปรเชิงคุณภาพ เช่น เพศ สถานภาพสมรส

ความสัมพันธ์กับหัวหน้าครอบครัว ฯลฯ แต่จะกำหนดให้แก่ตัวแปรเชิงปริมาณด้วยก็ได้ ถ้าสามารถแจกประชากรออกเป็น 2 กลุ่มที่มีความหมายได้

3. เหตุที่ตัวแปร x มีค่าได้ 2 ค่าคือ 0 และ 1 แสดงว่าฟังก์ชันการแจกแจงของตัวแปรสุ่มคือการแจกแจงแบบเบอร์นูลลี (Bernoulli Distribution) และตัวสถิติ $y = \sum_i^N x_i$ จะมีการแจกแจงแบบทวินาม (Binomial Distribution).

4. $P = \sum_i^N \frac{x_i}{N}$ เมื่อ $x = 0$ หรือ 1 แสดงให้เห็นว่า ถ้า $x_i = 0$ แล้ว $x_i^2 = 0$ และถ้า $x_i = 1$ แล้ว $x_i^2 = 1$ ดังนั้นเราสามารถสรุปผลทางคณิตศาสตร์ที่น่าสนใจจากความเป็นจริงต่อไปนี้

$$(1) \sum_i^N x_i^2 = \sum_i^N x_i$$

$$(2) P = \bar{X} = \text{สัดส่วนของประชากรที่ตกอยู่ใน Class C}$$

(3) $\sum_i^N x_i$ คือ ยอดรวมของประชากรที่ตกอยู่ใน Class C ถ้าให้ $\sum_i^N x_i = A$ จะสามารถเขียนสูตรใหม่ได้ดังนี้คือ

$$P = A/N$$

$$\text{นั่นคือ } A = NP \text{ หรือ } \sum_i^N x_i = NP \text{ หรือ } \sum_i^N x_i = NP = N\bar{X}$$

ในทำนองเดียวกัน จาก $p = \sum_i^n \frac{x_i}{n}$ ดังนั้น

$$(1) \sum_i^n x_i^2 = \sum_i^n x_i$$

$$(2) p = \bar{x} = \text{ค่าประมาณของสัดส่วนประชากรที่ตกอยู่ใน Class C}$$

(3) ถ้าให้ $\sum_i^n x_i = a$ แสดงว่า $p = a/n$ นั่นคือ $a = np = \text{ยอดรวมของหน่วยประชากรจาก Class C ที่ตกอยู่ในกลุ่มตัวอย่างขนาด } n$

ทฤษฎี 2.5 ถ้าสุ่มตัวอย่างหน่วยสำรวจมา n ชุดจากกลุ่มประชากรขนาด N โดยใช้แผนสำรวจแบบ SRS แล้ว

$$1. \hat{P} = p = \frac{\sum_i^n x_i}{n} = \text{จะใช้เป็นตัวประมาณค่าของ } P = \frac{\sum_i^N x_i}{N}$$

2. ความแปรปรวนของค่าประมาณ \hat{P} คือ

$$V(\hat{P}) = \frac{N-n}{N-1} \cdot \frac{PQ}{n}$$

3. ตัวประมาณค่าที่ปราศจากความอคติของ $V(\hat{P})$ คือ

$$\hat{V}(\hat{P}) = \frac{N-n}{N} \cdot \frac{pq}{n-1}$$

พิสูจน์

$$1. E(\hat{P}) \stackrel{?}{=} P$$

$$E(\hat{P}) = E\left(\frac{1}{n} \sum_i^n x_i\right)$$

$$= \frac{1}{n} \sum_i^n E(x_i)$$

$$= \frac{1}{n} \sum_i^n \left(\sum_i^N x_i \cdot \frac{1}{N} \right)$$

$$= \frac{1}{n} n \cdot \left(\frac{1}{N} \sum_i^N x_i \right) = P$$

แสดงว่า $\hat{P} = p = \frac{\sum_i^n x_i}{n}$ เป็นตัวประมาณค่าที่ปราศจากอคติของ P และเราสามารถ
ใช้ค่า \hat{p} เป็นค่าประมาณของ P ได้

2. เนื่องจาก $\hat{P} = p$ เป็นกรณีเฉพาะของ \bar{x} โดย $x_i = 0, 1$ ดังนั้นโครงสร้างของสูตรความแปรปรวนของค่าประมาณ $V(\hat{P})$ จึงเป็นโครงสร้างเดียวกัน

ดังนั้น

$$V(\hat{P}) = \frac{N-n}{N} \frac{S^2}{n}$$

พิจารณา S^2 สำหรับกรณีของสัดส่วนจะพบว่า

$$\begin{aligned} S^2 &= \frac{1}{N-1} \sum_i^N (x_i - \bar{X})^2 \\ &= \frac{1}{N-1} \sum_i^N (x_i^2 - 2x_i\bar{X} + \bar{X}^2) \\ &= \frac{1}{N-1} \left(\sum_i^N x_i^2 - N\bar{X}^2 \right) \end{aligned}$$

แต่ $\sum_i^N x_i^2 = \sum_i^N x_i = NP$ และ $\bar{X}^2 = P^2$ (ดูข้อสังเกตท้ายนิยาม 2.1)

$$\begin{aligned} \Rightarrow S^2 &= \frac{1}{N-1} (NP - NP^2) \\ &= \frac{1}{N-1} \cdot NP(1-P) \\ &= \frac{NPQ}{N-1} ; Q = 1-P \end{aligned}$$

$$\text{ดังนั้น } V(\hat{P}) = \frac{N-n}{N} \frac{S^2}{n} = \frac{N-n}{N} \frac{1}{n} \cdot \frac{NPQ}{N-1}$$

$$\Rightarrow V(\hat{P}) = \frac{N-n}{N-1} \frac{PQ}{n}$$

แต่เนื่องจาก P และ Q เป็นพารามิเตอร์ซึ่งเป็นตัวที่ไม่ทราบค่า ดังนั้นในทางปฏิบัติจึงไม่อาจใช้สูตรดังกล่าวนี้ในการคำนวณหาความแปรปรวนของค่าประมาณได้แต่ให้ใช้สูตรต่อไปนี้แทนคือ

$$\hat{V}(\hat{P}) = \frac{N-n}{N} \frac{pq}{n-1}$$

3. จาก $V(\hat{P}) = \frac{N-n}{N} \cdot \frac{S^2}{n}$

แต่เนื่องจาก $E(s^2) = S^2$ หรือ s^2 เป็นตัวประมาณค่าที่ปราศจากอคติของ S^2 ดังนั้นในทางปฏิบัติจึงใส่ s^2 ลงไปแทนที่ของ S^2 ทำให้ได้ค่าประมาณของ $V(\hat{P})$ เป็น

$$\hat{V}(\hat{P}) = \frac{N-n}{N} \frac{s^2}{n}$$

$$\begin{aligned} \text{แต่ } s^2 &= \frac{1}{n-1} \sum_i^n (x_i - \bar{x})^2 \\ &= \frac{1}{n-1} \left(\sum_i^n x_i^2 - n\bar{x}^2 \right) \\ &= \frac{1}{n-1} \left(\sum_i^n x_i - np^2 \right) \\ &= \frac{1}{n-1} (np - np^2) \\ &= \frac{np(1-p)}{n-1} \\ &= \frac{npq}{n-1} \end{aligned}$$

$$\text{ดังนั้น } \hat{V}(\hat{P}) = \frac{N-n}{N} \frac{s^2}{n} = \frac{N-n}{N} \cdot \frac{1}{n} \cdot \frac{npq}{n-1}$$

$$\Rightarrow \hat{V}(\hat{P}) = \frac{N-n}{N} \frac{pq}{n-1}$$

บทแทรก 2.3 เมื่อสุ่มตัวอย่างขนาด n มาจากกลุ่มประชากร N โดยแผนสำรวจแบบ SRS เพื่อประมาณสัดส่วน

1. $\hat{T} = N\hat{P} = Np$ เป็นค่าประมาณที่ปราศจากอคติของยอดรวมของประชากรที่อยู่ใน Class C

$$2. V(\hat{T}) = N^2 V(\hat{P}) = N^2 \frac{N-n}{N-1} \frac{PQ}{n}$$

$$3. \hat{V}(\hat{T}) = N^2 \hat{V}(\hat{P}) = N^2 \frac{N-n}{N-1} \frac{pq}{n-1}$$

4. ช่วงเชื่อมั่น $(1 - \alpha)$ 100% ที่คาดว่าค่าจริงของ P จะปรากฏอยู่คือ

$$\hat{P} - Z_{\alpha/2} \sqrt{\frac{N-n}{N} \frac{pq}{n-1}} < P < \hat{P} + Z_{1-\alpha/2} \sqrt{\frac{N-n}{N} \frac{pq}{n-1}}$$

5. ช่วงเชื่อมั่น $(1 - \alpha)$ 100% ที่คาดว่าค่าจริงของ T จะปรากฏอยู่คือ

$$Np + Z_{\alpha/2} \sqrt{\frac{N-n}{N} \frac{pq}{n-1}} < T < Np + Z_{1-\alpha/2} \sqrt{\frac{N-n}{N} \frac{pq}{n-1}}$$

พิสูจน์ ในที่นี้จะแสดงการพิสูจน์เฉพาะข้อ 1 เท่านั้นข้ออื่น ๆ จะขอเว้นไว้เป็นแบบฝึกหัด ซึ่งนักศึกษาสามารถพิสูจน์ได้เองโดยง่าย

$$\begin{aligned} 1. E(\hat{T}) &= E(Np) \\ &= NE(p) \\ &= NP \quad (\because E(p) = E(\hat{P}) = P) \\ &= N \cdot \frac{1}{N} \sum_i^N x_i \\ &= \sum_i^N x_i = T \end{aligned}$$

แสดงว่า $\hat{T} = Np$ เป็นตัวประมาณค่าที่ปราศจากอคติของยอดรวมประชากร

หมายเหตุ

1. ถ้า $n < 30$ ให้ใช้ $t_{n-1, 1-\alpha/2}$ และ $t_{n-1, \alpha/2}$ แทน $Z_{1-\alpha/2}$ และ $Z_{\alpha/2}$ ตามลำดับ
2. $Z_{\alpha/2} = -Z_{1-\alpha/2}$ เช่น กำหนดให้ $\alpha = 5\%$ ดังนั้น $Z_{1-\alpha/2} = Z_{.975} = 1.96$ และ $Z_{\alpha/2} = Z_{.025} = -1.96$ จะเห็นได้ว่า $Z_{.025} = -1.96 = Z_{.975}$

ดังนั้นทางซ้ายของสมการในข้อ 4 และ 5 เราจะใช้ $\hat{P} - Z_{1-\alpha/2} \sqrt{\frac{N-n}{N} \frac{pq}{n-1}}$

และ $Np - Z_{1-\alpha/2} \sqrt{\frac{N-n}{N} \frac{pq}{n-1}}$ ก็ได้ อื่น \hat{P} และ p มีความหมายเดียวกันจะเขียนโดยใช้ตัวใดก็ได้

3. ดังที่ได้กล่าวมาแล้วในตอนต้นว่า สำหรับกรณีของการประมาณค่าสัดส่วนนั้น เราถือว่าตัวแปรสุ่ม x มีค่าเพียง 2 ค่าคือ $x = 1$ เมื่อ $x \in C$ และ $x = 0$ เมื่อ $x \in C'$ ดังนั้นในทางทฤษฎีตัวแปรสุ่ม x จะมีการกระจายแบบเบอร์นูลลีคือ

$$f(x) = p^x(1-p)^{1-x} ; x = 0, 1$$

และเมื่อสุ่มตัวอย่างมา n หน่วยตัวแปรสุ่ม $Y = x_1 + x_2 + \dots + x_n$ จะมีการกระจายแบบทวินามคือ

$$f(y) = \binom{n}{y} p^y q^{n-y} ; y = 0, 1, \dots, n$$

ปัญหาที่พบก็คือ ตัวสถิติ $p = \frac{\sum_i x_i}{n}$ หรือ $\frac{y}{n}$ มีการแจกแจงแบบใดเพราะถ้าไม่ทราบ Sampling Distribution ของ p เราจะไม่สามารถประมาณค่าของ P ด้วยช่วงเชื่อมั่นหรือทดสอบสมมุติฐานเกี่ยวกับพารามิเตอร์ P ได้

¹ อานมนตรี พิริยะกุล ทฤษฎีสถิติ 2 (โรงพิมพ์มหาวิทยาลัยรามคำแหง, กทม. 2520) บทที่ 2

จากการศึกษาเรื่องทฤษฎีการโน้มเข้าสู่เกณฑ์กลาง (Central Limit Theorem)

เราสามารถพิสูจน์ได้ว่า ถ้า $n \rightarrow \infty$ แล้วตัวแปรสุ่ม $p = \frac{\sum_{i=1}^n x_i}{n}$ จะมีการแจกแจงแบบปกติ

โดยประมาณ มีค่าเฉลี่ย $\mu = P$ และความแปรปรวน $\sigma^2 = \frac{N-n}{N-1} \frac{PQ}{n}$

หรือ $p \sim N(P, \frac{N-n}{N-1} \frac{PQ}{n})$

ด้วยความจริงประการดังกล่าวจึงทำให้เราสามารถกะประมาณค่าของพารามิเตอร์ P ด้วยช่วงเชื่อมั่น $(1 - \alpha) 100\%$ ได้ดังบทแทรก 2.3

ปัญหาอีกประการหนึ่งก็คือ ตามปกติถ้า $p \approx .5$ แล้ว Sampling Distribution ของ p จึงจะเป็นแบบปกติได้โดยประมาณ แต่ถ้า $p \rightarrow 0$ หรือ $p \rightarrow 1$ แล้ว Sampling Distribution ของ p จะเบ้ไม่มีลักษณะเป็นแบบปกติได้ ปัญหาจึงขึ้นอยู่กับขนาดตัวอย่าง n ว่าควรจะเป็นเท่าไร จึงจะแก้ปัญหาคงเบ้ได้

เกี่ยวกับปัญหานี้เรามีทางแก้ปัญหของการคำนวณกำหนดขนาดตัวอย่างได้โดยใช้ขนาดตัวอย่างขั้นต่ำตามตารางซึ่งเสนอโดย W.G. Cochran ดังนี้

ค่าของ p	ค่าขนาดตัวอย่าง n ขั้นต่ำ
0.5	30
0.4 หรือ 0.6	50
0.3 หรือ 0.7	80
0.2 หรือ 0.8	200
0.1 หรือ 0.9	600
0.05 หรือ 0.95	1400