

Arrellar Estadística (The last One)

Prof: Rodrigo Assar (Sección 3) (Otoño-2008)

P.Aux: Andrés Iturriago, Víctor Riquelme

Tema: Regresión lineal.

Problema 1 El siguiente problema está relacionado a la producción de Biomasa de una planta (masa total de materia seca de una planta) y la radiación solar. A continuación, la tabla de datos:

Obs	Radiación solar (X)	Biomasa (Y)
1	28,7	16,6
2	68,4	48,1
3	120,7	121,7
4	217,2	219,6
5	313,5	374,5
6	419,1	570,8
7	535,9	648,2
8	641,5	758,6

a) Calcule $\hat{\beta}_0$ y $\hat{\beta}_1$ para la regresión lineal de Biomasa de una planta en función de la radiación solar (X). Escriba la ecuación de regresión e interprete.

sol la ecuación de regresión es

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2), \quad i=1, \dots, n \quad (n=8)$$

Por el método de mínimos cuadrados, encontramos $\hat{\beta}_0$ y $\hat{\beta}_1$:

$$\text{Sea } f(\beta_0, \beta_1) = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

$$\frac{\partial}{\partial \beta_0} f(\beta_0, \beta_1) = \sum_{i=1}^n 2(y_i - \beta_0 - \beta_1 x_i) \cdot (-1) = 0$$

$$\sum_{i=1}^n y_i - n\beta_0 - \beta_1 \sum_{i=1}^n x_i = 0 \Rightarrow \boxed{\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}}$$

$$\frac{\partial}{\partial \beta_1} f(\beta_0, \beta_1) = \sum_{i=1}^n 2(y_i - \beta_0 - \beta_1 x_i) \cdot (-x_i) = 0$$

$$\underbrace{\sum_{i=1}^n x_i y_i}_{xy} - \beta_0 \underbrace{\sum_{i=1}^n x_i}_{nx} - \beta_1 \underbrace{\sum_{i=1}^n x_i^2}_{xx} = 0$$

Tenemos que $xy - \hat{\beta}_0 x - \hat{\beta}_1 x^2 = 0$; $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ $\vec{e} = (1, -1, 1)$

$$xy - (\bar{y} - \hat{\beta}_1 \bar{x}) x - \hat{\beta}_1 x^2 = 0$$

$$xy - \bar{y} x + \hat{\beta}_1 \bar{x} x - \hat{\beta}_1 x^2 = 0$$

$$xy - \bar{y} x = \hat{\beta}_1 (x^2 - \bar{x} x)$$

$$\hat{\beta}_1 = \frac{xy - \bar{y} x}{x^2 - \bar{x} x} = \frac{\frac{1}{n} xy - \bar{y} \bar{x}}{\frac{1}{n} x^2 - \bar{x} \bar{x}} = \frac{\text{Cov}_{xy}}{S_x^2}$$

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

$$\hat{\beta}_0 = \bar{y} - \frac{\text{Cov}_{xy}}{S_x^2} \bar{x}$$

Entonces, con los datos dados, se tiene que

$$\hat{\beta}_1 = 1,269 \quad ; \quad \hat{\beta}_0 = -27,677$$

Si se aumenta 1 en Radiación, aumenta 1,269 en Biomasa.

(b) Construya intervalos de confianza de un 95% para β_0 y β_1 . Interprete los resultados.

Dem Recordemos que los estimadores son insesgados

$$E(\hat{\beta}_i) = \beta_i$$

Obs: Como $\varepsilon_i \sim N(0, \sigma^2)$, $\hat{\beta}_i \sim N(\beta_i, \text{Var}(\hat{\beta}_i))$ $(\text{Var}(\hat{\beta}) = \sigma^2 (X^T X)^{-1})$

$$\text{Var}(\hat{\beta}_0) = \frac{\sum x_i^2 \sigma^2}{n \sum x_i^2 - n^2 \bar{x}^2} \quad ; \quad \text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum x_i^2 - n \bar{x}^2}$$

$$\text{Var}(\hat{\beta}_0) = \frac{1}{n} \frac{\frac{1}{n} \sum x_i^2 \sigma^2}{S_x^2} \quad ; \quad \text{Var}(\hat{\beta}_1) = \frac{1}{n} \frac{\sigma^2}{S_x^2}$$

Entonces: $\hat{\beta}_i \sim N(\beta_i, \sigma^2 (X^T X)^{-1}_{ii})$

$$\Rightarrow \frac{\hat{\beta}_i - \beta_i}{\sigma \sqrt{(X^T X)^{-1}_{ii}}} \sim N(0,1)$$

Un estimador para σ^2 es (del método de max. verosimilitud)

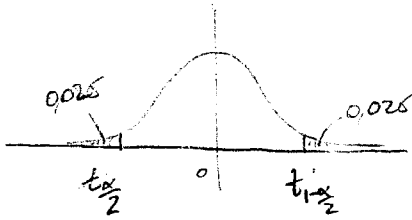
$$\hat{\sigma}^2 = \frac{1}{n} \sum_{j=1}^n (y_j - \hat{\beta}_0 - \hat{\beta}_1 x_j)^2 \quad (= \frac{1}{n} (Y - \hat{Y})^T (Y - \hat{Y}))$$

$$\boxed{\begin{array}{l} \text{Est. sesgado } \hat{\sigma}^2 \\ \hat{\sigma}^2 = \frac{1}{n-p-1} \hat{\sigma}^2 \text{ insesgado} \end{array}}$$

$$\frac{(n-2)\hat{\sigma}^2}{\sigma^2} = \frac{(n-2)}{n} \sum_{j=1}^n \frac{(y_j - \hat{\beta}_0 - \hat{\beta}_1 x_j)^2}{\sigma^2} \sim \chi_{n-2}^2$$

2

Entonces, $T = \frac{\hat{\beta}_i - \beta_i}{\frac{\hat{\sigma} \sqrt{(X^T X)^{-1}_{ii}}}{\sqrt{\frac{(n-2)\hat{\sigma}^2}{\sigma^2(n-2)}}}} = \frac{\hat{\beta}_i - \beta_i}{\hat{\sigma}_i} \sim t_{n-2}$ t-student de (n-2) grados de libertad



$$P\left(-t_{1-\frac{\alpha}{2}} < \frac{\hat{\beta}_i - \beta_i}{\hat{\sigma}_i} < t_{1-\frac{\alpha}{2}}\right) = 1-\alpha$$

$$P\left(-t_{1-\frac{\alpha}{2}} \hat{\sigma}_i < \hat{\beta}_i - \beta_i < t_{1-\frac{\alpha}{2}} \hat{\sigma}_i\right) = 1-\alpha$$

$$P\left(-t_{1-\frac{\alpha}{2}} \hat{\sigma}_i - \hat{\beta}_i < -\beta_i < t_{1-\frac{\alpha}{2}} \hat{\sigma}_i - \hat{\beta}_i\right) = 1-\alpha$$

$$P\left(\hat{\beta}_i - t_{1-\frac{\alpha}{2}} \hat{\sigma}_i < \beta_i < \hat{\beta}_i + t_{1-\frac{\alpha}{2}} \hat{\sigma}_i\right) = 1-\alpha$$

Para la t-student con $(n-2) = (8-2) = 6$ grados de libertad, tenemos que $(\alpha = 0.05)$
 $t_{1-0.025} = 2,447$

Ahora, $\hat{\sigma}_0 = \hat{\sigma} \sqrt{(X^T X)^{-1}_{00}} = \sqrt{\frac{(Y-\hat{Y})^T (Y-\hat{Y})}{n}} \sqrt{\frac{1}{n} \frac{\sum x_i^2}{S_x^2}} = 23,533 \cdot 0,61$

$$\hat{\sigma}_1 = \hat{\sigma} \sqrt{(X^T X)^{-1}_{11}} = \sqrt{\frac{(Y-\hat{Y})^T (Y-\hat{Y})}{n}} \sqrt{\frac{1}{n} \frac{1}{S_x^2}} = 23,533 \cdot 0,002$$

23,533

$$\hat{\sigma}_0 = 14,3 \quad ; \quad \hat{\sigma}_1 = 0,04. \quad //$$

los intervalos de confianza quedan

para β_0 : $(-27,67 - 2,447 \cdot 14,3; -27,67 + 2,447 \cdot 14,3) = (-62,668; 7,315)$

para β_1 : $(1,268 - 2,447 \cdot 0,04; 1,268 + 2,447 \cdot 0,04) = (1,171; 1,367)$

(c) Pruebe la significancia de los parámetros: considere

(i) * $H_0: \beta_1 = \beta_2 = 0$ vs H_1 : alguno no es nulo.

(ii) * $H_0: \beta_1 = 0$ vs $H_1: \beta_1 \neq 0$

(iii) * $H_0: \beta_2 = 0$ vs $H_1: \beta_2 \neq 0$.

Sol Primero, para contrastar $H_0: \beta_i = 0$ vs $\beta_i \neq 0$, usamos la t-student

$$\frac{\hat{\beta}_i - \beta_i}{\hat{\sigma}_i} \sim t_6 \quad ; \text{ bajo } H_0, \beta_i = 0. \quad \# \text{ grados libertad} = (n - (p+1))$$

$$(*) \quad \frac{\hat{\beta}_0}{\hat{\sigma}_0} = -1,935 \quad \leadsto \quad P(|t_6| > 1,935) > P(|t_6| > 2,447) = 0,05 \quad \begin{array}{l} \text{No rechazo} \\ \beta_0 = 0 \end{array}$$

$$(*) \quad \frac{\hat{\beta}_1}{\hat{\sigma}_1} = 31,725 \quad \leadsto \quad P(|t_6| > 31,725) < P(|t_6| > 2,447) = 0,05 \quad \begin{array}{l} \text{Rechazo} \\ \beta_1 = 0 \end{array}$$

Ahora, para el modelo global, comparamos los residuos ajustados por la regresión

$$F = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2 / p}{\sum_{i=1}^n \frac{(y_i - \hat{y}_i)^2}{E_i^2} / (n - (p+1))} \quad \begin{array}{l} \nearrow \text{Bajo } H_0 \\ \sim F_{(p, n-p-1)} \end{array} \quad \begin{array}{l} \rightarrow \text{varianza explicada / gl} \\ \rightarrow \text{varianza residual / gl} \end{array}$$

$$= \frac{R^2 / p}{(1-R^2) / (n - (p+1))}, \quad \text{con } R^2 = \frac{\sum (y_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} \quad \begin{array}{l} \text{coef. de} \\ \text{correlación múltiple.} \end{array}$$

$$F = 264,58$$

$$P(F_{2,6} > 5,14) = 0,05$$

$$\Rightarrow P(F_{2,6} > 264,58) \approx 0,00$$

Entonces, algún coeficiente es significativo ($\neq 0$).

$$\rightarrow \text{Var. Residual: } \frac{1}{n} \sum E_i^2 = \frac{1}{n} \sum (y_i - \hat{y}_i)^2 \rightarrow SSE \quad \begin{array}{l} (n - (p+1) = n - r) \\ SSTotal = SSE + SSR \end{array}$$

$$\rightarrow \text{Var. Explicada: } \frac{1}{n} \sum (\hat{y}_i - \bar{y})^2$$

$$\rightarrow \text{Var. Total: } \frac{1}{n} \sum (y_i - \bar{y})^2 \rightarrow SSR \quad (p+1)$$

Problema 2 El ministro de educación quiere estudiar de que depende el gasto anual en educación de un hogar, para ello, recolecta información en 100 hogares y plantea el modelo lineal $E(y) = \beta_0 + \beta_1 x^1 + \beta_2 x^2 + \beta_3 x^3 + \beta_4 x^4$ (1)
 con x^1 = ingreso del hogar, x^2 = # hijos; x^3 = talla del jefe de hogar, x^4 = # personas en casa.

(2) Complete los resultados de la regresión lineal (1) dados en las tablas 1 y 2. Interprete los resultados.

TABLA 1

Variable de	Estimación	Desviación típica	t-student	p-valor
de	20,387	20,384	= 1,000	0,319 → No Rechazo $\beta=0$ } X
Ingreso	0,189	0,021	= 9,242	0,000 → Rechazo $\beta=0$ } //
Nº hijos	17,378	2,978	= 5,836	0,000 → Rechazo $\beta=0$ } //
Talla jefe	8,869	6,176	= 1,436	0,154 → No Rechazo $\beta=0$ } X
Nº personas	0,187	0,107	1,749	0,083 → No Rechazo $\beta=0$ } X

TABLA 2

Fuente	G. libertad	Suma Cuadrados	F	p-valor
Regresión	4	129489,083	38,185	0,000 → Hay algún coeficiente significativo
Residuos	95	80539,635		
Total	99	210028,718		

$R^2 = 0,785$; $\hat{\sigma}_{marginal} = 28,12$ → El modelo es bueno (el R es alto)

Se pueden sacar del modelo x^3 y x^4 .

$$F = \frac{SC_{Reg} / GL_{Reg}}{SC_{Residuos} / GL_{Residuos}}$$

(ii) Se plantea un modelo con el ingreso y el # de hijos:

$$E(y) = \beta_0 + \beta_1 x^1 + \beta_2 x^2 \quad (2)$$

Se propone resolver el test H_0 : modelo (2) vs H_1 : modelo (1). Para ello se resuelve el modelo (2) obteniéndose las tablas (3) y (4). Conviene el cambio en la suma de los cuadrados de los residuos de los modelos (2) y (1).

Var	Est.	Des. Std	t-st	P-valor	Fuente	GL	SC	F	p-valor
de	81,035	8,575	6,301	0,000	Regresión	2	125292,857	71,71	0,000
Ingreso	0,196	0,019	10,019	0,000	Residuos	97	84735,867		
Nº hijos	17,804	2,969	5,995	0,000	Total	99	210028,718		

$$R^2 = 0,772, \hat{\sigma}_{marginal} = 28,56$$

Obj El cambio es pequeño entre los 2 este modelo explica más el gasto anual (en relación con el modelo anterior).

En resuma

(*) Test t-student: decide si el coeficiente β_j es significativo; para ello se plantea el test $\frac{\hat{\beta}_j}{\hat{\sigma}_j} \sim t_{n-p-1}$ bajo $H_0: \beta_j = 0$.

p-valor < 0.05 rechaza H_0 (el coef. es significativo en el modelo, es decir, la variable correspondiente debe ser considerada)

(*) Test Fisher: decide si el modelo es explicativo (si alguno de los β_j es distinto de 0)

$$\text{plantea el test } F = \frac{R^2/p}{(1-R^2)/(n-p-1)} = \frac{\text{suma cuadrados Regresión} / p}{\text{suma cuadrados Residuos} / (n-p-1)} \sim F_{p, n-p-1} \quad (\text{bajo } H_0: \text{ todos } \beta_j = 0)$$

p-valor < 0.05 rechaza H_0 (algunas variables explican a la variable Y)

(*) R^2 cercano a 1 \Rightarrow hay un buen ajuste del modelo hacia la variable explicada.

R^2 cercano a 0 \Rightarrow modelo malo.

(*) Modelo $Y = X\beta + \varepsilon$ con $\varepsilon \sim N_n(0, \sigma^2)$, σ^2 desconocida

$$\hat{\beta} = (X^t X)^{-1} X^t y \quad ; \quad \hat{y} = E(y) = X(X^t X)^{-1} X^t y$$

$$E(\hat{\beta}) = \beta \quad ; \quad \text{Var}(\hat{\beta}) = \sigma^2 (X^t X)^{-1}$$

$$\text{Est. Max. Vro: } \hat{\sigma}^2 = \frac{1}{n} (y - \hat{y})^t (y - \hat{y}) \quad \text{sesgado} \quad ; \quad \text{insesgado: } \tilde{\sigma}^2 = \frac{n}{n-p-1} \hat{\sigma}^2$$