

Capítulo 3

MODELOS DE REGRESIÓN

ESTIMACIÓN

3.1 Estimación puntual de los parámetros del modelo

Existen diversos métodos para hallar estimadores, los más usados son el de máxima de verosimilitud y el de mínimos cuadrados ordinarios. Para poder usar el método de máxima verosimilitud es necesario conocer la distribución de probabilidad de los y_i , los cuales en este caso se suponen normales. En el caso del método de mínimos cuadrados no es necesario conocer la distribución de probabilidad de los y_i . En primer lugar se determinan los estimadores por el método de máxima verosimilitud y luego se hará por el método de mínimos cuadrados.

3.1.1 Estimadores por el método de mínimos cuadrados.

3.1.1.1 Estimador de β En esta sección discutimos el método de los mínimos cuadrados para estimar β . Ningún supuesto sobre la distribución de \mathbf{y} es requerido para obtener los estimadores.

El método de mínimos cuadrados consiste en determinar los estimadores de los parámetros que minimizan la suma de los errores elevados al cuadrado. En este caso, dado que el modelo es $\mathbf{y} = \mathbf{X}\beta + \varepsilon$, el método de mínimos cuadrados consiste en determinar los valores de β , digamos $\hat{\beta}$, para los cuales $\sum_{i=1}^n \varepsilon_i^2$ es mínima. Dado que $\sum_{i=1}^n \varepsilon_i^2 = \varepsilon'\varepsilon$ y $\varepsilon = \mathbf{y} - \mathbf{X}\beta$, entonces se quiere

$$\min \varepsilon' \varepsilon = \min (\mathbf{y} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\beta) \quad (3.1)$$

El resultado está dado en el siguiente teorema

Teorema 3.1 Sea $\mathbf{y} = \mathbf{X}\beta + \varepsilon$ donde \mathbf{X} es una matriz $n \times p$ de rango $p < n$, entonces el valor de β que minimiza $(\mathbf{y} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\beta)$ es

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} = \mathbf{X}^{-}\mathbf{y} \quad (3.2)$$

Prueba. Queremos hallar el valor de β

$$\min \varepsilon' \varepsilon = \min (\mathbf{y} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\beta)$$

Desarrollando el producto $(\mathbf{y} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\beta)$ se tiene que

$$\varepsilon' \varepsilon = \mathbf{y}'\mathbf{y} - 2\mathbf{y}'\mathbf{X}\beta + \beta'\mathbf{X}'\mathbf{X}\beta$$

Ahora para obtener β que minimiza $\varepsilon' \varepsilon$ se deriva $\varepsilon' \varepsilon$ con respecto a β y se iguala a cero obteniéndose de esta manera la siguiente ecuación

$$\frac{\partial \varepsilon' \varepsilon}{\partial \beta} = -2\mathbf{X}'\mathbf{y} + 2\mathbf{X}'\mathbf{X}\hat{\beta} = 0$$

despejando se obtienen las siguientes ecuaciones

$$\mathbf{X}'\mathbf{X}\hat{\beta} = \mathbf{X}'\mathbf{y}$$

la cual es conocida como las **ecuaciones normales**. Ahora, suponiendo que la matriz \mathbf{X} es de rango completo, $\mathbf{X}'\mathbf{X}$ es no singular, se obtiene el estimador de β , el cual está dado por

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} = \mathbf{X}^{-}\mathbf{y}$$

lo cual demuestra el teorema.

Ya que $\hat{\beta}$ en 3.2 minimiza $\varepsilon' \varepsilon$, $\hat{\beta}$ es llamado el *estimador de mínimos cuadrados*. Note que cada $\hat{\beta}_j$ en $\hat{\beta}$ es una función lineal de \mathbf{y} ; esto es, $\hat{\beta}_j = \mathbf{a}'_j \mathbf{y}$ donde \mathbf{a}'_j es la j -ésima fila de $(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'$.

Ahora demostraremos que $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$ minimiza $\varepsilon' \varepsilon$.

$$\begin{aligned} \varepsilon' \varepsilon &= (\mathbf{y} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\beta) \\ &= (\mathbf{y} - \mathbf{X}\hat{\beta} + \mathbf{X}\hat{\beta} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\hat{\beta} + \mathbf{X}\hat{\beta} - \mathbf{X}\beta) \\ &= (\mathbf{y} - \mathbf{X}\hat{\beta})' (\mathbf{y} - \mathbf{X}\hat{\beta}) + (\hat{\beta} - \beta)' \mathbf{X}'\mathbf{X} (\hat{\beta} - \beta) \\ &\quad + 2(\hat{\beta} - \beta)' (\mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{X}\hat{\beta}) \end{aligned} \quad (3.3)$$

El tercer término de lado derecho de 3.3 desaparece ya que de las ecuaciones normales $\mathbf{X}'\mathbf{X}\hat{\beta} = \mathbf{X}'\mathbf{y}$. El segundo término es una forma cuadrática definida positiva (asumiendo que \mathbf{X} es de rango completo), por lo tanto $\varepsilon' \varepsilon$ es minimizada cuando $\beta = \hat{\beta}$.

Actividad: Determinar los estimadores de mínimos cuadrados en el caso del modelo lineal simple.

La estructura de la matriz $\mathbf{X}'\mathbf{X}$ y del vector $\mathbf{X}'\mathbf{y}$ es la siguiente

$$\begin{aligned}\mathbf{X}'\mathbf{X} &= \begin{pmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{pmatrix}' \begin{pmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{pmatrix} \\ &= \begin{pmatrix} n & \sum_i x_{i1} & \sum_i x_{i2} & \dots & \sum_i x_{ik} \\ \sum_i x_{i1} & \sum_i x_{i1}^2 & \sum_i x_{i1}x_{i2} & \dots & \sum_i x_{i1}x_{ik} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sum_i x_{ik} & \sum_i x_{i1}x_{ik} & \sum_i x_{i2}x_{ik} & \dots & \sum_i x_{ik}^2 \end{pmatrix}\end{aligned}$$

$$\begin{aligned}\mathbf{X}'\mathbf{y} &= \begin{pmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_k \end{pmatrix} \\ &= \begin{pmatrix} \sum_i y_i \\ \sum_i x_{i1}y_i \\ \vdots \\ \sum_i x_{ik}y_i \end{pmatrix}\end{aligned}$$

Si $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$, entonces

$$\hat{\boldsymbol{\varepsilon}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{y} - \hat{\mathbf{y}} \quad (3.4)$$

es el vector de residuales, $\hat{\varepsilon}_1 = y_1 - \hat{y}_1$, $\hat{\varepsilon}_2 = y_2 - \hat{y}_2, \dots, \hat{\varepsilon}_n = y_n - \hat{y}_n$. El vector de residuales $\hat{\boldsymbol{\varepsilon}}$ estima a $\boldsymbol{\varepsilon}$ en el modelo $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ y puede ser usado para chequear la validez del modelo y evaluar los supuestos. Esto se verá mas adelante.

■ EJEMPLO 3.1 Estimación de los parámetros en el Modelo de Regresión Lineal

Simple

Considere el modelo lineal simple

$$\begin{aligned}Y_i &= \beta_0 + \beta_1 x_i + \varepsilon_i; \quad i = 1, 2, \dots, n; \quad n \geq 2 \\ \varepsilon_i &\sim NID(0, \sigma^2) \\ \Omega &= \{(\beta_0, \beta_1, \sigma^2) : \beta_0 \in E_1, \beta_1 \in E_1, \sigma^2 > 0\}\end{aligned} \quad (3.5)$$

Cuya forma matricial es

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \quad \mathbf{X} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$$\mathbf{X}'\mathbf{X} = \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{pmatrix} \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} = \begin{pmatrix} n & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{pmatrix}$$

$$\mathbf{X}'\mathbf{y} = \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{pmatrix}$$

$$(\mathbf{X}'\mathbf{X})^{-1} = \frac{1}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \begin{pmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{pmatrix}$$

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \frac{1}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \begin{pmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{pmatrix} \begin{pmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{pmatrix} \\ &= \frac{1}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \begin{pmatrix} \sum_{i=1}^n x_i^2 \sum_{i=1}^n y_i - \sum_{i=1}^n x_i y_i \sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i \sum_{i=1}^n y_i + n \sum_{i=1}^n x_i y_i \end{pmatrix} \\ &= \begin{pmatrix} \bar{y} - \hat{\beta}_1 \bar{x} \\ \frac{\sum_{i=1}^n (x_i - \bar{x}) \sum_{i=1}^n (y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{pmatrix} \end{aligned}$$

Así,

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad \hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x}) \sum_{i=1}^n (y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (3.6)$$

■ EJEMPLO 3.2

Usando los datos de la tabla 2.2 ilustramos los cálculos de $\hat{\beta}$ usando (3.2). Los datos matricialmente son los siguientes

$$\mathbf{y} = \begin{pmatrix} 2 \\ 3 \\ 2 \\ 7 \\ 6 \\ 8 \\ 10 \\ 7 \\ 8 \\ 12 \\ 11 \\ 14 \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & 0 & 2 \\ 1 & 2 & 6 \\ 1 & 2 & 7 \\ 1 & 2 & 5 \\ 1 & 4 & 9 \\ 1 & 4 & 8 \\ 1 & 4 & 7 \\ 1 & 6 & 10 \\ 1 & 6 & 11 \\ 1 & 6 & 9 \\ 1 & 8 & 15 \\ 1 & 8 & 13 \end{pmatrix}$$

entonces,

$$\mathbf{X}'\mathbf{X} = \begin{pmatrix} 12 & 52 & 102 \\ 52 & 395 & 536 \\ 102 & 536 & 1004 \end{pmatrix}, \quad \begin{pmatrix} 90 \\ 482 \\ 872 \end{pmatrix}$$

$$(\mathbf{X}'\mathbf{X})^{-1} = \begin{pmatrix} 0.97476 & 0.24290 & -0.22871 \\ 0.24290 & 0.16207 & -0.11120 \\ -0.22871 & -0.11120 & 0.08360 \end{pmatrix}$$

por lo tanto,

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} = \begin{pmatrix} 5.3754 \\ 3.0118 \\ -1.2855 \end{pmatrix}$$

3.1.1.2 Propiedades del Estimador de Mínimos Cuadrados $\hat{\beta}$ El estimador de mínimos cuadrados $\hat{\beta}$ en el teorema 3.1 fue obtenido sin usar los supuestos de que $E(\mathbf{y}) = \mathbf{X}\beta$ y $cov(\mathbf{y}) = \sigma^2\mathbf{I}$ dados en la sección 2.2.2. Solamente se postuló un modelo $\mathbf{y} = \mathbf{X}\beta + \varepsilon$ y se ajustó. Si $E(\mathbf{y}) \neq \mathbf{X}\beta$, el modelo $\mathbf{y} = \mathbf{X}\beta + \varepsilon$ podría no ajustarse a los datos, en cuyo caso, $\hat{\beta}$ podría tener propiedades pobres. Si $cov(\mathbf{y}) \neq \sigma^2\mathbf{I}$, esto podría tener efectos

adversos adicionales sobre $\hat{\beta}$. Sin embargo, si $E(\mathbf{y}) = \mathbf{X}\beta$ y $cov(\mathbf{y}) = \sigma^2\mathbf{I}$, $\hat{\beta}$ tiene algunas propiedades buenas, como se notará en los cuatro teoremas de esta sección. Note que $\hat{\beta}$ es un vector aleatorio. Discutimos su vector de media y matriz de covarianza (sin el supuesto de la distribución de \mathbf{y}) y su distribución (asumiendo que las variables y son normales) se verá más adelante. Nuevamente en los siguientes teoremas, asumimos que \mathbf{X} es fija y de rango completo.

Teorema 3.2 Si $E(\mathbf{y}) = \mathbf{X}\beta$, entonces $\hat{\beta}$ es un estimador insesgado para β

Prueba.

$$E(\hat{\beta}) = E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}] = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E(\mathbf{Y}) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\beta = \beta$$

Teorema 3.3 Si $cov(\mathbf{y}) = \sigma^2\mathbf{I}$, entonces la matriz de covarianza para $\hat{\beta}$ está dado por $\sigma^2(\mathbf{X}'\mathbf{X})^{-1}$

Prueba.

$$\begin{aligned} cov(\hat{\beta}) &= cov[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}] \\ &= [(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']cov(\mathbf{y})[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']' \end{aligned} \quad (3.7)$$

$$= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\sigma^2\mathbf{I}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \quad (3.8)$$

$$= \sigma^2(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \quad (3.9)$$

$$= \sigma^2(\mathbf{X}'\mathbf{X})^{-1} \quad (3.10)$$

■ EJEMPLO 3.3

Para los datos de la tabla 2.2, $(\mathbf{X}'\mathbf{X})^{-1}$ esta dada en el ejemplo 3.2. Por lo tanto, $cov(\hat{\beta})$ está dada por

$$cov(\hat{\beta}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1} = \sigma^2 \begin{pmatrix} 0.97476 & 0.24290 & -0.22871 \\ 0.24290 & 0.16207 & -0.11120 \\ -0.22871 & -0.11120 & 0.08360 \end{pmatrix}$$

El valor negativo de $cov(\hat{\beta}_1, \hat{\beta}_2) = -0.11120$ indica que en muestreos repetidos (usando los mismos 12 valores de x_1 y x_2), $\hat{\beta}_1$ y $\hat{\beta}_2$ tenderán a moverse en direcciones opuestas; es decir, un incremento en una debería estar acompañada por un decremento en la otra.

Además de que $E(\hat{\beta}) = \beta$ y $cov(\hat{\beta}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}$, una tercera propiedad importante de $\hat{\beta}$ es que bajo los supuestos estandar, la varianza de cada $\hat{\beta}_j$ es mínima. Veamoslo en el siguiente teorema.

Teorema 3.4 (Teorema de Gauss-Markov) Si $E(\mathbf{y}) = \mathbf{X}\beta$ y $cov(\mathbf{y}) = \sigma^2\mathbf{I}$, el estimador de mínimos cuadrados $\hat{\beta}_j$, $j = 0, 1, \dots, k$ tiene varianza mínima entre todos los estimadores lineales insesgados.

Prueba. La prueba se consigue en la página 147 del Rencher (2008).

El teorema de Gauss-Markov es algunas veces establecido de la siguiente manera. Si $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$ y $cov(\mathbf{y}) = \sigma^2\mathbf{I}$, los estimadores de mínimo cuadrado $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$ son los mejores estimadores lineales insesgados (BLUE). En esta expresión, *mejor* significa varianza mínima y *lineal* indica que los estimadores son funciones lineales de \mathbf{y} .

La característica resaltante del teorema de Gauss-Markov es su generalidad distribucional. El resultado se obtiene para cualquier distribución de \mathbf{y} ; la normalidad no es requerida. Los únicos supuestos usados en la prueba son $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$ y $cov(\mathbf{y}) = \sigma^2\mathbf{I}$. Si estos supuestos no se cumplen, $\hat{\boldsymbol{\beta}}$ podría ser sesgado o cada $\hat{\beta}_j$ podría tener varianza más grande que algún otro estimador.

El teorema de Gauss-Markov se puede extender fácilmente a una combinación lineal de los $\hat{\boldsymbol{\beta}}$ como sigue

Corolario 3.1 Si $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$ y $cov(\mathbf{y}) = \sigma^2\mathbf{I}$, el mejor estimador lineal insesgado de $\mathbf{a}'\boldsymbol{\beta}$ es $\mathbf{a}'\hat{\boldsymbol{\beta}}$, donde $\hat{\boldsymbol{\beta}}$ es el estimador de mínimo cuadrados.

3.1.1.3 Estimador para σ^2 La minimización de la suma de cuadrados del error $\boldsymbol{\varepsilon}'\boldsymbol{\varepsilon}$ no provee un estimador de σ^2 ; sin embargo, podemos obtener un estimador insesgado de σ^2 basado en el estimador de mínimos cuadrados de $\boldsymbol{\beta}$. Por el supuesto de que $cov(\mathbf{y}) = \sigma^2\mathbf{I}$, se tiene que σ^2 es el mismo para cada $y_i, i = 1, 2, \dots, n$. Sabemos que σ^2 está definido por $\sigma^2 = E[y_i - E(y_i)]^2$, y como

$$E(y_i) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} = \mathbf{x}'_i \boldsymbol{\beta}$$

donde \mathbf{x}'_i es la i -ésima fila de \mathbf{X} . Por lo tanto

$$\sigma^2 = E[y_i - \mathbf{x}'_i \boldsymbol{\beta}]^2$$

Podemos estimar σ^2 por

$$s^2 = \frac{1}{n-p} \sum_{i=1}^n (y_i - \mathbf{x}'_i \hat{\boldsymbol{\beta}})^2 \quad (3.11)$$

donde n es el tamaño de la muestra y p es el número de parámetros. En términos matriciales se tiene

$$s^2 = \frac{1}{n-p} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})' (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \quad (3.12)$$

lo cual es igual a

$$s^2 = \frac{\mathbf{y}'\mathbf{y} - \hat{\boldsymbol{\beta}}' \mathbf{X}'\mathbf{y}}{n-p} = \frac{SCE}{n-p} \quad (3.13)$$

donde $SCE = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})' (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{y}'\mathbf{y} - \hat{\boldsymbol{\beta}}' \mathbf{X}'\mathbf{y}$. Con el denominador $n-p$, s^2 es un estimador insesgado de σ^2 , lo cual se demuestra a continuación.

Teorema 3.5 Si s^2 es definido por (3.11), (3.12) o (3.13) y si $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$ y $cov(\mathbf{y}) = \sigma^2\mathbf{I}$, entonces

$$E(s^2) = \sigma^2 \quad (3.14)$$

Prueba. SCE se puede escribir como la siguiente forma cuadrática

$$SCE = \mathbf{y}'\mathbf{y} - [(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}]' \mathbf{X}'\mathbf{y} \quad (3.15)$$

$$= \mathbf{y}'\mathbf{y} - \mathbf{y}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} = \mathbf{y}'[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}']\mathbf{y} \quad (3.16)$$

Aplicando la Esperanza de una forma cuadrática, se tiene que

$$\begin{aligned}
 E(SCE) &= \text{tr}\{[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\sigma^2\mathbf{I}\} + E(\mathbf{y}')[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']E(\mathbf{y}) \\
 &= \sigma^2 \text{tr}[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'] + \beta' \mathbf{X}' [\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'] \mathbf{X} \beta \\
 &= \sigma^2 \{n - \text{tr}[\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\} + \beta' \mathbf{X}' \mathbf{X} \beta - \beta' \mathbf{X}' \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \mathbf{X} \beta \\
 &= \sigma^2 \{n - \text{tr}[\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}]\} + \beta' \mathbf{X}' \mathbf{X} \beta - \beta' \mathbf{X}' \mathbf{X} \beta \\
 &= \sigma^2 [n - \text{tr}(\mathbf{I}_p)] = \sigma^2 (n - p)
 \end{aligned}$$

despejando se obtiene el resultado esperado.

Corolario 3.2 Un estimador insesgado de $\text{cov}(\hat{\beta})$ está dado por

$$\hat{\text{cov}}(\hat{\beta}) = s^2 (\mathbf{X}'\mathbf{X})^{-1} \quad (3.17)$$

Note la correspondencia entre $n - p$ y $\mathbf{y}'\mathbf{y} - \hat{\beta}'\mathbf{X}'\mathbf{y}$; hay n terminos en $\mathbf{y}'\mathbf{y}$ y p términos en $\hat{\beta}'\mathbf{X}'\mathbf{y} = \hat{\beta}'\mathbf{X}'\mathbf{X}\hat{\beta}$. Una propiedad de la muestra es que cada x (y $\hat{\beta}$) adicional en el modelo reduce la SCE .

Ya que SCE es una función cuadrática de \mathbf{y} este no es un mejor estimador lineal insesgado. La propiedad óptima de s^2 está dada en el siguiente teorema.

Teorema 3.6 Si $E(\varepsilon) = 0$, $\text{cov}(\varepsilon) = \sigma^2\mathbf{I}$ y $E(\varepsilon_i^4) = 3\sigma^4$ para el modelo lineal $\mathbf{y} = \mathbf{X}\beta + \varepsilon$, entonces s^2 en (3.12) o (3.13) es el mejor (mínima varianza) estimador cuadrático insesgado de σ^2 .

Prueba. Ver Graybill (1954).

■ EJEMPLO 3.4

Para los datos en 2.2, tenemos

$$\begin{aligned}
 SCE &= \mathbf{y}'\mathbf{y} - \hat{\beta}'\mathbf{X}'\mathbf{y} \\
 &= 840 - \begin{pmatrix} 5.3754 & 3.0118 & -1.2855 \end{pmatrix} \begin{pmatrix} 90 \\ 482 \\ 872 \end{pmatrix} \\
 &= 840 - 814.541 = 25.459
 \end{aligned}$$

entonces

$$s^2 = \frac{SCE}{n - p} = \frac{25.459}{12 - 3} = 2.829$$

Nota 3.1 Las ecuaciones para estimar β y σ^2 son a menudo llamadas ecuaciones normales o ecuaciones de mínimos cuadrados. el primer termino, comúnmente usado, se refiere a una interpretación geométrica de la maximización. El segundo termino surge al notar que estas son las ecuaciones estacionarias para la minimización de $Q(\beta) = (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta)$ con respecto a β y reconociendo que $Q(\beta)$ es la suma de las diferencias cuadráticas entre las observaciones y sus valores esperados. Así, $\hat{\beta}$ es también llamado el estimador de mínimos cuadrados de β . Una interpretación geométrica de la estimación de mínimos cuadrados está dada en el capítulo 2 en la discusión de modelos de regresión de Hocking.

3.1.2 Geometría de los Mínimos Cuadrados

En las secciones anteriores presentamos el modelo de regresión lineal múltiple como la ecuación matricial $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$. Definimos el principio de la estimación de mínimos cuadrados en términos de desviaciones desde el modelo, y luego usando cálculo matricial y álgebra de matrices derivamos los estimadores de $\boldsymbol{\beta}$ y σ^2 . Ahora presentamos una derivación alterna pero equivalente de estos estimadores basado completamente en ideas geométricas.

Es importante aclarar en primer lugar lo que no es el enfoque geométrico de los mínimos cuadrados. En dos dimensiones, ilustramos el principio de mínimos cuadrados al crear un gráfico de dispersión de dos dimensiones de los n puntos $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$. Luego visualizamos la recta de regresión de mínimos cuadrados como la línea recta que mejor se ajusta a los datos. Este enfoque puede generalizarse para presentar la estimación de mínimos cuadrados en regresión lineal múltiple sobre la base del hiperplano en un espacio p -dimensional que mejor se ajusta a los n puntos $(x_{11}, x_{12}, \dots, x_{1k}, y_1), (x_{21}, x_{22}, \dots, x_{2k}, y_2), \dots, (x_{n1}, x_{n2}, \dots, x_{nk}, y_n)$. Aunque este enfoque es algo usado en visualizar regresión lineal múltiple, el enfoque geométrico para la estimación de mínimos cuadrados en regresión lineal múltiple no envuelve esta generalización de alta dimensión.

El enfoque geométrico discutido más abajo es atractivo debido a su elegancia matemática. Por ejemplo, el estimador es derivado sin el uso de cálculo matricial. Además, el enfoque geométrico proporciona una visión más profunda de la inferencia estadística. Varios métodos de estadística avanzada incluidos suavización kernel (Eubank y Eubank 1999), análisis de Fourier (Bloomfield 2000), y análisis Wavelet (Ogden 1997) pueden entenderse como generalizaciones de este enfoque geométrico. El enfoque geométrico de modelos lineales fue primero propuesto por Fisher (Mahalanobis 1964). Christensen (1996) y Jammalamadaka y Sengupta (2003) discuten el modelo estadístico lineal casi completamente desde la perspectiva geométrica.

3.1.2.1 Espacio del Parámetro, Espacio de Datos y Espacio de Predicción El enfoque geométrico de mínimos cuadrados comienza con dos espacios de altas dimensiones, un espacio p -dimensional y un espacio n -dimensional. El vector de parámetros desconocido $\boldsymbol{\beta}$ puede verse como un punto en el espacio p -dimensional, con ejes correspondientes a los p coeficientes de regresión $\beta_0, \beta_1, \dots, \beta_k$. Por lo tanto llamamos a este espacio el *espacio del parámetro* (ver figura ***). De manera similar, el vector de datos \mathbf{y} puede verse como un punto en el espacio n -dimensional con ejes correspondientes a las n observaciones. Llamamos este espacio el *espacio de datos*.

La matriz \mathbf{X} del modelo de regresión múltiple puede escribirse como una matriz particionada en términos de sus p columnas como

$$\mathbf{X} = \begin{pmatrix} \mathbf{j} & \mathbf{x}_1 & \mathbf{x}_2 & \mathbf{x}_3 & \dots & \mathbf{x}_k \end{pmatrix}$$

Las columnas de \mathbf{X} , incluyendo a \mathbf{j} , son todos vectores n -dimensionales y son, por lo tanto, puntos en el espacio de datos. Note que debido a que asumimos que \mathbf{X} es de rango p , estos vectores son linealmente independientes. El conjunto de todas las posibles combinaciones lineales de las columnas de \mathbf{X} constituyen un subconjunto del espacio de datos (ver sección 2.3 del Rencher). Elementos de este subconjunto pueden escribirse como

$$\mathbf{X}\mathbf{b} = b_0\mathbf{j} + b_1\mathbf{x}_1 + b_2\mathbf{x}_2 + \dots + b_k\mathbf{x}_k$$

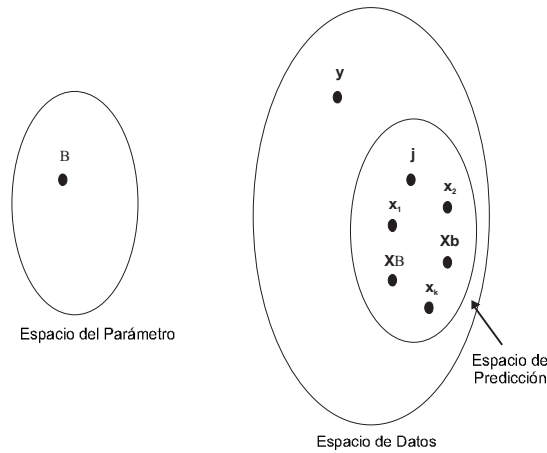


Figure 3.1 Espacio del parámetro, espacio de datos y espacio de predicción con elementos representativos

donde \mathbf{b} es cualquier vector $p \times 1$, es decir, es cualquier vector en el espacio del parámetro. Este subconjunto tiene el estatus de subespacio ya que este es cerrado bajo la adición y la multiplicación escalar (Harville 1997, pp. 28 - 29). Este subconjunto se dice ser el subespacio generado por las columnas de \mathbf{X} , y nosotros llamaremos este subespacio como el *espacio de predicción*. Las columnas de \mathbf{X} constituyen un *conjunto base* para el espacio de predicción.

3.1.2.2 Interpretación Geométrica del Modelo de Regresión Lineal Múltiple El modelo de regresión lineal múltiple establece que \mathbf{y} es igual a un vector en el espacio de predicción, $E(\mathbf{y}) = \mathbf{X}\beta$, mas un vector de errores aleatorios, ε (figura 3.2). El problema es que ni β ni ε son conocidos. Sin embargo, el vector de datos \mathbf{y} , el cual no está en el espacio de predicción, es conocido. Y se sabe que $E(\mathbf{y})$ está en el espacio de predicción.

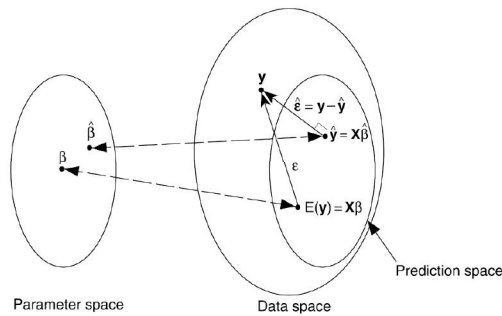


Figure 3.2 Relación geométrica de los vectores asociados con el modelo de regresión lineal múltiple

La regresión lineal múltiple puede entenderse geoméricamente como el proceso de hallar una estimación sensible de $E(\mathbf{y})$ en el espacio de predicción y luego determinar el vector

en el espacio del parámetro que esta asociado con esta estimación (figura 3.2). La estimación de $E(\mathbf{y})$ es denotada como $\hat{\mathbf{y}}$, y el vector asociado en el espacio del parámetro es denotado por $\hat{\boldsymbol{\beta}}$.

Una idea geométrica razonable es estimar $E(\mathbf{y})$ usando el punto en el espacio de predicción que esta mas cercano a \mathbf{y} . Resulta en que $\hat{\mathbf{y}}$, el punto más cercano en el espacio de predicción a \mathbf{y} , puede hallarse al notar que el vector diferencia $\hat{\boldsymbol{\varepsilon}} = \mathbf{y} - \hat{\mathbf{y}}$ debe ser ortogonal (perpendicular) al espacio de predicción (Harville 1997, p. 170). Además, ya que el espacio de predicción es generado por las columnas de \mathbf{X} , el punto $\hat{\mathbf{y}}$ debe ser tal que $\hat{\boldsymbol{\varepsilon}}$ es ortogonal a las columnas de \mathbf{X} . Es decir, que $\hat{\mathbf{y}}$ es tal que

$$\mathbf{X}'\hat{\boldsymbol{\varepsilon}} = 0$$

o

$$\mathbf{X}'(\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{X}'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = 0$$

lo cual implica que

$$\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{y}$$

Por lo tanto, usando ideas solamente geométricas, obtenemos las ecuaciones normales y en consecuencia el usual estimado de mínimos cuadrados $\hat{\boldsymbol{\beta}}$. Podemos entonces calcular $\hat{\mathbf{y}}$ como $\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{H}\mathbf{y}$. Además, $\hat{\boldsymbol{\varepsilon}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = (\mathbf{I} - \mathbf{H})\mathbf{y}$ puede ser tomado como un estimador de $\boldsymbol{\varepsilon}$. Ya que $\hat{\boldsymbol{\varepsilon}}$ es un vector en el espacio $(n - p)$ dimensional, es razonable estimar a σ^2 como el cuadrado de la longitud de $\hat{\boldsymbol{\varepsilon}}$ dividido por $n - p$. En otras palabras, un estimador sensible de σ^2 es $s^2 = \mathbf{y}'(\mathbf{I} - \mathbf{H})\mathbf{y}/(n - p)$.

3.1.3 Estimadores por el método de máxima verosimilitud

Hasta el momento no se ha hecho ningún supuesto sobre la distribución de las variables aleatorias y_1, y_2, \dots, y_n .

Supongamos que $\boldsymbol{\varepsilon} \sim N(0, \sigma^2\mathbf{I})$ lo cual es equivalente a $\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$. Bajo este supuesto podemos entonces hallar los estimadores máximo verosimiles de $\boldsymbol{\beta}$ y σ^2 .

Como \mathbf{Y} se distribuye $N(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$, la función de verosimilitud es:

$$\begin{aligned} L &= L(\boldsymbol{\beta}, \sigma^2 : y_1, y_2, \dots, y_n; x_1, x_2, \dots, x_n) = L(\boldsymbol{\beta}, \sigma^2 : \mathbf{y}; \mathbf{X}) \\ &= \left(\frac{1}{2\pi\sigma^2}\right)^{n/2} \exp\left\{-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right\} \end{aligned} \quad (3.18)$$

Ya que L es una función cuyo rango es positivo y la función logaritmo es una función monótona que mantiene el orden, es decir si $f(x)$ alcanza un máximo en x_0 , entonces, $\log f(x)$ alcanza también el máximo en x_0 , por lo tanto, en vez de hallar el máximo en L , se hallará el máximo del $\log L$, pues los cálculos son más fáciles. Así,

$$\begin{aligned}
\log L(\beta, \sigma^2 : \mathbf{y}; \mathbf{X}) &= \log \left\{ \left(\frac{1}{2\pi\sigma^2} \right)^{n/2} \exp \left[-\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\beta) \right] \right\} \\
&= \log \left\{ \left(\frac{1}{2\pi\sigma^2} \right)^{n/2} \right\} + \log \exp \left[-\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\beta) \right] \\
&= \frac{n}{2} \log \left(\frac{1}{2\pi\sigma^2} \right) - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\beta) \\
&= -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\beta) \quad (3.19)
\end{aligned}$$

El espacio del parámetro es

$$\Omega = \{(\beta, \sigma^2) : \sigma^2 > 0, -\infty < \beta_i < \infty; i_1, 2, \dots, p\} \quad (3.20)$$

Para hallar los valores de β y σ^2 en Ω que maximizan la función de verosimilitud se iguala a cero las derivadas parciales con respecto a β y σ^2 .

Antes de derivar se reescribe la ecuación

$$\log L = \frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y}'\mathbf{y} - \mathbf{y}'\mathbf{X}\beta - \beta'\mathbf{X}'\mathbf{y} + \beta'\mathbf{X}'\mathbf{X}\beta)$$

Haciendo $\mathbf{V} = \mathbf{y}'\mathbf{X}$ y $\mathbf{W} = \mathbf{X}'\mathbf{X}$

$$\log L = \frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y}'\mathbf{y} - \mathbf{V}'\beta - \beta'\mathbf{V} + \beta'\mathbf{W}\beta)$$

Derivando con respecto a β y σ^2 e igualando a cero (ver ??), se obtiene

$$\begin{aligned}
\frac{\partial \log L}{\partial \beta} &= -\frac{1}{2\sigma^2} (-\mathbf{V} - \mathbf{V} + 2\mathbf{W}\beta) = -\frac{1}{2\sigma^2} (-2\mathbf{V} + 2\mathbf{W}\beta) \\
&= \frac{1}{\sigma^2} (\mathbf{V} - \mathbf{W}\beta) = \frac{1}{\sigma^2} (\mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{X}\beta) = 0 \quad (3.21)
\end{aligned}$$

$$\frac{\partial \log L}{\partial \sigma^2} = -\frac{n}{2} \frac{1}{\sigma^2} + \frac{1}{2\sigma^4} (\mathbf{y} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\beta) = 0 \quad (3.22)$$

Sean $\tilde{\beta}$ y $\tilde{\sigma}^2$ las soluciones de las ecuaciones para β y σ^2 cuando las derivadas son en conjunto iguales a cero. Entonces

$$\frac{1}{\tilde{\sigma}^2} (\mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{X}\tilde{\beta}) = 0 \Rightarrow \mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{X}\tilde{\beta} = 0 \Rightarrow \mathbf{X}'\mathbf{X}\tilde{\beta} = \mathbf{X}'\mathbf{y} \quad (3.23)$$

$$-\frac{n}{2} \frac{1}{\tilde{\sigma}^2} = -\frac{1}{2\tilde{\sigma}^4} (\mathbf{y} - \mathbf{X}\tilde{\beta})' (\mathbf{y} - \mathbf{X}\tilde{\beta}) \Rightarrow \tilde{\sigma}^2 = \frac{1}{n} (\mathbf{y} - \mathbf{X}\tilde{\beta})' (\mathbf{y} - \mathbf{X}\tilde{\beta}) \quad (3.24)$$

Ahora como \mathbf{X} es de rango p , entonces $\mathbf{X}'\mathbf{X}$ es también de rango p la cual es su dimensión. Así, $\mathbf{X}'\mathbf{X}$ es invertible. (ver ?? y ??)

Por lo tanto, los estimadores máximos verosímiles de β y σ^2 son

$$\tilde{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \quad (3.25)$$

$$\begin{aligned} \tilde{\sigma}^2 &= \frac{1}{n} [\mathbf{y} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}]' [\mathbf{y} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}] \\ &= \frac{1}{n} [\mathbf{y}' - \mathbf{y}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'] [\mathbf{y} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}] \\ &= \frac{1}{n} [\mathbf{y}' - \mathbf{y}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'] [\mathbf{y} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}] \\ &= \frac{1}{n} \mathbf{y}' [\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'] [\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'] \mathbf{y} \end{aligned} \quad (3.26)$$

Sea $\mathbf{K} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$. Se puede notar que \mathbf{K} es idempotente (ver ??), pues

$$\mathbf{K}\mathbf{K} = \mathbf{X} \overbrace{(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}}^{\mathbf{I}} (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' = \mathbf{K}$$

Entonces,

$$\tilde{\sigma}^2 = \frac{1}{n} \mathbf{y}' (\mathbf{I} - \mathbf{K})' (\mathbf{I} - \mathbf{K}) \mathbf{y}$$

Además,

$$(\mathbf{I} - \mathbf{K})' (\mathbf{I} - \mathbf{K}) = \mathbf{I} - \mathbf{K} - \mathbf{K} + \mathbf{K}\mathbf{K} = \mathbf{I} - \mathbf{K} - \mathbf{K} + \mathbf{K} = \mathbf{I} - \mathbf{K}$$

Así,

$$\tilde{\sigma}^2 = \frac{1}{n} \mathbf{y}' (\mathbf{I} - \mathbf{K}) \mathbf{y}$$

Haciendo $\mathbf{M} = \mathbf{I} - \mathbf{K}$ se tiene que

$$\tilde{\sigma}^2 = \frac{1}{n} \mathbf{y}' \mathbf{M} \mathbf{y}$$

Para demostrar que $\tilde{\beta}$ y $\tilde{\sigma}^2$ son los estimadores máximo verosímiles de β y σ^2 se debe demostrar que ellos maximizan la función de verosimilitud (La prueba esta dada en la sección 9.3 del Graybill).

3.1.3.1 Propiedades de los estimadores máximo verosímiles Ahora se examinan las propiedades de los estimadores $\tilde{\beta}$ y $\tilde{\sigma}^2$ y se determinan sus distribuciones.

3.1.3.2 Suficiencia Primero se demostrará que los $p + 1$ estimadores $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_p, \tilde{\sigma}^2$ son estadísticos suficientes. Para ello se usa el siguiente teorema

Teorema 3.7 Sea $f_{\mathbf{y}}(y; \theta)$ la fdp de una muestra observable y_1, y_2, \dots, y_n . Los estadísticos S_1, S_2, \dots, S_r son estadísticos suficientes si y sólo si la fdp puede factorizarse como

$$f_{\mathbf{Y}}(y; \theta) = g(S_1, S_2, \dots, S_r; \theta) h(y)$$

donde $g(S_1, S_2, \dots, S_r; \theta)$ es una función no negativa de solamente los estadísticos S_1, S_2, \dots, S_r y el parámetro θ , y $h(\mathbf{y})$ no es una función de θ .

Entonces se quiere demostrar que la fdp de \mathbf{y} , llámese $f(\mathbf{y} : \mathbf{X}; \beta; \sigma^2)$, se puede factorizar como $g(\tilde{\beta}, \tilde{\sigma}^2 : \beta; \sigma^2)h(\mathbf{y})$ donde $g(\tilde{\beta}, \tilde{\sigma}^2 : \beta; \sigma^2)$ contiene las observaciones \mathbf{y} en la forma de $\tilde{\beta}$ y $\tilde{\sigma}^2$ solamente, y $h(\mathbf{y})$ no contiene los parámetros β y σ^2 .

La fdp de \mathbf{y} está dada por

$$f(\mathbf{y} : \mathbf{X}; \beta; \sigma^2) = \left(\frac{1}{2\pi\sigma^2} \right)^{n/2} \exp \left[-\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\beta) \right] \quad (3.27)$$

Trabajando con una porción del exponente se obtiene que

$$\begin{aligned} (\mathbf{y} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\beta) &= (\mathbf{y} - \mathbf{X}\hat{\beta} + \mathbf{X}\hat{\beta} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\hat{\beta} + \mathbf{X}\hat{\beta} - \mathbf{X}\beta) \\ &= [(\mathbf{y} - \mathbf{X}\hat{\beta}) + \mathbf{X}\hat{\beta} - \mathbf{X}\beta]' [(\mathbf{y} - \mathbf{X}\hat{\beta}) + \mathbf{X}\hat{\beta} - \mathbf{X}\beta] \\ &= [(\mathbf{y} - \mathbf{X}\hat{\beta}) - \mathbf{X}(\beta - \hat{\beta})]' [(\mathbf{y} - \mathbf{X}\hat{\beta}) - \mathbf{X}(\beta - \hat{\beta})] \\ &= [(\mathbf{y} - \mathbf{X}\hat{\beta})' - (\beta - \hat{\beta})' \mathbf{X}'] [(\mathbf{y} - \mathbf{X}\hat{\beta}) - \mathbf{X}(\beta - \hat{\beta})] \\ &= (\mathbf{y} - \mathbf{X}\hat{\beta})' (\mathbf{y} - \mathbf{X}\hat{\beta}) - (\mathbf{y} - \mathbf{X}\hat{\beta})' \mathbf{X} (\beta - \hat{\beta}) \\ &\quad - (\beta - \hat{\beta})' \mathbf{X}' (\mathbf{y} - \mathbf{X}\hat{\beta}) + (\beta - \hat{\beta})' \mathbf{X}' \mathbf{X} (\beta - \hat{\beta}) \\ &= (\mathbf{y} - \mathbf{X}\hat{\beta})' (\mathbf{y} - \mathbf{X}\hat{\beta}) + (\beta - \hat{\beta})' \mathbf{X}' \mathbf{X} (\beta - \hat{\beta}) \\ &= n\tilde{\sigma}^2 + (\beta - \hat{\beta})' \mathbf{X}' \mathbf{X} (\beta - \hat{\beta}) \end{aligned} \quad (3.28)$$

ya que,

$$\begin{aligned} (\mathbf{y} - \mathbf{X}\hat{\beta})' \mathbf{X} &= \mathbf{y}' \mathbf{X} - \hat{\beta}' \mathbf{X}' \mathbf{X} = \mathbf{y}' \mathbf{X} - [(\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y}]' \mathbf{X}' \mathbf{X} \\ &= \mathbf{y}' \mathbf{X} - \mathbf{y}' \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{X} = \mathbf{y}' \mathbf{X} - \mathbf{y}' \mathbf{X} = 0 \end{aligned}$$

y de la misma manera, $\mathbf{X}' (\mathbf{y} - \mathbf{X}\hat{\beta}) = 0$.

Sustituyendo 3.28 en 3.27 se obtiene

$$\begin{aligned} f(\mathbf{y} : \mathbf{X}; \beta; \sigma^2) &= \left(\frac{1}{2\pi\sigma^2} \right)^{n/2} \exp \left[-\frac{1}{2\sigma^2} [n\tilde{\sigma}^2 + (\beta - \hat{\beta})' \mathbf{X}' \mathbf{X} (\beta - \hat{\beta})] \right] \\ &= g(\tilde{\beta}, \tilde{\sigma}^2 : \beta; \sigma^2) h(\mathbf{y}) \end{aligned}$$

donde $h(\mathbf{Y}) = 1$. Por lo tanto, $\tilde{\beta}$ y $\tilde{\sigma}^2$ son estadísticos suficientes.

3.1.3.3 Completitud. Ahora se probará que $\tilde{\beta}$ y $\tilde{\sigma}^2$ son completos. Para ello es necesario recordar la siguiente definición y el subsiguiente teorema

Definición 3.1 (Familia Exponencial de Densidades) Sea $f_{\mathbf{y}}(y; \theta)$, $\theta \in \Omega$ la fdp conjunta de una muestra observable y_1, y_2, \dots, y_n ; donde $\theta' = [\theta_1, \theta_2, \dots, \theta_k]$. La fdp $f_{\mathbf{y}}(y; \theta)$

se define a pertenecer a una familia exponencial de densidades si esta puede escribirse como

$$f_{\mathbf{y}}(\mathbf{y}; \theta) = h(\theta)g(\mathbf{y}) \exp \left[\sum_{j=1}^k S_j(\mathbf{y})P_j(\theta) \right] \quad a_i < y_i < b_i; i = 1, 2, \dots, n \quad (3.29)$$

donde

1. a_i y b_i no dependen de los θ_j , a_i puede ser $-\infty$ y b_i puede ser $+\infty$.
2. Ω , el espacio del parámetro, contiene un espacio K -dimensional no degenerado (rectángulo).
3. $h(\theta)$ y los $P_j(\theta)$ no dependen de los y_j .
4. $g(\mathbf{y})$ y los $S_j(\mathbf{y})$ no dependen de los θ_i ; $g(\mathbf{y})$ es no negativa.

Teorema 3.8 Sea Y_1, Y_2, \dots, Y_n una muestra observable con fdp $f_Y(\mathbf{y}; \theta)$, $\theta \in \Omega$. Supóngase que $f_Y(\mathbf{y}; \theta)$ pertenece a una familia exponencial, esto es, $f_Y(\mathbf{y}; \theta)$ puede escribirse en la forma dada en la ecuación 3.29 y satisface a, b, c y d. Supóngase que las siguientes condiciones también se satisfacen

1. Los $S_j(\mathbf{y})$ son funcionalmente independientes para $j = 1, 2, \dots, k$
2. $\partial[S_j(\mathbf{y})]/\partial y_i$ existe y es continua para todos los $j = 1, 2, \dots, k$ y $i = 1, 2, \dots, n$
3. $P_j(\theta)$ es una función continua de θ para $j = 1, 2, \dots, k$
4. Si se define $\gamma' = [P_1(\theta), P_2(\theta), \dots, P_k(\theta)]$, entonces el conjunto de valores de γ (Como θ varía sobre Ω) contiene un rectángulo K dimensional no degenerado.

Entonces el conjunto de K estadísticos $S_1(\mathbf{y}), S_2(\mathbf{y}), \dots, S_k(\mathbf{y})$, es un conjunto de estadísticos mínimo suficientes.

Ahora, para establecer la completitud se enuncia el siguiente teorema

Teorema 3.9 Sean y_1, y_2, \dots, y_n una muestra observable con fdp $f_Y(\mathbf{y}; \theta)$, $\theta \in \Omega$. Si $f_Y(\mathbf{y}; \theta)$ pertenece a una familia exponencial como se define en 3.1 y si las condiciones a, b, c, d del teorema 3.8 se satisfacen, entonces los K estadísticos $S_1(\mathbf{y}), S_2(\mathbf{y}), \dots, S_k(\mathbf{y})$, son estadísticos suficientes completos.

Por lo tanto, para probar la completitud primero se debe probar que $f(\mathbf{y} : \mathbf{X}; \beta; \sigma^2)$ pertenece a una familia exponencial, es decir, $f(\mathbf{y} : \mathbf{X}; \beta; \sigma^2)$ debe tener la forma

$$f(\mathbf{y} : \mathbf{X}; \beta; \sigma^2) = h(\beta, \sigma^2)g(\mathbf{y}, \mathbf{X}) \exp [S_1(\mathbf{y}, \mathbf{X})P_1(\beta, \sigma^2) + S_2(\mathbf{y}, \mathbf{X})P_2(\beta, \sigma^2)] \quad (3.30)$$

La fdp de \mathbf{y} es

$$\begin{aligned}
f(\mathbf{y}; \mathbf{X}; \boldsymbol{\beta}; \sigma^2) &= \left(\frac{1}{2\pi\sigma^2}\right)^{n/2} \exp\left[-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right] \\
&= \left(\frac{1}{2\pi\sigma^2}\right)^{n/2} \exp\left[-\frac{\mathbf{y}'\mathbf{y}}{2\sigma^2} + \frac{2\mathbf{y}'\mathbf{X}\boldsymbol{\beta}}{2\sigma^2} - \frac{\boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta}}{2\sigma^2}\right] \\
&= \left(\frac{1}{2\pi}\right)^{n/2} \left(\frac{1}{\sigma^2}\right)^{n/2} \exp\left[-\frac{\boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta}}{2\sigma^2}\right] \exp\left[-\frac{\mathbf{y}'\mathbf{y}}{2\sigma^2} + \frac{2\mathbf{y}'\mathbf{X}\boldsymbol{\beta}}{2\sigma^2}\right] \\
&= \left(\frac{1}{2\pi}\right)^{n/2} \left(\frac{1}{\sigma^2}\right)^{n/2} \exp\left[-\frac{\boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta}}{2\sigma^2}\right] \\
&\quad \exp\left[-\mathbf{y}'\mathbf{y}\frac{1}{2\sigma^2} + \mathbf{y}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\frac{\boldsymbol{\beta}}{\sigma^2}\right] \\
&= \left(\frac{1}{2\pi}\right)^{n/2} \left(\frac{1}{\sigma^2}\right)^{n/2} \exp\left[-\frac{\boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta}}{2\sigma^2}\right] \exp\left[-\mathbf{y}'\mathbf{y}\frac{1}{2\sigma^2} + \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{X}\frac{\boldsymbol{\beta}}{\sigma^2}\right]
\end{aligned}$$

Así,

$$\begin{aligned}
h(\boldsymbol{\beta}, \sigma^2) &= \left(\frac{1}{\sigma^2}\right)^{n/2} \exp\left[-\frac{\boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta}}{2\sigma^2}\right], g(\mathbf{y}, \mathbf{X}) = \left(\frac{1}{2\pi}\right)^{n/2}, S_1(\mathbf{y}) = \mathbf{y}'\mathbf{y}, P_1(\boldsymbol{\beta}, \sigma^2) = \\
&-\frac{1}{2\sigma^2}, S_2(\mathbf{Y}) = \hat{\boldsymbol{\beta}}' \text{ y } P_2(\boldsymbol{\beta}, \sigma^2) = \mathbf{X}'\mathbf{X}\frac{\boldsymbol{\beta}}{\sigma^2}.
\end{aligned}$$

Por lo tanto, $\mathbf{y}'\mathbf{y}$ y $\hat{\boldsymbol{\beta}}$ son completos.

3.1.3.4 Distribución de $\hat{\boldsymbol{\beta}}$ y $\hat{\sigma}^2$

Distribución de $\hat{\boldsymbol{\beta}}$. Para hallar la distribución de $\hat{\boldsymbol{\beta}}$ es necesario primero enunciar el siguiente teorema

Teorema 3.10 Sea \mathbf{X} un vector aleatorio $p \times 1$ que se distribuye $N(\mu, \Sigma)$, donde Σ tiene rango k , Sea \mathbf{B} cualquier matriz $q \times p$ de constantes, y sea \mathbf{b} cualquier vector $q \times 1$ de constantes. Entonces el vector \mathbf{y} $q \times 1$ definido por $\mathbf{y} = \mathbf{B}\mathbf{X} + \mathbf{b}$ (\mathbf{Y} es una función lineal de \mathbf{X}) se distribuye $N(\mathbf{B}\mu + \mathbf{b}, \mathbf{B}\Sigma\mathbf{B}')$

Como $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$, donde $\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}; \sigma^2\mathbf{I})$ y haciendo $\mathbf{B} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ (\mathbf{B} es una matriz $p \times n$) y $\mathbf{b} = 0$, entonces $\hat{\boldsymbol{\beta}} = \mathbf{B}\mathbf{Y}$ es una función lineal de \mathbf{Y} , por lo tanto, por el teorema 3.10 se tiene que $\hat{\boldsymbol{\beta}}$ se distribuye normal, cuyos parámetros son

$$E(\hat{\boldsymbol{\beta}}) = E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}] = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E(\mathbf{Y}) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\boldsymbol{\beta} = \boldsymbol{\beta} \tag{3.31}$$

lo cual indica que $\hat{\boldsymbol{\beta}}$ es un estimador insesgado de $\boldsymbol{\beta}$,

$$\begin{aligned}
Cov(\tilde{\beta}) &= E[(\hat{\beta} - E(\hat{\beta}))(\hat{\beta} - E(\hat{\beta}))'] = E[(\hat{\beta} - \beta)(\hat{\beta} - \beta)'] \\
&= E(\hat{\beta}\hat{\beta}' - \hat{\beta}\beta' - \beta\hat{\beta}' + \beta\beta') = E(\hat{\beta}\hat{\beta}') - \beta\beta' - \beta\beta' + \beta\beta' \\
&= E(\hat{\beta}\hat{\beta}') - \beta\beta' = E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}\mathbf{Y}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}] - \beta\beta' \\
&= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E(\mathbf{Y}\mathbf{Y}')\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} - \beta\beta' \\
&= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'[Cov(\mathbf{Y}) + E(\mathbf{y})E(\mathbf{y}')]\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} - \beta\beta' \\
&= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'[\sigma^2\mathbf{I} + (\mathbf{X}\beta)(\mathbf{X}\beta)']\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} - \beta\beta' \\
&= \sigma^2(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\beta\beta'\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} - \beta\beta' \\
&= (\mathbf{X}'\mathbf{X})^{-1}\sigma^2 + \beta\beta' - \beta\beta' = (\mathbf{X}'\mathbf{X})^{-1}\sigma^2 \tag{3.32}
\end{aligned}$$

Por lo tanto,

$$\hat{\beta} \sim N(\beta, (\mathbf{X}'\mathbf{X})^{-1}\sigma^2)$$

Distribución de $\hat{\sigma}^2$. Para hallar la distribución de $\hat{\sigma}^2$ se usa el siguiente teorema

Teorema 3.11 Sea \mathbf{y} un vector aleatorio $n \times 1$ que se distribuye $N(\mu, \Sigma)$, donde Σ tiene rango n . Entonces la forma cuadrática $U = \mathbf{y}'\mathbf{A}\mathbf{y} \sim \chi^2(p, \lambda)$, donde $\lambda = \frac{1}{2}\mu'\mathbf{A}\mu$, si y sólo si cualquiera de las siguientes 3 condiciones se satisfacen:

1. $\mathbf{A}\Sigma$ es una matriz idempotente de rango p .
2. $\Sigma\mathbf{A}$ es una matriz idempotente de rango p .
3. Σ es una c -inversa de \mathbf{A} y \mathbf{A} tiene rango p .

Como $\tilde{\sigma}^2$ está dado por

$$\tilde{\sigma}^2 = \frac{1}{n}\mathbf{y}'\mathbf{M}\mathbf{y}$$

donde $\mathbf{y} \sim N(N(\mathbf{X}\beta; \sigma^2\mathbf{I}))$ y \mathbf{M} es una matriz simétrica e idempotente $n \times n$ cuyo rango es

$$r(\mathbf{M}) = r(\mathbf{I} - \mathbf{K}) = tr(\mathbf{I} - \mathbf{K}) = tr(\mathbf{I}) - tr(\mathbf{K}) = r(\mathbf{I}) - r(\mathbf{K}) = n - p$$

entonces por el teorema anterior

$$U = \frac{\mathbf{y}'\mathbf{M}\mathbf{y}}{\sigma^2} \sim \chi^2(n - p, \lambda)$$

donde,

$$\begin{aligned}
\lambda &= \frac{1}{2}[E(\mathbf{y})]'M[E(\mathbf{y})] = \frac{1}{2}(\mathbf{X}\boldsymbol{\beta})'(\mathbf{I}-\mathbf{K})(\mathbf{X}\boldsymbol{\beta}) \\
&= \frac{1}{2}\boldsymbol{\beta}'\mathbf{X}'[\mathbf{I}-\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{X}\boldsymbol{\beta} \\
&= \frac{1}{2}\boldsymbol{\beta}'[\mathbf{X}'\mathbf{X}-\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}] \\
&= \frac{1}{2}\boldsymbol{\beta}'[\mathbf{X}'\mathbf{X}-\mathbf{X}'\mathbf{X}]\boldsymbol{\beta} = 0
\end{aligned}$$

Por lo tanto,

$$U = \frac{\mathbf{y}'M\mathbf{y}}{\sigma^2} \sim \chi^2(n-p)$$

Dado que el primer momento de la distribución chi-cuadrado central son sus grados de libertad, entonces $E(U) = n-p$, así

$$E(U) = E\left(\frac{1}{\sigma^2}\mathbf{y}'M\mathbf{y}\right) = E\left(\frac{1}{\sigma^2}n\hat{\sigma}^2\right) = n-p$$

entonces,

$$E(\hat{\sigma}^2) = \frac{n-p}{n}\sigma^2$$

lo cual indica que $\hat{\sigma}^2$ es un estimador sesgado de σ^2 . Por lo tanto,

$$\hat{\sigma}^2 = \frac{n}{n-p}\hat{\sigma}^2 = \frac{1}{n-p}\mathbf{y}'M\mathbf{y}$$

es un estimador insesgado de σ^2 .

Nota 3.2 $\hat{\sigma}^2$ tiene distintas formas de calcularse

1. $\hat{\sigma}^2 = \frac{1}{n-p}(\mathbf{y}-\mathbf{X}\hat{\boldsymbol{\beta}})'(\mathbf{y}-\mathbf{X}\hat{\boldsymbol{\beta}})$ donde $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{X}^{-}\mathbf{y}$
2. $\hat{\sigma}^2 = \frac{1}{n-p}\mathbf{y}'M\mathbf{y}$ donde $M = \mathbf{I}-\mathbf{K} = \mathbf{I}-\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$
3. $\hat{\sigma}^2 = \frac{1}{n-p}\mathbf{y}'\mathbf{y} - \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{y}$

3.1.3.5 $\hat{\boldsymbol{\beta}}$ y $\hat{\sigma}^2$ son independientes Para demostrar la independencia se hace uso del siguiente teorema

Teorema 3.12 Sea \mathbf{y} un vector aleatorio $n \times 1$ que se distribuye $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, donde $\boldsymbol{\Sigma}$ tiene rango n . Si $\mathbf{B}\boldsymbol{\Sigma}\mathbf{A} = 0$. La forma cuadrática $\mathbf{y}'\mathbf{A}\mathbf{y}$ es independiente de la forma lineal $\mathbf{B}\mathbf{y}$, donde \mathbf{B} es una matriz $q \times n$.

Como

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{X}^{-}\mathbf{y} \\ \hat{\sigma}^2 &= \frac{1}{n-p}\mathbf{y}'\mathbf{M}\mathbf{y}\end{aligned}$$

Entonces haciendo $\mathbf{B} = \mathbf{X}^{-}$, $\Sigma = \sigma^2\mathbf{I}$ y $\mathbf{A} = \mathbf{M}$ entonces,

$$\begin{aligned}\mathbf{B}\Sigma\mathbf{A} &= \mathbf{X}^{-}\sigma^2\mathbf{I}\mathbf{M} = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{I} - \mathbf{K}) \\ &= \sigma^2 [(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'] \\ &= \sigma^2 [(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'] \\ &= \sigma^2(0) = 0\end{aligned}$$

Por lo tanto, $\hat{\beta}$ y $\hat{\sigma}^2$ son independientes.

3.1.3.6 $\hat{\beta}$ y $\hat{\sigma}^2$ son estimadores insesgados de mínima varianza Se puede probar que $\hat{\beta}$ y $\hat{\sigma}^2$ son estimadores insesgados de mínima varianza pero se escapa del objetivo de este curso la demostración de este resultado.

Las propiedades de $\hat{\beta}$ y $\hat{\sigma}^2$ se resumen en el siguiente teorema

Teorema 3.13 Sea $\mathbf{y} = \mathbf{X}\beta + \varepsilon$ donde $\varepsilon \sim N(\mathbf{0}, \sigma^2\mathbf{I})$. Se cumplen las siguientes propiedades

1. $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{X}^{-}\mathbf{y}$ es el estimador de máxima verosimilitud de β .
2. $\hat{\sigma}^2 = \frac{1}{n-p}\mathbf{y}'\mathbf{M}\mathbf{y}$ es el estimador de máxima verosimilitud de σ^2 .
3. $\hat{\beta} \sim N(\beta, (\mathbf{X}'\mathbf{X})^{-1}\sigma^2)$
4. $\frac{(n-p)\hat{\sigma}^2}{\sigma^2} \sim \chi^2(n-p)$
5. $\hat{\beta}$ y $\hat{\sigma}^2$ son independientes.
6. $\hat{\beta}$ y $\hat{\sigma}^2$ son estadísticos suficientes.
7. $\hat{\beta}$ y $\hat{\sigma}^2$ son estadísticos completos.

La prueba de este teorema se desarrollo durante toda la sección.

A continuación se muestra un resultado muy importante sobre las propiedades de los estimadores, previo es necesario enunciar el siguiente teorema

Teorema 3.14 Sea Y_1, Y_2, \dots, Y_n una muestra observable con fdp $f_Y(\mathbf{y} : \theta)$ perteneciente a una familia exponencial, si las condiciones a, b, c y d del teorema 3.8 se satisfacen, entonces los K estadísticos $s_1\mathbf{y}, s_2\mathbf{y}, \dots, s_K\mathbf{y}$ son estadísticos suficientes completos.

Teorema 3.15 Sea $\mathbf{y} = \mathbf{X}\beta + \varepsilon$, donde ε se distribuye $N(\mathbf{0}, \sigma^2\mathbf{I})$, dado en la definición ???. Sea $t(\beta, \sigma^2)$ cualquier función de los parámetros β y σ^2 para los cuales un estimador

insesgado existe. Entonces existe una función de los estadísticos suficientes $\hat{\beta}$ y $\hat{\sigma}^2$, digamos $q(\hat{\beta}, \hat{\sigma}^2)$, que es también un estimador insesgado de $t(\beta, \sigma^2)$. En adición, $q(\hat{\beta}, \hat{\sigma}^2)$ es el estimador insesgado uniformemente de varianza mínima para $t(\beta, \sigma^2)$.

Nota: Por la condición (6) del teorema 3.13 se sabe que no se necesita el vector \mathbf{y} ni la matriz \mathbf{X} excepto para obtener los estadísticos suficientes $\hat{\beta}$ y $\hat{\sigma}^2$. Ya que los $p + 1$ variables aleatorias $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_p, \hat{\sigma}^2$ son estadísticos suficientes, ellos contienen toda la "información" acerca del modelo (para el caso 1) que esta contenida en los $n + np$ elementos de \mathbf{y} y \mathbf{X} . Ya que se conoce la distribución conjunta de las $p + 1$ variables aleatorias $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_p, \hat{\sigma}^2$, se puede usar esa distribución para la estimación puntual, estimación por intervalo y prueba de hipótesis, sin sacrificar "información" por no usar las observaciones originales \mathbf{y} y \mathbf{X} .

3.1.4 Propiedades de los estimadores para muestras grandes

Teorema 3.16 Considere una sucesión de modelos lineal general

$$\mathbf{y}_n = \mathbf{X}_n \beta + \varepsilon_n \quad \varepsilon_n \sim N(\mathbf{0}_n, \sigma^2 \mathbf{I}_n) \quad n = p + 1, p + 2, \dots$$

donde \mathbf{y}_n es un vector $n \times 1$ de variables aleatorias observables, \mathbf{X}_n es una matriz de rango p (para cada n) de variables no aleatorias observables, β es un vector p veces 1 de constantes no observables, y ε_n es un vector $n \times 1$ de variables aleatorias no observables. Sean $\hat{\beta}_n$ y $\hat{\sigma}_n^2$, respectivamente, los estimadores de máxima verosimilitud ($\hat{\sigma}_n^2$ ajustado por sesgo) de β y σ^2 en el n -ésimo modelo. Por lo tanto,

$$\begin{aligned} \hat{\beta}_n &= (\mathbf{X}'_n \mathbf{X}_n)^{-1} \mathbf{X}'_n \mathbf{y}_n \\ \hat{\sigma}_n^2 &= (n - p)^{-1} \mathbf{Y}'_n [\mathbf{I}_n - \mathbf{X}_n (\mathbf{X}'_n \mathbf{X}_n)^{-1} \mathbf{X}'_n] \mathbf{y}_n \quad n = p + 1, p + 2, \dots \end{aligned}$$

1. Si $\lim_{n \rightarrow \infty} (\mathbf{X}'_n \mathbf{X}_n)^{-1} = 0$, entonces la sucesión de estimadores $\{\ell' \hat{\beta}_n\}$ es un estimador consistente en media cuadrática (y simple) de $\{\ell' \beta_n\}$ (donde ℓ es un vector $p \times 1$ de constantes).
2. La sucesión de estimadores $\{\hat{\sigma}_n^2\}$ es un estimador consistente en media cuadrática (y simple) de σ^2 .

Prueba 3.1 Hay que demostrar que $\hat{\sigma}_n^2$ y $\ell' \hat{\beta}_n$ son insesgados en el límite y que las varianzas de $\hat{\sigma}_n^2$ y $\ell' \hat{\beta}_n$ son cero en el límite. Se dejan los detalles como ejercicio.

Corolario 3.3 Bajo las condiciones del teorema 3.16, cada elemento de $\hat{\beta}_n$, es un estimador consistente en media cuadrática (y simple) del correspondiente elemento en β .

3.1.5 Estimación Puntual de funciones lineales de los parámetros

Para un vector \mathbf{x} en el dominio \mathbb{D} , se sigue que

$$\mu(\mathbf{x}) = \beta' \mathbf{x} = \sum_{i=1}^p \beta_i x_i$$

y esta es una combinación lineal de los β_i . Por lo tanto por las partes 6 y 7 del teorema 3.13, y más específicamente por el teorema 3.15, se puede hacer $t(\beta, \sigma^2)$ igual a

$\beta'x$, y se sigue que el estimador insesgado de varianza uniformemente mínima de $\mu(x)$ es $\hat{\mu}(x) = \hat{\beta}'x$ para cualquier vector x en el dominio \mathbb{D} , donde $\hat{\beta}$ está dado como en la parte 1 del teorema 3.13.

Para cualquier vector de constantes ℓ $p \times 1$, el EIVUM de $\ell'\beta$, cualquier combinación lineal de β , está dado por $\ell'\hat{\beta}$. Esto es,

$$\widehat{\ell'\beta} = \ell'\hat{\beta}$$

Por supuesto, $\mu(x)$ es un caso especial de $\ell'\beta$, como se puntualizará en el siguiente ejemplo

■ EJEMPLO 3.5

Considere el ejemplo ??, donde $\mu(x) = \beta_0 + \beta_1 x$ para x en el intervalo $50 \leq x \leq 300$ millas. Si se recolectan los datos y se usan las formulas del ejemplo 3.1, se puede estimar β_0 y β_1 . Suponga que se desea estimar la media de la variable aleatoria $Y_{(60)}$, es decir, $\mu(60)$. Ya que $\mu(60) = \beta_0 + \beta_1$, el EIVUM de $\mu(60)$ es

$$\hat{\mu}(60) = \hat{\beta}_0 + 60\hat{\beta}_1$$

Por otro lado, suponga que el investigador quiere estimar $30\beta_0 - 80\beta_1$, lo cual se denota por $\ell'\beta$ donde $\ell' = [30, -80]$. De acuerdo con el modelo de la situación del mundo real, $30\beta_0 - 80\beta_1$ no es la media $\mu(x)$ de cualquiera de las distribuciones guajiras. Sin embargo, se puede estimar $\ell'\beta$, y por el teorema 6.2.2 el EIVUM es $\ell'\hat{\beta} = 30\hat{\beta}_0 - 80\hat{\beta}_1$.

3.2 Valores Ajustados y Residuales

Una vez obtenido el estimador $\hat{\beta}$ de β . Si sustituimos esos valores en la ecuación del modelo

$$E[\mathbf{y}] = \mathbf{X}\beta$$

obtenemos la ecuación de regresión estimada

$$E[\widehat{\mathbf{y}}] = \mathbf{X}\hat{\beta} \quad (3.33)$$

que por costumbre usamos $\hat{\mathbf{y}}$ en vez de $E[\widehat{\mathbf{y}}]$.

Cuando sustituimos los valores observados en la ecuación 3.33, obtenemos el **vector de valores ajustados**, representado por $\hat{\mathbf{y}}$. Otra manera de representar el vector de valores ajustados es

$$\hat{\mathbf{y}} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{H}\mathbf{y} \quad (3.34)$$

donde

$$\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \quad (3.35)$$

La diferencia entre el vector de valores observados \mathbf{y} y el vector de valores ajustados $\hat{\mathbf{y}}$ es conocido como el **vector de residuales** o simplemente **residuales**, y se representa con la letra e , es decir

$$e = y - \hat{y} \quad (3.36)$$

el cual se puede escribir como

$$e = y - Hy = (I - H)y \quad (3.37)$$

3.3 Evaluación del Modelo

La regresión lineal calcula una ecuación que minimiza la distancia entre la recta ajustada y todos los puntos de los datos. Técnicamente, el método de mínimos cuadrados minimiza la suma de cuadrados de los residuos. En general, un modelo se ajusta bien a los datos si las diferencias entre los valores observados y los valores predichos (ajustados) son pequeños e insesgados.

Antes de mirar las medidas estadísticas para evaluar la bondad del ajuste, se deben chequear los gráficos de los residuos. Los mismos pueden revelar patrones que indican resultados sesgados de manera más eficaz que los números. Cuando los gráficos de residuos pasan la prueba, se deben usar números y comprobar las estadísticas de bondad de ajuste.

3.3.1 Partición de las Sumas de Cuadrados

La variación de un conjunto de datos es convencionalmente medida en términos de las desviaciones de los y_i alrededor de la media \bar{y} . Estas desviaciones se muestran por la línea vertical en la figura 3.3.

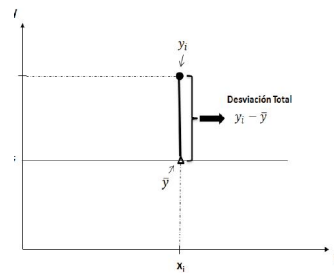


Figure 3.3 Desviación Total

La medida de la variación total, denotada por SCT , es la suma de las desviaciones al cuadrado, es decir

$$SCT = \sum_{i=1}^n (y_i - \bar{y})^2 \quad (3.38)$$

Aquí SCT denota la *suma de cuadrados total*. Si todas las observaciones y_i son las mismas, $SCT = 0$. Cuanto mayor es la variación entre las observaciones y_i , mayor es SCT . Por lo tanto, SCT es la variación de las observaciones y_i cuando no se toman en cuenta las variables explicatorias.

Cuando utilizamos la variable predictora x , la variación que refleja la incertidumbre con respecto a la variable y es la de las observaciones y_i alrededor de la regresión ajustada:

$$y_i - \hat{y}_i$$

Estas desviaciones se muestran en la línea vertical en la figura 3.4.

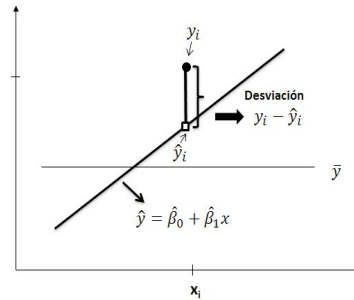


Figure 3.4 Desviación de las observaciones con respecto a la regresión

La medida de la variación en las observaciones y_i que está presente cuando la variable predictora x se tiene en cuenta es la suma de las desviaciones $y_i - \hat{y}_i$ al cuadrado, es decir

$$SCE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (3.39)$$

Aquí SCE denota la *suma de cuadrado del error*. Si todas las observaciones y_i caen sobre la recta de regresión ajustada, $SCE = 0$. Cuanto mayor sea la variación de las observaciones y_i alrededor de la recta de regresión ajustada, mayor será la SCE .

La diferencia entre estas dos medida también es una suma de cuadrados, la cual está dada por

$$SCR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \quad (3.40)$$

Note que la SCR es una suma de desviaciones al cuadrado, las desviaciones son

$$\hat{y}_i - \bar{y}$$

Estas desviaciones se muestran por la línea vertical en la figura 3.5. Cada desviación es simplemente la diferencia entre el valor ajustado de la recta de regresión y la media de las observaciones (la cual es la misma media de los valores ajustados). Si la recta de regresión es horizontal de manera que $\hat{y}_i - \bar{y} = 0$, entonces $SCR = 0$, de lo contrario, SCR será positiva.

La SCR puede considerarse como una medida de la parte de la variabilidad de los y_i la cual está asociada con la recta de regresión. Cuanto mayor sea SCR en relación a la SCT , mayor será el efecto de la relación de regresión para explicar la variación total de las observaciones y_i .

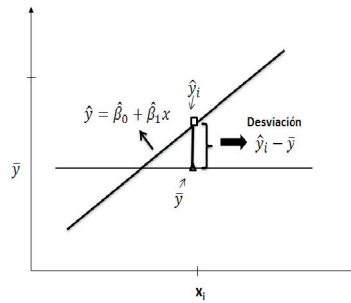


Figure 3.5 Desviación de la regresión con respecto a la media

3.3.1.1 Desarrollo Formal de la Partición La desviación total, usada en la medida de la variación total de las observaciones $y_i - \bar{y}$, usada en la medida de la variación total de las observaciones y_i sin tomar en cuenta la variable predictora, se puede descomponer en dos componentes:

$$y_i - \bar{y} = \hat{y}_i - \bar{y} + y_i - \hat{y}_i \tag{3.41}$$

Los dos componentes son:

1. La desviación del valor ajustado \hat{y}_i alrededor de la media \bar{y} , $(\hat{y}_i - \bar{y})$
2. La desviación de la observación y_i alrededor de la recta de regresión ajustada, $(y_i - \hat{y}_i)$

Es una propiedad remarcable que la suma de cuadrados de estas desviaciones cuadrados tienen la misma relación:

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2 \tag{3.42}$$

o, usando la notación en (3.44), (3.39) y (3.40):

$$SCT = SCR + SCE \tag{3.43}$$

Para probar este resultado procedemos de la siguiente manera:

$$\begin{aligned} \sum_{i=1}^n (y_i - \bar{y})^2 &= \sum_{i=1}^n [(\hat{y}_i - \bar{y}) + (y_i - \hat{y}_i)]^2 \\ &= \sum_{i=1}^n [(\hat{y}_i - \bar{y})^2 + (y_i - \hat{y}_i)^2 + 2(\hat{y}_i - \bar{y})(y_i - \hat{y}_i)] \\ &= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2 + 2 \sum_{i=1}^n (\hat{y}_i - \bar{y})(y_i - \hat{y}_i) \end{aligned}$$

Se puede probar fácilmente que el tercer termino del lado derecho (la suma cruzadas es igual a cero).

3.3.1.2 Sumas de cuadrados en términos matriciales En la ecuación (3.44) tenemos una expresión de la SCT , pero la misma se puede escribir como

$$SCT = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n y_i^2 - \frac{(\sum_{i=1}^n y_i)^2}{n} \quad (3.44)$$

Pero $\sum_{i=1}^n y_i^2$ se puede escribir en términos matricial como $\mathbf{y}'\mathbf{y}$. Además

$$\frac{(\sum_{i=1}^n y_i)^2}{n} = \left(\frac{1}{n}\right) \mathbf{y}'\mathbf{J}\mathbf{y}$$

donde \mathbf{J} es una matriz de unos. Por lo tanto

$$SCT = \mathbf{y}'\mathbf{y} - \left(\frac{1}{n}\right) \mathbf{y}'\mathbf{J}\mathbf{y} = \mathbf{y}' \left[\mathbf{I} - \left(\frac{1}{n}\right) \mathbf{J} \right] \mathbf{y} \quad (3.45)$$

Así mismo, la SCE esta dada por la ecuación (3.39). Si denotamos $e_i = y_i - \hat{y}_i$ entonces $SCE = \sum_{i=1}^n e_i^2$, lo cual en terminos matriciales sería

$$SCE = \mathbf{e}'\mathbf{e} = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \quad (3.46)$$

que también se puede escribir como

$$SCE = \mathbf{e}'\mathbf{e} = \mathbf{y}'\mathbf{y} - \hat{\boldsymbol{\beta}}' \mathbf{X}'\mathbf{y} \quad (3.47)$$

Finalmente, la suma de cuadrados de regresión se puede escribir como

$$SCR = \hat{\boldsymbol{\beta}}' \mathbf{X}'\mathbf{y} - \left(\frac{1}{n}\right) \mathbf{y}'\mathbf{J}\mathbf{y} \quad (3.48)$$

3.4 Coeficiente de Determinación

El coeficiente de determinación, denotado por R^2 es una medida estadística de cuan cerca de los datos está la ecuación ajustada. La definición de R^2 es muy simple, este es el porcentaje de la variación de la variable respuesta que es explicada por el modelo lineal, es decir,

$$R^2 = \frac{\text{Variación Explicada}}{\text{Variación Total}}$$

De acuerdo con lo estudiado anteriormente la variación explicada está dada por la *suma de cuadrados de regresión*, $SCR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$, y la variación total está dada por la *suma de cuadrados total*, $SCT = \sum_{i=1}^n (y_i - \bar{y})^2$. Por lo tanto,

$$R^2 = \frac{SCR}{SCT} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (3.49)$$

Forma matricial: En la sección anterior se llegó a las siguientes expresiones de las sumas de cuadrados en forma matricial

- $SCT = \mathbf{y}' \left[\mathbf{I} - \left(\frac{1}{n}\right) \mathbf{J} \right] \mathbf{y}$
- $SCE = \mathbf{e}'\mathbf{e} = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$

$$\blacksquare SCR = \hat{\beta}' X' y - \left(\frac{1}{n}\right) y' J y$$

Por lo tanto el coeficiente de terminación en forma matricial está dado por

$$R^2 = \frac{SCR}{SCT} = \frac{\hat{\beta}' X' y - \left(\frac{1}{n}\right) y' J y}{y' \left[\mathbf{I} - \left(\frac{1}{n}\right) J \right] y} = \frac{\hat{\beta}' X' y - n\bar{y}^2}{y' y - n\bar{y}^2} \quad (3.50)$$

3.4.1 Propiedades del coeficiente de determinación

Algunas propiedades de R^2 se listan a continuación.

1. El rango de R^2 es $0 \leq R^2 \leq 1$. Si todos los $\hat{\beta}_j$ fuesen cero, excepto $\hat{\beta}_0$, R^2 debería ser cero. Si todos los valores de y caen sobre la superficie ajustada, es decir $y_i = \hat{y}_i$, $i = 1, 2, \dots, n$, entonces R^2 debería ser 1.
2. Adicionar una variable x al modelo incrementa (nunca decrementa) el valor de R^2 .
3. Si $\hat{\beta}_1 = \hat{\beta}_2 = \dots = \hat{\beta}_k = 0$, entonces

$$E(R^2) = \frac{k}{n-1} \quad (3.51)$$

Note que los $\hat{\beta}_j$ no serían 0 cuando los β_j son 0.

4. R^2 es invariante a transformaciones lineales de rango completo sobre los x y a un cambio de escala sobre y (pero no invariante a transformaciones lineales conjuntas de y y las x).

En las propiedades 2 y 3 vemos que si k es una fracción relativamente grande de n , es posible tener un valor grande de R^2 que no es significativo. En este caso, los x que no contribuyen a predecir y pueden aparecer para hacerlo en un ejemplo particular, y la ecuación de regresión estimada puede no ser un estimador útil del modelo poblacional. Para corregir esta tendencia, un R^2 ajustado, denotado por R_a^2 fue propuesto por Ezekiel (1930), el cual veremos más adelante.

3.4.2 Observaciones sobre el coeficiente de determinación

1. No permite determinar si los coeficientes estimados y las predicciones son sesgadas.
2. No indica si un modelo de regresión es adecuado. Se puede tener un valor bajo de R^2 para un modelo bueno, o un valor alto de R^2 para un modelo que no se ajusta a los datos.
 - (a) R^2 **bajos no necesariamente son malos.** En algunos campos, es completamente esperado que el valor de R^2 sea bajo. Por ejemplo, al intentar predecir el comportamiento humano, tal como en psicología, típicamente se tienen valores de R^2 por debajo del 50%. Los humanos simplemente son más difíciles de predecir que, digamos, los procesos físicos.

Si el valor de R^2 es bajo pero se tienen predictores estadísticamente significativos, se pueden sacar conclusiones importantes sobre como cambios en los valores de los predictores están asociados con cambios en la variable respuesta. Independientemente del R^2 , los coeficientes significativos aún representan el cambio medio

en la respuesta por el cambio en una unidad en el predictor mientras los otros predictores se mantienen constantes. Obviamente, este tipo de información puede ser extremadamente valiosa.

Por lo tanto, no debe sorprendernos tener resultados con R^2 bajos y p – *valor* bajos.

- (b) R^2 **altos no necesariamente son buenos**. Un R^2 alto no necesariamente indica que el modelo tiene un buen ajuste, esto puede ser una sorpresa. Pero existen situaciones en las que se puede observar gráficamente que la línea ajustada (modelo de regresión lineal simple) tiene un comportamiento alejado de los datos y aún así el R^2 puede ser alto. Esto es posible por muchas razones, una de ellas es que la mejor manera de modelar la variable respuesta es a través de regresión no lineal; otra es que en la construcción del modelo no se están usando los predictores importantes o que se han incluido muchos predictores. Existen otras razones que se dejan como tema de investigación a los estudiantes.
3. Una pregunta común al realizar un análisis de regresión es ¿Qué tan alto debe ser el R^2 ? La respuesta a esta pregunta es... **depende**.. Depende del objetivo principal que se tenga al realizar el análisis. Dos objetivos frecuentes en el análisis de regresión son
- (a) Describir la relación entre los predictores y la variable respuesta, o
- (b) Predecir la variable respuesta

Si el principal objetivo es determinar cuales predictores son estadísticamente significativos y como cambios en los predictores se relacionan a cambios en la variable respuesta (objetivo (a)), entonces el valor de R^2 es irrelevante. Por ejemplo, suponga que se modela la relación entre X y Y , se encuentra que el p –valor para X es significativo, además que su coeficiente es 2 y que los supuestos se cumplen, estos resultados indican que un incremento de una unidad en X esta asociado a un incremento promedio de 2 unidades en Y . Esta interpretación es correcta aún si el R^2 es 0.95 o 0.25. Por lo tanto preguntarse ¿cuán grande debe ser el R^2 ? en este contexto no tiene sentido porque este no es relevante. Un R^2 bajo ni niega un predictor significativo, ni cambia el sentido de su coeficiente. El R^2 es simplemente cualquier valor, y no necesita ser un valor particular para permitir una interpretación válida.

Por el contrario, si el objetivo principal es producir predicciones precisas, el valor de R^2 es importante. Las predicciones no son tan simples como solo predecir un valor, ya que ellas incluyen un margen de error, predicciones más precisas tienen menos error. El R^2 entra en escena porque un R^2 bajo indica que el modelo tiene mucho error. Por lo tanto, un R^2 puede advertir de predicciones imprecisas. Sin embargo, no se puede usar el R^2 para determinar si las predicciones son lo suficientemente precisas para las necesidades.

3.4.3 Coeficiente de determinación ajustado

Hemos visto que una de las propiedades que tiene el R^2 es que este aumenta (o por lo menos no disminuye) a medida que se van añadiendo variables al modelo. En consecuencia un modelo con más coeficientes puede parecer tener un mejor ajuste simplemente porque tiene más términos, situación que no necesariamente es cierta.

El coeficiente de determinación ajustado es una versión modificada del R^2 que se ajusta por el número de predictores en el modelo. El R^2 ajustado se incrementa solo si el nuevo término mejora al modelo más lo que se esperaba por azar. Este disminuye cuando un predictor mejora al modelo por menos de lo esperado al azar. Siempre es más bajo que el R^2 .

Para obtener R_a^2 , primero restamos $k/(n-1)$ en 3.51 desde R^2 para corregir por el sesgo cuando $\beta_1 = \beta_2 = \dots = \beta_k = 0$. Esta corrección, sin embargo, debería hacer R_a^2 más pequeña cuando los $\hat{\beta}$ son grandes, por lo que una modificación adicional se hace para que $R_a^2 = 1$ cuando $R^2 = 1$. Por lo tanto R_a^2 es definido como

$$R_a^2 = \frac{(R^2 - \frac{k}{n-1})(n-1)}{n-k-1} = \frac{(n-1)R^2 - k}{n-k-1} \quad (3.52)$$