

Gauss propuso que se considerase ε una variable aleatoria.

Por otra parte, se ha considerado asumir que

$$E(\varepsilon) = 0 \quad \text{VAR}(\varepsilon) = \sigma^2 \quad \text{y que}$$

$$(RL) \quad \dots \quad Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \quad ; \quad i=1, 2, \dots, N.$$

Si se asume que la variable explicativa X tiene valores que son fijados en forma determinística por quienes estudian el experimento, entonces podemos considerar que $\beta_0 + \beta_1 X_i$ no es un término aleatorio en el modelo y por tanto el carácter aleatorio de ε_i se hereda de forma directa a Y_i , de forma que

$$\begin{aligned} E(Y_i) &= E(\beta_0 + \beta_1 X_i + \varepsilon_i) \\ &= E(\beta_0 + \beta_1 X_i) + E(\varepsilon_i) \\ &= \beta_0 + \beta_1 X_i + 0 = \beta_0 + \beta_1 X_i ; \forall i=1, \dots, N. \end{aligned}$$

También

$$\begin{aligned} \text{VAR}(Y_i) &= \text{VAR}(\beta_0 + \beta_1 X_i + \varepsilon_i) \\ &= \text{VAR}(\varepsilon_i) = \sigma^2 ; \forall i=1, \dots, N \end{aligned}$$

Para simplificar el modelo, también se considera que el error en la medición Y_i no tiene que ver, en un sentido estocástico, con el error para la medición Y_j , para cada $i \neq j$. Aunque la correlación cero no implica independencia de variables aleatorias, por el momento asumiremos que $\text{COV}(\epsilon_i, \epsilon_j) = 0 \quad \forall i \neq j$.

Como consecuencia de este supuesto

$$\begin{aligned}
 \text{COV}(Y_i, Y_j) &= \text{COV}(\beta_0 + \beta_1 X_i + \epsilon_i, \beta_0 + \beta_1 X_j + \epsilon_j) \\
 &= \text{COV}(\beta_0 + \beta_1 X_i, \beta_0 + \beta_1 X_j + \epsilon_j) + \text{COV}(\epsilon_i, \beta_0 + \beta_1 X_j + \epsilon_j) \\
 &= \text{COV}(\epsilon_i, \beta_0 + \beta_1 X_j + \epsilon_j) \\
 &= \text{COV}(\epsilon_i, \epsilon_j) + \text{COV}(\epsilon_i, \beta_0 + \beta_1 X_j) \\
 &= \text{COV}(\epsilon_i, \epsilon_j) = 0 \quad ; \quad \forall i \neq j
 \end{aligned}$$

Usaremos estas propiedades para evaluar la calidad de los coeficientes ajustados por mínimos cuadrados $\hat{\beta}_0$ y $\hat{\beta}_1$. Recordemos que

$$\hat{\beta}_1 = \frac{\sum_{i=1}^N (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^N (x_i - \bar{x})^2} \quad \text{y} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x},$$

Como ya mencionamos antes, los valores de la variable explicativa se consideran como constantes conocidas, de forma que $\sum_{i=1}^N (x_i - \bar{x})^2$ también es una constante conocida para nosotros. Por otra parte $\sum_{i=1}^N (y_i - \bar{y})(x_i - \bar{x})$ involucre valores de la variable de respuesta, dado el modelo (RL), estos valores son realizaciones de las variables aleatorias Y_1, \dots, Y_N y por ende, $\hat{\beta}_1$ también se puede pensar como una variable aleatoria. A continuación analizaremos a $\hat{\beta}_1$ con más detalle, observemos que

$$\begin{aligned} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y}) &= \sum_{i=1}^N (x_i - \bar{x}) y_i - \sum_{i=1}^N (x_i - \bar{x}) \bar{y} \\ &= \sum_{i=1}^N (x_i - \bar{x}) y_i - \bar{y} \sum_{i=1}^N (x_i - \bar{x}) \\ &= \sum_{i=1}^N (x_i - \bar{x}) y_i - \bar{y} \cdot 0 \\ &= \sum_{i=1}^N (x_i - \bar{x}) y_i. \end{aligned}$$

De forma que $\hat{\beta}_1 = \sum_{i=1}^N a_i Y_i$, donde $a_i = \frac{x_i - \bar{x}}{\sum_{i=1}^N (x_i - \bar{x})^2}$; $i=1, 2, \dots, N$. Lo anterior nos dice que $\hat{\beta}_1$ es una función lineal de las variables aleatorias Y_1, \dots, Y_N , como consecuencia de esto se tiene que

$$\begin{aligned}
\mathbb{E}(\hat{\beta}_1) &= \mathbb{E}\left(\sum_{i=1}^N a_i Y_i\right) = \sum_{i=1}^N \mathbb{E}(a_i Y_i) \\
&= \sum_{i=1}^N a_i \mathbb{E}(Y_i) = \sum_{i=1}^N a_i (\beta_0 + \beta_1 X_i) \\
&= \sum_{i=1}^N a_i \beta_0 + \sum_{i=1}^N a_i \beta_1 X_i \\
&= \beta_0 \sum_{i=1}^N a_i + \beta_1 \sum_{i=1}^N a_i X_i.
\end{aligned}$$

Pero $\sum_{i=1}^N a_i = \frac{1}{\sum_{i=1}^N (X_i - \bar{X})^2} \cdot \sum_{i=1}^N (X_i - \bar{X}) = 0$ y además

$$\begin{aligned}
\sum_{i=1}^N a_i X_i &= \frac{1}{\sum_{i=1}^N (X_i - \bar{X})^2} \sum_{i=1}^N (X_i - \bar{X}) X_i \\
&= \frac{1}{\sum_{i=1}^N (X_i - \bar{X})^2} \left\{ \sum_{i=1}^N (X_i - \bar{X})(X_i - \bar{X}) + \sum_{i=1}^N (X_i - \bar{X}) \bar{X} \right\} \\
&= \frac{1}{\sum_{i=1}^N (X_i - \bar{X})^2} \cdot \sum_{i=1}^N (X_i - \bar{X})^2 = 1
\end{aligned}$$

$$\therefore \mathbb{E}(\hat{\beta}_1) = \beta_1$$

" $\hat{\beta}_1$ es un estimador
insesgado y lineal
de β_1 "

Por otra parte, debido a que $\text{COV}(Y_i, Y_j) = 0 \quad \forall i \neq j$

$$\text{VAR}(\hat{\beta}_1) = \text{VAR}\left(\sum_{i=1}^N a_i Y_i\right) = \sum_{i=1}^N a_i^2 \text{VAR}(Y_i)$$

$$\begin{aligned}
 \text{VAR}(\hat{\beta}_1) &= \sum_{i=1}^N \frac{(x_i - \bar{x})^2}{\left(\sum_{i=1}^N (x_i - \bar{x})^2 \right)^2} \sigma^2 \\
 &= \sigma^2 \frac{1}{\left(\sum_{i=1}^N (x_i - \bar{x})^2 \right)^2} \sum_{i=1}^N (x_i - \bar{x})^2 \\
 &= \frac{\sigma^2}{\sum_{i=1}^N (x_i - \bar{x})^2}
 \end{aligned}$$

De este cálculo, observamos que: La precisión de la estimación aumenta en la medida en que se incrementa el valor de $\sum_{i=1}^N (x_i - \bar{x})^2$ ⁽¹⁾. Entonces, cuando es posible elegir los valores de la variable explicativa, en el diseño del experimento, resulta conveniente que estos valores sean tales que su varianza sea lo más grande posible.

El caso de $\hat{\beta}_0$ es similar

$$\begin{aligned}
 \hat{\beta}_0 &= \bar{Y} - \hat{\beta}_1 \bar{X} = \bar{Y} - \left(\sum_{i=1}^N a_i Y_i \right) \bar{X} \\
 &= \sum_{i=1}^N \frac{1}{N} Y_i - \left(\sum_{i=1}^N a_i Y_i \right) \bar{X} = \sum_{i=1}^N \left(\frac{1}{N} - a_i \bar{X} \right) Y_i
 \end{aligned}$$

(1) Si $\sum_{i=1}^N (x_i - \bar{x})^2$ crece, entonces $\text{VAR}(\hat{\beta}_1) = E((\hat{\beta}_1 - \beta_1)^2)$ decrece y en promedio la "distancia" $(\hat{\beta}_1 - \beta_1)^2$ es más pequeña.

$$\hat{\beta}_0 = \sum_{i=1}^N b_i Y_i, \quad \text{donde } b_i = \frac{1}{N} - a_i \bar{X}$$

Por tanto, $\hat{\beta}_0$ es también, un estimador lineal de β_0 .

Como variable aleatoria tenemos que $\hat{\beta}_0$ satisface

$$\begin{aligned} E(\hat{\beta}_0) &= \sum_{i=1}^N b_i E(Y_i) \\ &= \sum_{i=1}^N b_i (\beta_0 + \beta_1 X_i) \\ &= \sum_{i=1}^N b_i \beta_0 + \sum_{i=1}^N b_i \beta_1 X_i \\ &= \beta_0 \sum_{i=1}^N b_i + \beta_1 \sum_{i=1}^N b_i X_i \\ &= \beta_0 \sum_{i=1}^N \left(\frac{1}{N} - a_i \bar{X} \right) + \beta_1 \sum_{i=1}^N \left(\frac{1}{N} - a_i \bar{X} \right) X_i. \end{aligned}$$

Pero como $\sum_{i=1}^N a_i X_i = 1$ entonces

$$\sum_{i=1}^N \left(\frac{1}{N} - a_i \bar{X} \right) X_i = \sum_{i=1}^N \frac{1}{N} X_i - \bar{X} \sum_{i=1}^N a_i X_i = \bar{X} - \bar{X} = 0.$$

Por otra parte $\sum_{i=1}^N a_i = 0$ implica que

$$\sum_{i=1}^N \left(\frac{1}{N} - a_i \bar{X} \right) = 1 - \bar{X} \sum_{i=1}^N a_i = 1. \quad \text{Por tanto}$$

$$E(\hat{\beta}_0) = \beta_0 \times 1 + \beta_1 \times 0 = \beta_0,$$

de manera que $\hat{\beta}_0$ es un estimador lineal e insesgado

de β_0 . Por otra parte (debido a que $\text{cov}(Y_i, Y_j) = 0$
 $\forall i \neq j$)

$$\text{VAR}(\hat{\beta}_0) = \sum_{i=1}^N b_i^2 \text{VAR}(Y_i)$$

$$= \sum_{i=1}^N b_i^2 \sigma^2$$

$$= \sigma^2 \sum_{i=1}^N b_i^2$$

Pero $\sum_{i=1}^N b_i^2 = \sum_{i=1}^N \left(\frac{1}{N} - a_i \bar{X} \right)^2 = \sum_{i=1}^N \left(\frac{1}{N^2} - \frac{2a_i \bar{X}}{N} + a_i^2 \bar{X}^2 \right)$

$$= \frac{1}{N} + \bar{X}^2 \sum_{i=1}^N a_i^2 = \frac{1}{N} + \bar{X}^2 \sum_{i=1}^N \left\{ \frac{(X_i - \bar{X})^2}{\left(\sum_{i=1}^N (X_i - \bar{X})^2 \right)^2} \right\}$$

$$= \frac{1}{N} + \bar{X}^2 \frac{1}{\sum_{i=1}^N (X_i - \bar{X})^2}$$

$$= \frac{\frac{1}{N} \sum_{i=1}^N (X_i - \bar{X})^2 + \bar{X}^2}{\sum_{i=1}^N (X_i - \bar{X})^2} = \frac{\frac{1}{N} \left(\sum_{i=1}^N X_i^2 - N \bar{X}^2 \right) + \bar{X}^2}{\sum_{i=1}^N (X_i - \bar{X})^2}$$

$$= \frac{\frac{1}{N} \sum_{i=1}^N X_i^2 - \bar{X}^2 + \bar{X}^2}{\sum_{i=1}^N (X_i - \bar{X})^2} = \frac{1}{N} \frac{\left(\sum_{i=1}^N X_i^2 \right)}{\sum_{i=1}^N (X_i - \bar{X})^2}$$

$$\therefore \text{VAR}(\hat{\beta}_0) = \sigma^2 \left(\sum_{i=1}^N X_i^2 \right) / N \left(\sum_{i=1}^N (X_i - \bar{X})^2 \right)$$

Esta variancia es más pequeña, a medida que los valores de la variable explicativa producen un