

IECD-325: Regresión lineal simple

Felipe Osorio

felipe.osorio@uv.cl

Suponga el **modelo de regresión**

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad i = 1, \dots, n,$$

donde y es la variable dependiente (respuesta) y x es la variable regresora o explicativa.

Los parámetros β_0 y β_1 corresponden al **intercepto** y la **pendiente**, respectivamente, mientras que $\{\epsilon_i\}$ representa disturbios aleatorios.

Para propósitos de inferencia estadística, supondremos

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

para $i, j = 1, \dots, n$ ($i \neq j$).

De este modo, dispondremos de $(x_1, y_1), \dots, (x_n, y_n)$ observaciones desde (x, y) .

Regresión lineal (simple)

Deseamos hallar β_0 y β_1 tal que produzcan el **mejor ajuste** a los datos¹. En este curso usaremos el **método de mínimos cuadrados** ordinarios, dado por

$$\min_{\theta} Q(\beta),$$

con $\beta = (\beta_0, \beta_1)^\top$ y

$$Q(\beta) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 = \sum_{i=1}^n \epsilon_i^2$$

Observación:

La función $Q(\beta)$ es conocida como **suma de cuadrados de los errores**.

¹Es decir, deseamos obtener estimadores para α y β

Regresión lineal (simple)

Usando el **método de mínimos cuadrados** obtenemos:

$$\begin{aligned}\hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x}, \\ \hat{\beta}_1 &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}.\end{aligned}$$

La **recta de regresión** es dada por:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i, \quad i = 1, \dots, n,$$

y llamamos a \hat{y}_i es **valor predicho** (o valor ajustado). Además,

$$e_i = y_i - \hat{y}_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i, \quad i = 1, \dots, n,$$

es conocido como el i -ésimo **residuo**.

Regresión lineal (simple)

Una **medida de variabilidad** es dada por:

$$s^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2.$$

Mientras que

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

se denomina **coeficiente de determinación**.²

Interpretación:

R^2 es la varianza de los datos **que puede ser explicada** por el modelo.

²permite medir la calidad (bondad) del ajuste

Regresión lineal (simple)

Es posible notar que (cuando el modelo tiene intercepto):

$$R^2 = 1 - \frac{RSS}{s_{\text{DATOS}}^2},$$

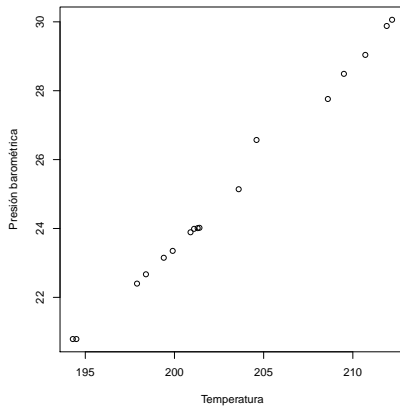
con

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2, \quad s_{\text{DATOS}}^2 = \sum_{i=1}^n (y_i - \bar{y})^2.$$

Interpretación:

En efecto, $0 \leq R^2 \leq 1$ permite medir la **calidad** (o bondad) **del ajuste**.

Presión barométrica en pulgadas de mercurio y temperatura de ebullición del agua en grados Fahrenheit para 17 diferentes altitudes.



Datos de Forbes

```
1 > library(MASS)
2 > data(forbes) # disponibiliza los datos en la sesión
3
4 > forbes
5      bp  pres
6 1  194.5 20.79
7 2  194.3 20.79
8 3  197.9 22.40
9 4  198.4 22.67
10 5  199.4 23.15
11 6  199.9 23.35
12 7  200.9 23.89
13 8  201.1 23.99
14 9  201.4 24.02
15 10 201.3 24.01
16 11 203.6 25.14
17 12 204.6 26.57
18 13 209.5 28.49
19 14 208.6 27.76
20 15 210.7 29.04
21 16 211.9 29.88
22 17 212.2 30.06
23
```

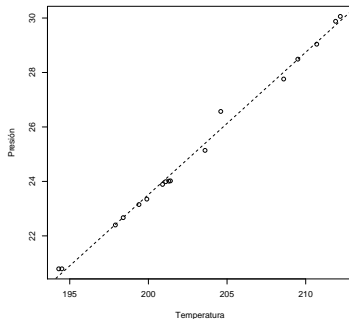

Datos de Forbes

```
1 # ajuste de un modelo de regresión
2 > fm <- lm(pres ~ bp, data = forbes)
3
4 # salida:
5 > fm
6
7 Call:
8 lm(formula = pres ~ bp, data = forbes)
9
10 Coefficients:
11 (Intercept)          bp
12   -81.0637       0.5229
13
14 # residuos y valores ajustados
15 > res <- residuals(fm)
16 > fit <- fitted(fm)
17
18 # otra forma de calcular R^2
19 > cor(fit, forbes$pres)^2
20 [1] 0.9944282
21
```

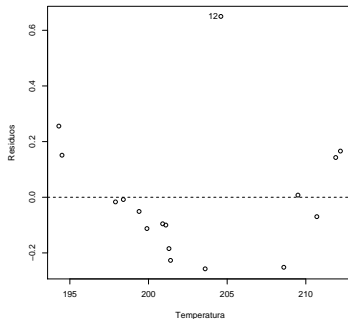
Datos de Forbes

```
1 # salida un poco más extensa
2 > summary(fm)
3
4 Call:
5 lm(formula = pres ~ bp, data = forbes)
6
7 Residuals:
8      Min       1Q   Median       3Q      Max
9 -0.25717 -0.11246 -0.05102  0.14283  0.64994
10
11 Coefficients:
12             Estimate Std. Error t value Pr(>|t|)
13 (Intercept) -81.06373    2.05182   -39.51  <2e-16 ***
14 bp           0.52289    0.01011    51.74  <2e-16 ***
15
16 Residual standard error: 0.2328 on 15 degrees of freedom
17 Multiple R-squared:  0.9944, Adjusted R-squared:  0.9941
18 F-statistic: 2677 on 1 and 15 DF, p-value: < 2.2e-16
19
```

Datos de Forbes



(a) recta ajustada



(b) residuos vs. ajuste

Ahora consideramos el modelo

$$100 \times \log_{10}(\text{Presión}_i) = \beta_0 + \beta_1 \text{Temperatura}_i + \epsilon_i,$$

para $i = 1, \dots, n$.

Se obtuvo (usando función `lm` de **R**)

$$\hat{\beta} = (-42.1378, 0.8955)^\top \quad \text{y} \quad s^2 = 0.1438$$

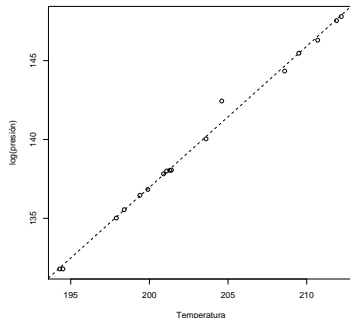
Además, $R^2 = 0.9950$.

Datos de Forbes

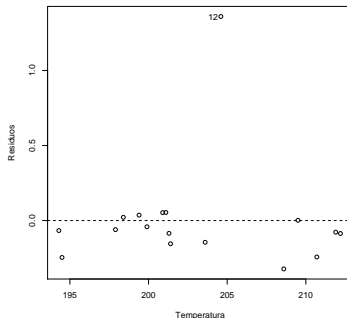
```
1 # modelo con datos transformados
2 > f1 <- lm(100 * log10(pres) ~ bp, data = forbes)
3
4 > summary(f1)
5
6 Call:
7 lm(formula = 100 * log10(pres) ~ bp, data = forbes)
8
9 Residuals:
10      Min       1Q   Median       3Q      Max
11 -0.31974 -0.14707 -0.06890  0.01877  1.35994
12
13 Coefficients:
14             Estimate Std. Error t value Pr(>|t|)
15 (Intercept) -42.16418    3.34136  -12.62 2.17e-09 ***
16 bp           0.89562    0.01646   54.42 < 2e-16 ***
17
18 Residual standard error: 0.3792 on 15 degrees of freedom
19 Multiple R-squared:  0.995, Adjusted R-squared:  0.9946
20 F-statistic: 2962 on 1 and 15 DF, p-value: < 2.2e-16
21
```

Datos de Forbes

Recta de regresión y gráfico de residuos para los datos de Forbes³.



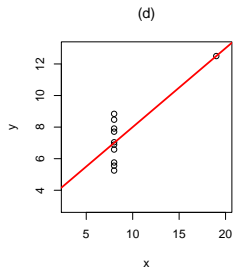
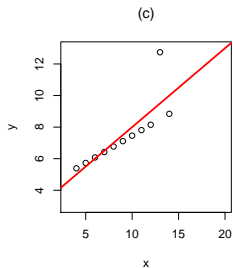
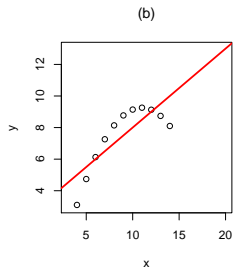
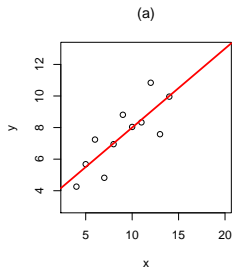
(a) recta ajustada



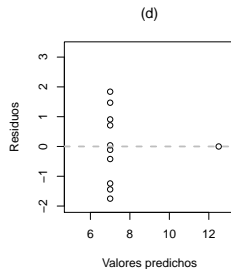
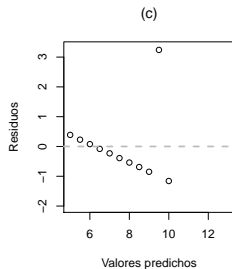
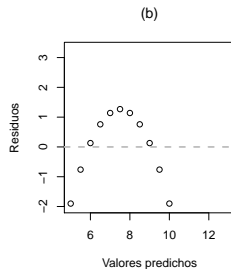
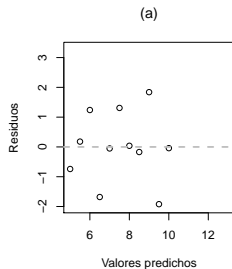
(b) residuos vs. ajuste

³datos transformados

Cuarteto de regresiones “idénticas” de Anscombe (1973)



Cuarteto de regresiones “idénticas” de Anscombe (1973)



Cuarteto de regresiones “idénticas” de Anscombe (1973)

Observaciones:

- ▶ Para el cuarteto de regresiones de Anscombe se obtiene (para todos los modelos):

$$\hat{\beta}_0 = 3.001, \quad \hat{\beta}_1 = 0.500, \quad s^2 = 1.528, \quad R^2 = 0.666, \quad F = 17.97, \quad p = 0.002$$

confiar **solamente** en medidas globales puede ser **engañoso**.

- ▶ En efecto, note que

$$\sum_{i=1}^n e_i = 0, \quad \sum_{i=1}^n e_i \hat{y}_i = 0,$$

de modo que el gráfico de dispersión de **residuos vs. valores predichos** no debería presentar algún **comportamiento sistemático**.

- ▶ Es recomendable realizar un **análisis de residuos** o de **diagnóstico**.