

2.2 다음 진술들에 대하여 동의하는지 혹은 그렇지 않은지를 그 이유와 함께 설명하여라.

- (a) $Cov(Y, X)$ 와 $Cor(Y, X)$ 는 $-\infty$ 와 $+\infty$ 사이의 값을 가질 수 있다.
- (b) $Cov(Y, X) = 0$ 또는 $Cor(Y, X) = 0$ 이면, Y 와 X 사이에 아무런 관계가 없다고 결론지을 수 있다.
- (c) Y 대 \hat{Y} 의 산점도에 있는 점들에 적합된 최소 제곱회귀선은 절편항 0과 기울기 1을 가진다.

Solve)

$$(a) Cov(Y, X) = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{n-1} \text{ 이다.}$$

$$-\infty < \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) < \infty \quad (\because -\infty < (y_i - \bar{y}), (x_i - \bar{x}) < \infty) \text{ 이므로,}$$

$Cov(Y, X)$ 는 $-\infty$ 와 $+\infty$ 사이의 값을 가지게 된다.

$Cor(Y, X)$ 또한 $-\infty$ 와 $+\infty$ 사이의 값을 갖는지 확인해 보자.

우선 $Cov(U, V)$, $U = \frac{Y - E(Y)}{\sqrt{Var(Y)}}$, $V = \frac{X - E(X)}{\sqrt{Var(X)}} \sim N(0, 1)$ 에 관한 아래식을 보자.

$$Cov(U, V) = E(UV) - E(U)E(V)$$

$$= E\left(\frac{(Y - E(Y))(X - E(X))}{\sqrt{Var(Y)}\sqrt{Var(X)}}\right) - 0, \quad (\because U, V \sim N(0, 1) \text{ 이므로 } E(U), E(V) = 0)$$

$$= \frac{Cov(Y, X)}{\sqrt{Var(Y)}\sqrt{Var(X)}}, \quad Cov(Y, X) = E(Y - E(Y))(X - E(X))$$

$$= Cor(Y, X)$$

위 식을 통해 $Cov(U, V) = Cor(Y, X)$ 임을 알 수 있다.

U, V 의 분산에 대해 다음이 성립한다.

$$1) Var(U + V) = Var(U) + Var(V) + 2Cov(U, V)$$

$$2) Var(U - V) = Var(U) + Var(V) - 2Cov(U, V)$$

위의 결과 $Cov(U, V) = Cor(Y, X)$ 을 식에 대입해보면 다음과 같다.

$$\begin{aligned} 1) Var(U + V) &= Var(U) + Var(V) + 2Cov(U, V) \\ &= 2 + 2Cov(U, V), \quad (\because U, V \sim N(0, 1) \text{ 이므로 } Var(U), Var(V) = 1) \\ &= 2 + 2Cor(Y, X) \geq 0, \quad (\because Var(U, V) \geq 0, Cov(U, V) = Cor(Y, X)) \\ &\Leftrightarrow Cor(Y, X) \geq -1 \end{aligned}$$

$$\begin{aligned} 2) Var(U - V) &= Var(U) + Var(V) - 2Cov(U, V) \\ &= 2 - 2Cov(U, V), \quad (\because U, V \sim N(0, 1) \text{ 이므로 } Var(U), Var(V) = 1) \end{aligned}$$

따라서 $-1 \leq Cor(Y, X) \leq 1$ 이므로 (a)의 진술에 동의할 수 없다.

$$\therefore -\infty < Cov(Y, X) < \infty, \quad -1 < Cor(Y, X) < 1$$

(b) $Z \sim Unif(0, 2\pi)$ 이고, $X = \sin z$, $Y = \cos z$ 일 때, $X^2 + Y^2 = \sin^2 z + \cos^2 z = 1$ 이다.

이를 통해 X 와 Y 는 서로 독립이 아님을 알 수 있다.

그리고 아래의 식을 보자.

$$E(XY) = \int_0^{2\pi} \frac{1}{2\pi} \sin z \cos z dz = 0$$

$$E(X) = \int_0^{2\pi} \frac{1}{2\pi} \sin z dz = 0$$

$$E(Y) = \int_0^{2\pi} \frac{1}{2\pi} \cos z dz = 0$$

$$\therefore \text{Cov}(Y, X) = E(XY) - E(X)E(Y) = 0$$

$$\text{Cor}(Y, X) = \frac{\text{Cov}(Y, X)}{S_y \cdot S_x} = 0$$

각각 구한 값으로 $\text{Cov}(Y, X)$, $\text{Cor}(Y, X)$ 가 모두 0임을 알 수 있다.

따라서 위의 예를 통해, “ $\text{Cov}(Y, X)$ 또는 $\text{Cor}(Y, X)$ 가 0이면, Y 와 X 사이에 아무런 관계가 없다.”에 대한 반례가 성립하므로 위의 진술에 동의할 수 없다.

$\therefore \text{Cov}(Y, X)$ 또는 $\text{Cor}(Y, X)$ 가 0이면, Y 와 X 사이에 아무런 관계가 없다고 할 수 없다.

(c) Y 와 \hat{Y} 의 관계는 $Y = \beta_0 + \beta_1 \hat{Y} + \epsilon$, ($\hat{Y} = \hat{\beta}_0^* + \hat{\beta}_1^* X$)의 선형모형이라고 가정할 수 있다.

따라서 각 관측개체는 $y_i = \beta_0 + \beta_1 \hat{y}_i + \epsilon_i$ 로 표현이 가능하고, 다음과 같이 바꿔 쓸 수 있다.

$$\epsilon_i = y_i - \beta_0 - \beta_1 \hat{y}_i, (\hat{y}_i = \hat{\beta}_0^* + \hat{\beta}_1^* x_i)$$

위 식의 제곱합을 최소화하는 직선인 최소제곱회귀선을 구하면 된다.

다음은 제곱합 $S(\beta_0, \beta_1)$ 을 최소화 하는 β_0 , β_1 을 추정하는 과정이다.

$$S(\beta_0, \beta_1) = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 \hat{y}_i)^2$$

$$\frac{\partial S(\beta_0, \beta_1)}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 \hat{y}_i) = 0$$

$$\Leftrightarrow \beta_0 = \bar{y} - \beta_1 \bar{\hat{y}} \quad \dots (1)$$

$$\frac{\partial S(\beta_0, \beta_1)}{\partial \beta_1} = -2 \sum_{i=1}^n \hat{y}_i (y_i - \beta_0 - \beta_1 \hat{y}_i) = 0$$

$$\Leftrightarrow \sum_{i=1}^n y_i \hat{y}_i - \beta_0 \sum_{i=1}^n \hat{y}_i - \beta_1 \sum_{i=1}^n \hat{y}_i^2 = 0 \quad \dots (2)$$

(1)을 (2)에 대입하여 풀면, 다음과 같다.

$$\begin{aligned} & \sum_{i=1}^n y_i \hat{y}_i - \beta_0 \sum_{i=1}^n \hat{y}_i - \beta_1 \sum_{i=1}^n \hat{y}_i^2 \\ &= \sum_{i=1}^n y_i \hat{y}_i - (\bar{y} - \beta_1 \bar{\hat{y}}) \sum_{i=1}^n \hat{y}_i - \beta_1 \sum_{i=1}^n \hat{y}_i^2, (\because \beta_0 = \bar{y} - \beta_1 \bar{\hat{y}}) \end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^n y_i \hat{y} - n \bar{y} \bar{\hat{y}} + n \beta_1 \bar{\hat{y}}^2 - \beta_1 \sum_{i=1}^n \hat{y}_i^2 = 0 \\
&\Leftrightarrow \sum_{i=1}^n y_i \hat{y}_i - n \bar{y} \bar{\hat{y}} = \beta_1 (\sum_{i=1}^n \hat{y}_i^2 - n \bar{\hat{y}}^2) \\
&\Leftrightarrow \beta_1 = \frac{\sum_{i=1}^n y_i \hat{y}_i - n \bar{y} \bar{\hat{y}}}{\sum_{i=1}^n \hat{y}_i^2 - n \bar{\hat{y}}^2} = \frac{\sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2} \\
&\Leftrightarrow \beta_1 = \frac{\hat{\beta}_1 \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\hat{\beta}_1^2 \sum_{i=1}^n (x_i - \bar{x})^2}, \hat{y}_i = \hat{\beta}_0^* + \hat{\beta}_1^* x_i, \bar{\hat{y}} = \hat{\beta}_0^* + \hat{\beta}_1^* \bar{x} \\
&\Leftrightarrow \beta_1 = 0, \hat{\beta}_1^* = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \dots (3)
\end{aligned}$$

따라서 β_1 의 최소제곱추정치는 $\hat{\beta}_1 = 1$ 이다.

(3)을 다시 (1)에 대입해보면 다음과 같다.

$$\begin{aligned}
\beta_0 &= \bar{y} - \hat{\beta}_1 \bar{\hat{y}} \\
&= \bar{y} - \bar{\hat{y}}, \hat{\beta}_1 = 1 \\
&= 0, \hat{\beta}_0^* = \bar{y} - \hat{\beta}_1^* \bar{x} \Leftrightarrow \bar{y} = \hat{\beta}_0^* + \hat{\beta}_1^* \bar{x} \\
&\quad, \bar{\hat{y}} = \hat{\beta}_0^* + \hat{\beta}_1^* \bar{x}
\end{aligned}$$

이를 통해 β_0 의 최소제곱추정치는 $\hat{\beta}_0 = 0$ 임을 알 수 있다.

따라서 $\hat{\beta}_1 = 1, \hat{\beta}_0 = 0$ 이므로, “ Y 대 \hat{Y} 의 산점도에 있는 점들에 적합된 최소제곱회귀선은 절편항 0과 기울기 1을 가진다.”라고 할 수 있다.

2.3 표 2.9에 있는 회귀분석 결과를 이용하여 다음 가설들에 대한 검정을 수행하여라 ($\alpha=0.1$).

변수	계수	표준오차	t-검정	p-값
상수	4.162	3.355	1.24	0.2385
Units	15.509	0.505	30.71	<0.0001

(a) $H_0 : \beta_1 = 15$ 대 $H_0 : \beta_1 \neq 15$

(b) $H_0 : \beta_1 = 15$ 대 $H_0 : \beta_1 > 15$

(c) $H_0 : \beta_0 = 0$ 대 $H_0 : \beta_0 \neq 0$

(d) $H_0 : \beta_0 = 5$ 대 $H_0 : \beta_0 \neq 5$

Solve)

(a) 1. 가설설정: $H_0 : \beta_1 = 15$ vs $H_0 : \beta_1 \neq 15$

2. 유의수준: $\alpha = 0.1$

3. 검정통계량

$$t_1 = \frac{\hat{\beta}_1 - \beta_1}{s.e(\hat{\beta}_1)} = \frac{15.509 - 15}{0.505} = 1.00792$$

$$(\hat{\beta}_1 = 15.509, \beta_1 = 15, s.e(\hat{\beta}_1) = 0.505)$$

4. 기각역: $t_1 = 1.00792 < 1.78 = t_{\alpha/2}(12)$

5. 유의수준 $\alpha = 0.1$ 하에서 검정통계량 t_1 이 임계값보다 작으므로 귀무가설을 기각할 수 없다.

$\Leftrightarrow H_0$ 채택

(b) 1. 가설설정: $H_0 : \beta_1 = 15$ vs $H_0 : \beta_1 > 15$

2. 유의수준: $\alpha = 0.1$

3. 검정통계량

$$t_1 = \frac{\hat{\beta}_1 - \beta_1}{s.e(\hat{\beta}_1)} = \frac{15.509 - 15}{0.505} = 1.00792$$

$$(\hat{\beta}_1 = 15.509, \beta_1 = 15, s.e(\hat{\beta}_1) = 0.505)$$

4. 기각역: $t_1 = 1.00792 < 1.36 = t_{\alpha}(12)$

5. 유의수준 $\alpha = 0.1$ 하에서 검정통계량 t_1 이 임계값보다 작으므로 귀무가설을 기각할 수 없다.

$\Leftrightarrow H_0$ 채택

(c) 1. 가설설정: $H_0 : \beta_0 = 0$ vs $H_0 : \beta_0 \neq 0$

2. 유의수준: $\alpha = 0.1$

3. 검정통계량

$$t_0 = \frac{\hat{\beta}_0 - \beta_0}{s.e.(\hat{\beta}_0)} = \frac{4.162 - 0}{3.355} = 1.240536$$

$$(\hat{\beta}_0 = 4.162, \beta_0 = 0, s.e.(\hat{\beta}_0) = 3.355)$$

4. 기각역: $t_0 = 1.240536 < 1.78 = t_{\alpha/2}(12)$

5. 유의수준 $\alpha = 0.1$ 하에서 검정통계량 t_0 이 임계값보다 작으므로 귀무가설을 기각할 수 없다.

$\Leftrightarrow H_0$ 채택

(d) 1. 가설설정: $H_0 : \beta_0 = 5$ vs $H_0 : \beta_0 \neq 5$

2. 유의수준: $\alpha = 0.1$

3. 검정통계량

$$t_0 = \frac{\hat{\beta}_0 - \beta_0}{s.e.(\hat{\beta}_0)} = \frac{4.162 - 5}{3.355} = -0.249776$$

$$(\hat{\beta}_0 = 4.162, \beta_0 = 5, s.e.(\hat{\beta}_0) = 3.355)$$

4. 기각역: $t_0 = -0.249776 < 1.78 = t_{\alpha/2}(12)$

5. 유의수준 $\alpha = 0.1$ 하에서 검정통계량 t_0 이 임계값보다 작으므로 귀무가설을 기각할 수 없다.

$\Leftrightarrow H_0$ 채택

2.4 표 2.9에 있는 회귀분석 결과를 이용하여 β_0 에 대한 99% 신뢰구간을 구축하여라.

변수	계수	표준오차	t-검정	p-값
상수	4.162	3.355	1.24	0.2385
Units	15.509	0.505	30.71	<0.0001

Solve)

$$\beta_0 \text{의 신뢰구간 } (\alpha = 0.01) : \hat{\beta}_0 \pm t_{\alpha/2}(12) \cdot s.e.(\hat{\beta}_0)$$

$$= 4.162 \pm (3.06 \cdot 3.355), \hat{\beta}_0 = 4.162, t_{\alpha/2}(12) = 3.06, s.e.(\hat{\beta}_0) = 3.355$$

$$\Leftrightarrow (-6.1043, 14.4283)$$

$\therefore \beta_0$ 의 신뢰구간 ($\alpha = 0.01$) : $(-6.1043, 14.4283)$

2.6 표 2.5의 데이터와 표 2.7의 적합값 및 잔차를 이용하여 다음을 보여라.

[표 2.5] 수리시간(Minutes)과 수리될 부품(Units)의 수

번호	Minutes	Units	번호	Minutes	Units
1	23	1	8	97	6
2	29	2	9	109	7
3	49	3	10	119	8
4	64	4	11	149	9
5	74	4	12	145	9
6	87	5	13	154	10
7	96	6	14	166	10

[표 2.7] 컴퓨터 수리시간 데이터에 대한 적합값 \hat{y}_i 와 보통의 최소제곱잔차 e_i

i	x_i	y_i	\hat{y}_i	e_i	i	x_i	y_i	\hat{y}_i	e_i
1	1	23	19.67	3.33	8	6	97	97.21	-0.21
2	2	29	35.18	-6.18	9	7	109	112.72	-3.72
3	3	49	50.69	-1.69	10	8	119	128.23	-9.23
4	4	64	66.20	-2.20	11	9	149	143.74	5.26
5	4	74	66.20	7.80	12	9	145	143.74	1.26
6	5	87	81.71	5.29	13	10	154	159.25	-5.25
7	6	96	97.21	-1.21	14	10	166	159.25	6.75

(a) $Cor(Y, X) = Cor(Y, \hat{Y}) = 0.994$

(b) $SST = 27768.348$

(c) $SSE = 348.848$

Solve)

$$(a) Cor(Y, X) = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2} \sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} \text{ 이다.}$$

위의 식에 표 2.5의 데이터를 대입하면 다음과 같다.

$$\begin{aligned} & \frac{\sum_{i=1}^{14} (y_i - \bar{y})(x_i - \bar{x})}{\sqrt{\sum_{i=1}^{14} (y_i - \bar{y})^2} \sqrt{\sum_{i=1}^{14} (x_i - \bar{x})^2}}, \bar{y} = 97.21, \bar{x} = 6 \\ &= \frac{(-74.21) \cdot (-5) + (-68.21) \cdot (-4) + \dots + 68.79 \cdot 4}{\sqrt{(-74.21)^2 + (-68.21)^2 + \dots + (68.79)^2} \sqrt{((-5)^2 + (-4)^2 + \dots + 4^2)}} \\ &= 0.993698 \approx 0.994 \end{aligned}$$

따라서 $Cor(Y, X) = 0.994$ 을 구할 수 있다.

$Cor(Y, \hat{Y})$ 또한 $Cor(Y, X)$ 와 같은 값을 갖는지 확인해보자.

$$Cor(Y, \hat{Y}) = \frac{\sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2} \sqrt{\sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2}} \text{ 이다.}$$

위의 식에 표 2.7의 데이터를 대입하면 다음과 같다.

$$\begin{aligned} & \frac{\sum_{i=1}^{14} (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sqrt{\sum_{i=1}^{14} (y_i - \bar{y})^2} \sqrt{\sum_{i=1}^{14} (\hat{y}_i - \bar{\hat{y}})^2}}, \bar{y} = 97.21, \bar{\hat{y}} = 97.21 \\ &= \frac{(-74.21) \cdot (-77.54) + (-68.21) \cdot (-62.03) + \dots + 68.79 \cdot 62.03}{\sqrt{(-74.21)^2 + (-68.21)^2 + \dots + (68.79)^2} \sqrt{((-77.54)^2 + (-62.03)^2 + \dots + 62.03^2)}} \\ &= 0.993710 \approx 0.994 \end{aligned}$$

따라서 $Cor(Y, \hat{Y})$ 또한 0.994를 구할 수 있다.

위의 두 결과를 통해 $Cor(Y, X) = Cor(Y, \hat{Y}) = 0.994$ 임을 보일 수 있다.

(b) $SST = \sum_{i=1}^n (y_i - \bar{y})^2$ 이다.

위의 식에 표 2.5의 데이터를 대입하면 다음과 같다.

$$\begin{aligned} & \sum_{i=1}^{14} (y_i - \bar{y})^2, \bar{y} = 97.21 \\ &= (23 - 97.21)^2 + (29 - 97.21)^2 + \dots + (166 - 97.21)^2 \\ &= (-74.21)^2 + (-68.21)^2 + \dots + (68.78)^2 \\ &= 27768.36 \approx 27768.348 \end{aligned}$$

따라서 $SST = 27768.348$ 임을 보일 수 있다.

(c) $SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n e_i^2$ 이다.

위의 식에 표 2.7의 데이터를 대입하면 다음과 같다.

$$\begin{aligned} & \sum_{i=1}^{14} (y_i - \hat{y}_i)^2 = \sum_{i=1}^{14} e_i^2 \\ &= (3.33)^2 + (-6.18)^2 + \dots + (6.75)^2 \\ &= 348.7212 \approx 348.848 \end{aligned}$$

따라서 $SSE = 348.848$ 임을 보일 수 있다.

2.8 최소제곱법을 이용하여 단순선형회귀모형 $Y = \beta_0 + \beta_1 X + \epsilon$ 을 데이터에 적합할 때,

$H_0 : \beta_1 = 0$ 가 기각되지 않는다고 가정하자. 이것은 모형을 $Y = \beta_0 + \epsilon$ 과 같이 단순하게 쓸 수 있음을 의미한다. β_0 의 최소제곱추정치는 $\hat{\beta}_0 = \bar{y}$ 이다(이를 증명할 수 있는가?).

(*) β_0 의 최소제곱추정치는 $\hat{\beta}_0 = \bar{y}$ 이다(이를 증명할 수 있는가?).

(a) 이 경우 최소제곱잔차는 무엇인가?

(b) 최소제곱잔차의 합계가 0임을 보여라.

Solve)

(*) 모형 $Y = \beta_0 + \epsilon$ 에서 각 관측개체는 $y_i = \beta_0 + \epsilon_i$ 로 표현될 수 있다.

β_0 의 최소제곱추정치는, $\sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \beta_0)^2$ 을 최소로 하는 $\hat{\beta}_0$ 임을 알 수 있다.

위의 식을 β_0 에 대해 미분하면 다음과 같다.

$$\frac{d}{d\beta_0} \sum_{i=1}^n (y_i - \beta_0)^2 = -2 \sum_{i=1}^n (y_i - \beta_0) = 0$$

$$\Leftrightarrow \sum_{i=1}^n y_i = \sum_{i=1}^n \beta_0$$

$$\Leftrightarrow n\bar{y} = n\beta_0$$

$$\therefore \bar{y} = \beta_0$$

이를 통해 최소제곱추정치 $\hat{\beta}_0 = \bar{y}$ 임을 알 수 있다.

(a) 최소제곱잔차는 $e_i = y_i - \hat{y}_i$ 으로, 최소제곱법으로 구한 적합값 \hat{y}_i 을 통해 구할 수 있다.

β_0 의 최소제곱추정치는 $\hat{\beta}_0 = \bar{y}$ 임을 알고 있으므로, $\hat{y}_i = \hat{\beta}_0 = \bar{y}$ 임을 쉽게 구할 수 있다.

따라서 최소제곱잔차 $e_i = y_i - \hat{y}_i = y_i - \bar{y}$ 이다.

$$\therefore e_i = y_i - \bar{y} \quad (e_i : \text{최소제곱잔차})$$

(b) 위의 (a)의 풀이를 통해 $\sum_{i=1}^n e_i = \sum_{i=1}^n (y_i - \bar{y})$ 라는 사실을 알 수 있다.

식을 전개하면 다음과 같다.

$$\sum_{i=1}^n e_i = \sum_{i=1}^n (y_i - \bar{y})$$

$$= \sum_{i=1}^n y_i - n\bar{y} = n\bar{y} - n\bar{y} \quad (\because \sum_{i=1}^n y_i = n\bar{y}) = 0$$

따라서 위의 식을 통해 최소제곱잔차 e_i 의 합이 0임을 알 수 있다.

2.13 다음의 y_1, y_2, \dots, y_n 은 알려지지 않은 평균 μ 와 분산 σ^2 으로부터 추출된 표본이다.

평균 μ 를 추정하는 하나의 방법은 다음의 선형모형을 적합시키고 제곱합 $\sum_{i=1}^n (y_i - \mu)^2$ 을 최소화하는 최소제곱방법을 이용하는 것이다.

$$y_i = \mu + \epsilon; \quad i = 1, 2, \dots, n$$

다른 하나의 방법은 수직거리의 합 $\sum_{i=1}^n |y_i - \mu|$ 을 최소화하는 최소절댓값(LAV: least absolute value) 방법을 이용하는 것이다.

- (a) μ 의 최소제곱 추정치는 표본평균 \bar{y} 임을 증명하여라.
- (b) μ 의 최소절댓값 추정치는 표본중위수임을 증명하여라.
- (c) 표본평균의 장점과 단점을 하나씩 설명하여라.
- (d) 표본중위수의 장점과 단점을 하나씩 설명하여라.
- (e) μ 에 대한 위 두 개의 추정치 중 어느 것을 선택할 것인가? 그 이유는 무엇인가?

Solve)

(a) 주어진 선형모형 $y_i = \mu + \epsilon_i$ 는 $\epsilon_i = y_i - \mu$ 로 나타 낼 수 있다.

여기서 ϵ_i 의 제곱합, 즉 $\sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \mu)^2$ 를 최소화하는 μ 가 최소제곱추정치가 된다.

따라서 위의 식을 μ 에 대해 미분하면 다음과 같다.

$$\frac{d}{d\mu} \sum_{i=1}^n (y_i - \mu)^2 = -2 \sum_{i=1}^n (y_i - \mu) = 0$$

$$\Leftrightarrow \sum_{i=1}^n y_i = \sum_{i=1}^n \mu$$

$$\Leftrightarrow n\bar{y} = n\mu$$

$$\therefore \mu = \bar{y}$$

이를 통해 μ 의 최소제곱추정치는 \bar{y} 임을 알 수 있다.

(b) 주어진 선형모형 $y_i = \mu + \epsilon_i$ 는 $\epsilon_i = y_i - \mu$ 로 나타 낼 수 있다.

여기서 $|\epsilon_i|$ 의 합, 즉 $\sum_{i=1}^n |\epsilon_i| = \sum_{i=1}^n |y_i - \mu|$ 를 최소화하는 μ 가 최소절댓값추정치가 된다.

이를 구하기 위해, y_i 를 크기순으로 $y_{(1)}, y_{(2)} \dots y_{(n)}$ 로 다음과 같이 나타내었다.

$$y_{(1)} < y_{(2)} < \cdots < y_{(n-1)} < y_n, \quad (n=1, 2, \cdots, n)$$

$\sum_{i=1}^n |y_i - \mu| = f(\mu)$ 라고 했을 때, $f(\mu)$ 의 최소를 만족하는 μ 는 다음과 같이 구할 수 있다.

$$y_{(k)} \leq \mu \leq y_{(k+1)}, \quad (1 \leq k \leq n-1) \text{ 일 때,}$$

$$\begin{aligned} f(\mu) &= \sum_{i=1}^k (\mu - y_{(i)}) + \sum_{i=k+1}^n (y_{(i)} - \mu) \\ &= k\mu - (n-k)\mu + \sum_{i=1}^k y_{(i)} + \sum_{i=k+1}^n y_{(i)} \\ &= (2k-n)\mu + \sum_{i=1}^k y_{(i)} + \sum_{i=k+1}^n y_{(i)} \end{aligned}$$

1) $2k-n < 0$ 일 경우, 함수 $f(x)$ 는 $[y_{(k)}, y_{(k+1)}]$ 에서 감소한다. , $(k=1, 2, \cdots, \frac{n}{2})$

2) $2k-n > 0$ 일 경우, 함수 $f(x)$ 는 $[y_{(k)}, y_{(k+1)}]$ 에서 증가한다. , $(k=\frac{n}{2}+1, \frac{n}{2}+2, \cdots, n-1)$

1), 2)를 모두 고려했을 때, $\mu = y_{(\frac{n}{2}+1)}$ 일 때, $f(\mu)$ 의 기울기가 0이므로, $f(\mu)$ 가

위의 n 이 짝수인 경우, $y_{(\frac{n}{2})}$ 와 $y_{(\frac{n}{2}+1)}$ 사이의 적당한 실수 하나를 중앙값으로 정한다.

예를 들면, $y_{(k)}$ 와 $y_{(k+1)}$ 의 산술평균을 중앙값으로 할 수 있다.

따라서, μ 의 최소절댓값추정치는 표본 y_i 의 중위수임을 알 수 있다.

(c) 장점: 모든 자료의 값을 이용하여 나타낸다.

단점: 모든 자료의 값을 이용하므로, 이상치의 영향을 받는다.

(d) 장점: 자료의 값의 분포가 치우쳐 있거나, 이상치가 있어도, 그 영향을 덜 받기 때문에 표본평균보다 유용하게 사용할 수 있다.

단점: 모든 자료의 값을 활용하지 못한다.

(e) 표본평균을 사용할 것이다. 이상치가 흔하게 나오지 않는다고 생각하고, 또한 표본자료의 모든 값을 이용하여 잘 나타내므로, 표본중위값보다 대체적으로 μ 의 추정치로 적합하다고 생각했기 때문이다.