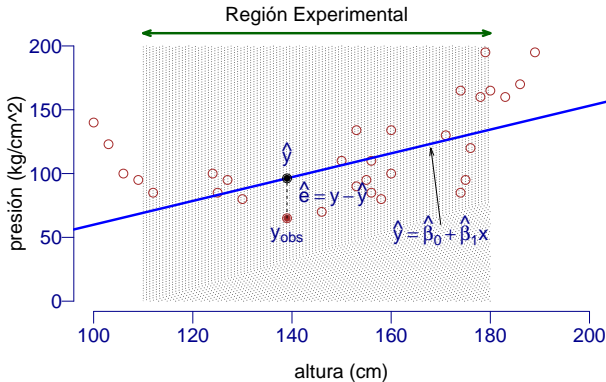


2 Regresión Lineal Simple



Contenido

1 Introducción

- Modelo y supuestos

2 Estimación de parámetros

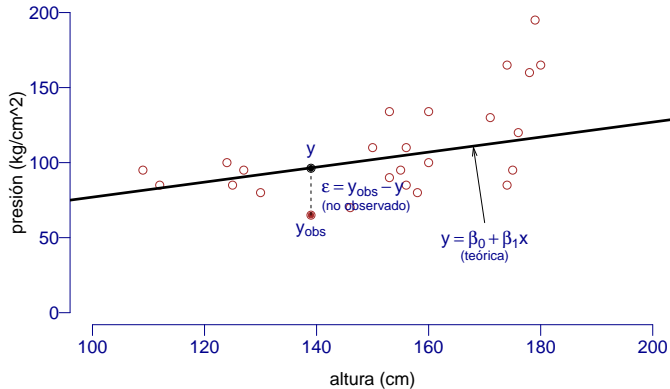
- Estimadores de mínimos cuadrados
- Ejemplo
- Propiedades

3 Inferencia

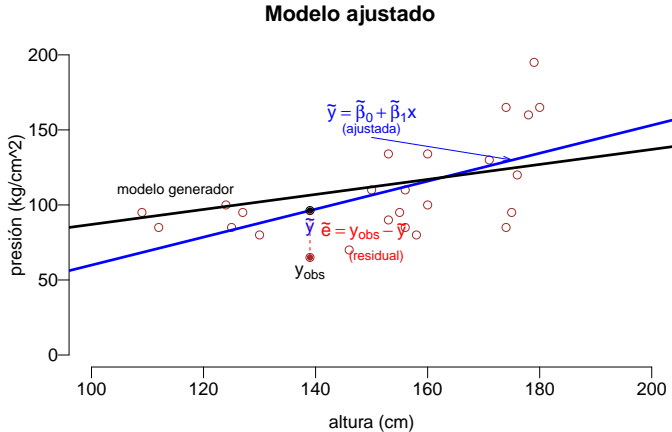
- Distribuciones
- Pruebas de hipótesis
- Regiones de Confianza
- Análisis de varianza
- Coeficiente de determinación
- Regresión por el origen

Modelo generador teórico

Modelo generador (teórico)



Modelo ajustado



El problema de mínimos cuadrados

Considere la *suma de cuadrados*: $S(\beta) \doteq \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$

Criterio

$$\min_{\beta} S(\beta) \equiv \min_{\beta_0, \beta_1} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

Estimadores Mínimos Cuadrados (EMC)

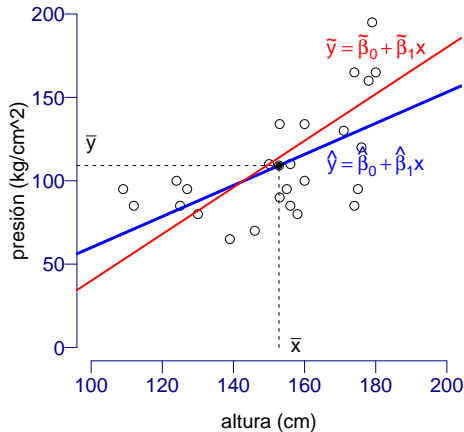
$$\frac{dS}{d\beta} \equiv 0 \Rightarrow \left\{ \begin{array}{l} \text{Ecuaciones normales (ortogonales) :} \\ n\beta_0 + \beta_1 \sum x_i = \sum y_i \\ \beta_0 \sum x_i + \beta_1 \sum x_i^2 = \sum x_i y_i \end{array} \right.$$

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Recta ajustada por mínimos cuadrados

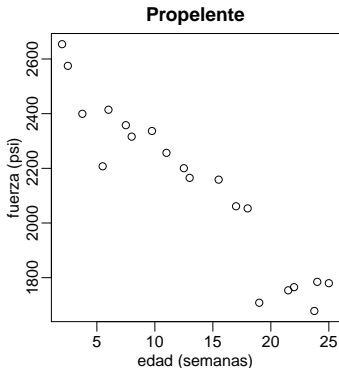
Recta de mínimos cuadrados



Ejemplo propelente¹

Datos Se considera un motor de cohete estudiando el propelente de encendido dentro de un depósito de metal. La fuerza para separar la unión entre las componentes del combustible es la respuesta, que depende de la edad de propelente.

obs	fuerza	edad
1	2158.70	15.50
2	1678.15	23.75
3	2316.00	8.00
4	2061.30	17.00
5	2207.50	5.50
6	1708.30	19.00
7	1784.70	24.00
8	2575.00	2.50
9	2357.90	7.50
10	2256.70	11.00
11	2165.20	13.00
12	2399.55	3.75
13	1779.80	25.00
14	2336.75	9.75
15	1765.30	22.00
16	2053.50	18.00
17	2414.40	6.00
18	2200.50	12.50
19	2654.20	2.00
20	1753.70	21.50



¹Montgomery et al. (2001)

Ejemplo propelente (cont.)

Ajuste

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	2627.822	44.184	59.48	< 2e-16
edad	-37.154	2.889	-12.86	1.64e-10

Residual standard error: 96.11 on 18 degrees of freedom

Multiple R-squared: 0.9018, Adjusted R-squared: 0.8964

F-statistic: 165.4 on 1 and 18 DF, p-value: 1.643e-10

Analysis of Variance Table

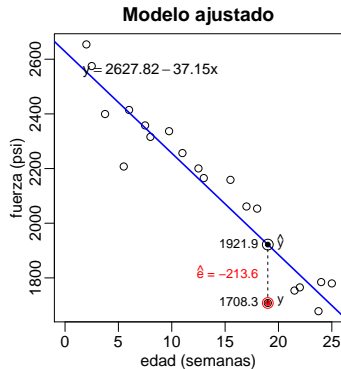
Response: fuerza

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
edad	1	1527483	1527483	165.38	1.643e-10
Residuals	18	166255	9236		

Ejemplo propolente (cont.)

Datos observados, ajustados y residuales

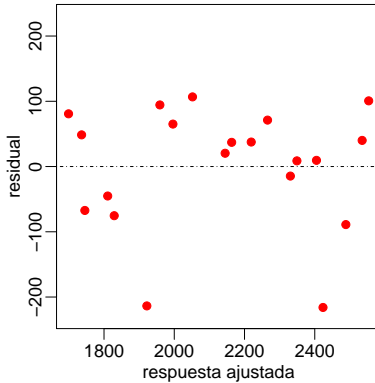
obs	yobs	yhat	res
1	2158.70	2051.9	106.8
2	1678.15	1745.4	-67.3
3	2316.00	2330.6	-14.6
4	2061.30	1996.2	65.1
5	2207.50	2423.5	-216.0
6	1708.30	1921.9	-213.6
7	1784.70	1736.1	48.6
8	2575.00	2534.9	40.1
9	2357.90	2349.2	8.7
10	2256.70	2219.1	37.6
11	2165.20	2144.8	20.4
12	2399.55	2488.5	-88.9
13	1779.80	1699.0	80.8
14	2336.75	2265.6	71.2
15	1765.30	1810.4	-45.1
16	2053.50	1959.1	94.4
17	2414.40	2404.9	9.5
18	2200.50	2163.4	37.1
19	2654.20	2553.5	100.7
20	1753.70	1829.0	-75.3



Ejemplo propolente (cont.)

Gráfica de residuales

Residuales vs. respuesta ajustada



Propiedades de los estimadores de mínimos cuadrados

$\sum(x_i - \bar{x}) = 0$ y por lo mismo $S_{xy} = \sum(x_i - \bar{x})y_i$. Luego $\hat{\beta}_1 = \sum c_i y_i$, donde $c_i = (x_i - \bar{x})/S_{xx}$.

$$\hat{\beta}_1 = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\mathbb{E}[\hat{\beta}_1] = \beta_1$$

$$\mathbb{E}[\hat{\beta}_0] = \beta_0$$

$$\text{var}(\hat{\beta}_1) = \sigma^2 \frac{1}{S_{xx}}$$

$$\text{var}(\hat{\beta}_0) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)$$

Teorema Gauss-Markov

Bajo los supuestos

$$\mathbb{E}[\epsilon_i] = 0, \quad \text{var}(\epsilon_i) = \sigma^2, \quad \text{cov}(\epsilon_i, \epsilon_j) = 0$$

los *estimadores de mínimos cuadrados* son los *mejores* estimadores lineales insesgados, en el sentido de que tienen varianza mínima.

el modelo de regresión:

$$\begin{aligned} L(\beta_0, \beta_1, \sigma^2; \mathbf{x}, \mathbf{y}) &= \prod_{i=1}^n f_{\epsilon}(x_i, y_i; \beta_0, \beta_1, \sigma^2) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2}(y_i - \beta_0 - \beta_1 x_i)^2\right\} \\ &= (2\pi)^{-n/2} \sigma^{-n} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2\right\} \\ \ell = \log L &= -\frac{n}{2} \log(2\pi) - n \log(\sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \end{aligned}$$

Pruebas de Hipótesis sobre los coeficientes

Suponiendo σ^2 conocida

- Recta que pasa por β_0^0 al origen:

$$H_0 : \beta_0 = \beta_0^0 \text{ vs. } H_1 : \beta_0 \neq \beta_0^0$$

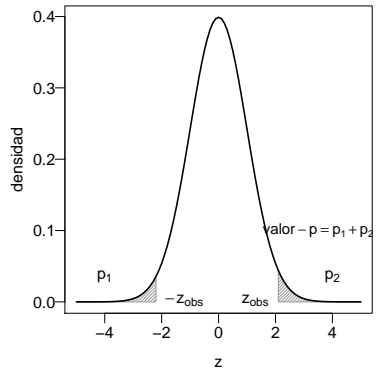
$$z_0 = \frac{\hat{\beta}_0 - \beta_0^0}{\text{de}(\hat{\beta}_0)} = \frac{\hat{\beta}_0 - \beta_0^0}{\sigma \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}} \sim N(0, 1)$$

- Recta de pendiente β_1^0 :

$$H_0 : \beta_1 = \beta_1^0 \text{ vs. } H_1 : \beta_1 \neq \beta_1^0$$

$$z_1 = \frac{\hat{\beta}_1 - \beta_1^0}{\text{de}(\hat{\beta}_1)} = \frac{\hat{\beta}_1 - \beta_1^0}{\sigma \sqrt{\frac{1}{S_{xx}}}} \sim N(0, 1)$$

Distribución nula



Pruebas de Hipótesis sobre los coeficientes

Suponiendo σ^2 desconocida

- Recta que pasa por β_0^0 al origen:

$$H_0 : \beta_0 = \beta_0^0 \text{ vs. } H_1 : \beta_0 \neq \beta_0^0$$

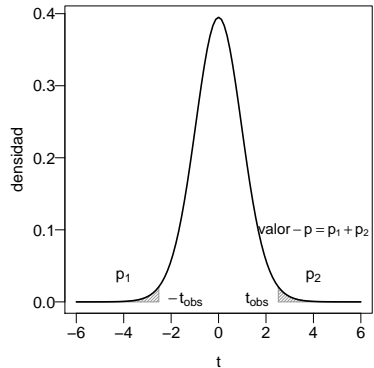
$$t_0 = \frac{\hat{\beta}_0 - \beta_0^0}{\text{ee}(\hat{\beta}_0)} = \frac{\hat{\beta}_0 - \beta_0^0}{s\sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}} \sim t_{n-2}$$

- Recta de pendiente β_1^0 :

$$H_0 : \beta_1 = \beta_1^0 \text{ vs. } H_1 : \beta_1 \neq \beta_1^0$$

$$t_1 = \frac{\hat{\beta}_1 - \beta_1^0}{\text{ee}(\hat{\beta}_1)} = \frac{\hat{\beta}_1 - \beta_1^0}{s\sqrt{\frac{1}{S_{xx}}}} \sim t_{n-2}$$

Distribución nula

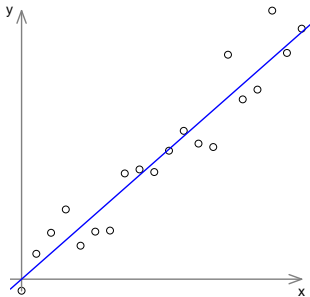


Pruebas de Hipótesis sobre los coeficientes

Casos particulares

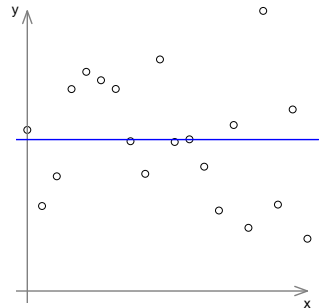
Recta que pasa por el origen:

$$H_0 : \beta_0 = 0 \quad \text{vs.} \quad H_1 : \beta_0 \neq 0$$



Significancia de la regresión:

$$H_0 : \beta_1 = 0 \quad \text{vs.} \quad H_1 : \beta_1 \neq 0$$



Este gráfico ilustra la regresión lineal con sus respectivas bandas de confianza y predicción. Se muestran los datos observados (puntos blancos) y la recta ajustada (línea azul). Las bandas de confianza (líneas marrones) representan la incertidumbre en la estimación de la media de la variable dependiente, mientras que las bandas de predicción (líneas rojas) representan la incertidumbre en la predicción de una nueva observación. El eje horizontal está etiquetado como \bar{x} y el eje vertical como $y_{\bar{x}}$. La ecuación de la recta ajustada se muestra como $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$ y se la denomina "Recta Ajustada".

Regiones de Confianza para $\beta = (\beta_0, \beta_1)$

Considere la reparametrización del modelo de regresión

$$y = \beta_0 + \beta_1 x + \epsilon = \beta'_0 + \beta_1(x - \bar{x}) + \epsilon$$

Entonces

$$\left. \begin{aligned} \hat{\beta}'_0 &= \bar{y} \quad y \quad \text{var}(\hat{\beta}'_0) = \sigma^2/n \\ \left(\frac{\hat{\beta}'_0 - \beta'_0}{\sigma/\sqrt{n}} \right)^2 &= \frac{n(\hat{\beta}'_0 - \beta'_0)^2}{\sigma^2} \sim \chi^2_1 \\ \left(\frac{\hat{\beta}_1 - \beta_1}{\sigma/\sqrt{S_{xx}}} \right)^2 &= \frac{S_{xx}(\hat{\beta}_1 - \beta_1)^2}{\sigma^2} \sim \chi^2_1 \end{aligned} \right\} \text{Independientes}$$

Luego,

$$\frac{n(\hat{\beta}'_0 - \beta'_0)^2}{\sigma^2} + \frac{S_{xx}(\hat{\beta}_1 - \beta_1)^2}{\sigma^2} \sim \chi^2_2$$

Y por otro lado,

$$(n-2)s^2/\sigma^2 \sim \chi^2_{n-2}$$

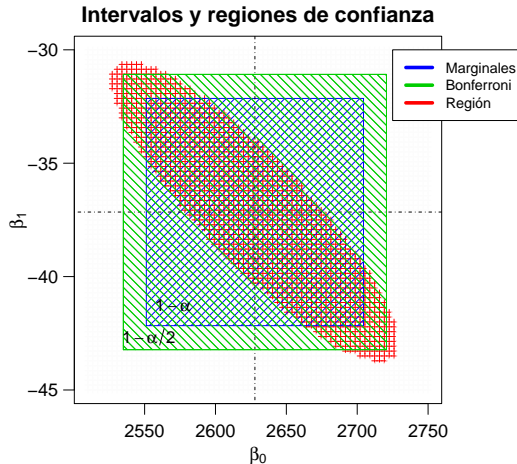
independientemente de $\hat{\beta}'_0$ y $\hat{\beta}_1$. Entonces,

$$\frac{n(\hat{\beta}'_0 - \beta'_0)^2 + S_{xx}(\hat{\beta}_1 - \beta_1)^2}{2s^2} \sim F_{2, n-2}$$

Por lo que una región del $(1 - \alpha)$ nivel de confianza para $\beta = (\beta_0, \beta_1)$ es:

$$\mathcal{RC}_{(1-\alpha)} = \left\{ \beta \in \mathbb{R}^2 : n(\hat{\beta}_0 - \beta_0)^2 + 2(\hat{\beta}_0 - \beta_0)(\hat{\beta}_1 - \beta_1) \sum x_i + (\hat{\beta}_1 - \beta_1)^2 \sum x_i^2 \leq 2s^2 F_{(1-\alpha; 2, n-2)} \right\}$$

Regiones de Confianza para $\beta = (\beta_0, \beta_1)$



Ejemplo propylene (Montgomery et al)

Regresión Simple por el Origen

El modelo de una línea recta que pasa por el origen es:

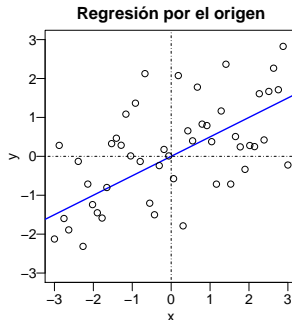
$$y_i = \beta x_i + \epsilon_i,$$

tiene los siguientes estimadores (insesgados) de mínimos cuadrados:

$$\hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2} = \frac{S_{xy}}{S_{xx}}$$

y

$$s^2 = \frac{1}{n-1} \sum (y_i - \hat{y}_i)^2$$



Los estimadores $\hat{\beta}$ y s^2 tiene propiedades similares a las correspondientes en los modelos con ordenada al origen. A saber,

$$\hat{\beta} \sim N(\beta, \sigma^2 \frac{1}{S_{xx}}) \quad \text{y} \quad s^2 \sim \frac{\sigma^2}{n-1} \chi_{n-1}^2$$

Recuerde que en este caso la R^2 no es fácilmente interpretable.