

Proceso de modelización estadística

lunes, 5 de agosto de 2024 9:29

1er. Especificación del modelo:

- Qué variables intervienen, que supuestos se deben cumplir, forma funcional, cuáles son los parámetros a estimar.

2do. Estimación del modelo:

- Construir estimaciones para los parámetros a partir de una muestra

3er. Inferencia del modelo

- Pruebas de hipótesis sobre los parámetros y la bondad de ajuste del modelo.

4to. Validación del modelo

- Diagnóstico y comprobación de supuestos del modelo.

Especificación de un modelo de regresión lineal simple

lunes, 5 de agosto de 2024 9:38

Variables: Dos variables (simple)

- X = variable independiente (numéricas continuas)
- Y = variable dependiente (numéricas continuas)

Forma funcional:

$$Y = F(X) + E$$

Donde:

F(X) = Función de regresión

E = término de error

APLICACIÓN DE SUPUESTOS EN EL MODELO:

1. Linealidad: $F(X) = a + bX$

a = Intersección,

b = pendiente

Como a y b son desconocidos, la idea es estimarlos a partir de una muestra.

PARAMETROS	INTERPRETACION EN EL MODELO
β_0	Intersección en el eje Y
β_1	Pendiente de la recta

Forma funcional:

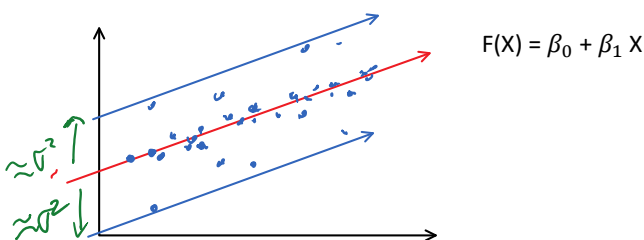
$$Y = F(X) + E$$

$$Y = \beta_0 + \beta_1 X + E$$

2. Homocedasticidad:

Los errores al ser v.a. pueden tomar cualquier valor, independientemente del valor de X, pero, deben tener una magnitud de variabilidad constante

Gráficamente, deben formar una "banda de ancho fijo" alrededor de la función de regresión



4. Independencia:

Los errores deben ser independientes entre sí.

Gráficamente, los errores no deben presentar ningún tipo de ordenamiento o patrón.

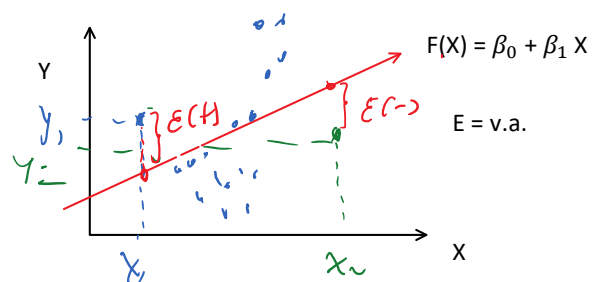
Supuestos estructurales:

1. Linealidad: F(X) es una ecuación de una recta
2. Homocedasticidad: $\text{Var}[E | X] = \text{cte} = \sigma^2$
3. Normalidad: E sigue una distribución $N(0, \sigma^2)$
4. Independencia: No debe existir relación entre un error y otro

NOTA: A partir de ahora vamos a SUPONER que los supuestos se cumplen. Pero, una vez estimado el modelo, el cumplimiento de estos supuestos se deben VERIFICAR!!!

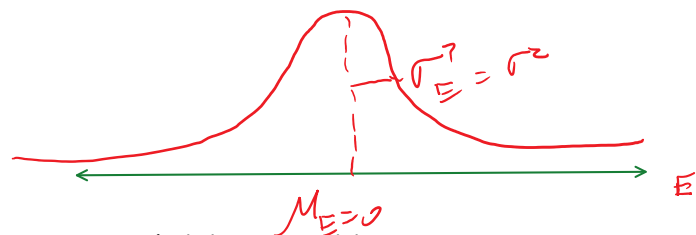
¿Qué representa el término de error E en el modelo?

- Si el término de error es CERO; los pares (x, y) estarían exactamente sobre la RECTA.
- Como asumimos que los errores no son cero, entonces los puntos pueden estar sobre (errores positivos) o bajo (errores negativos) la recta.
- Los errores pueden tomar cualquier valor (grande o pequeño)
- Entonces los Errores son una VARIABLE ALEATORIA!!



3. Normalidad:

Los errores son v.a. que siguen una distribución normal, con media cero y varianza homocedástica σ^2

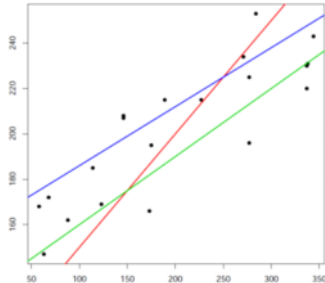


La mayoría de los errores deben ser cero o muy cercanos a cero, una minoría pueden ser muy positivo o negativos, pero sin salirse de la banda homocedástica!

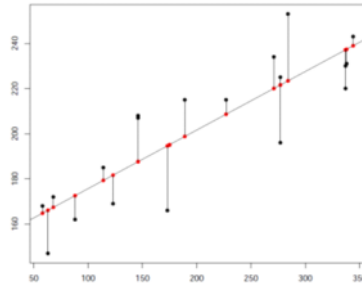
Además, se espera que la mitad de los puntos tengan errores positivos, y la otra mitad tengan errores negativos

En todo modelo de RLS, se deben estimar tres parámetros: $\beta_0, \beta_1, \sigma^2$

Para estimar los parámetros, se toma una muestra de n pares observados de (X, Y) y se aplica un método de estimación: MINIMOS CUADRADOS ORDINARIOS.



Cada recta tiene sus propios residuos



Los residuos pueden ser (+) o (-)

Medida de ajuste = residuos al cuadrado

$$\min \sum \hat{\varepsilon}^2$$

Problema de optimización numérica!!

OBJETIVO: Construir una recta que minimice la suma de residuos cuadrados!!

Para hallar la estimación de los parámetros, se debe derivar y luego igualar a cero, la primera derivada del criterio de optimización

Estimación por Mínimos Cuadrados

Consiste en seleccionar $\hat{\beta}_0$ y $\hat{\beta}_1$ tales que:

$$(\hat{\beta}_0, \hat{\beta}_1) = \arg \min_{(\beta_0, \beta_1)} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_i)^2$$

Estimadores

$$\hat{\beta}_0 = \bar{Y} - \frac{S_{xy}}{S_x^2} \bar{x} \quad \hat{\beta}_1 = \frac{S_{xy}}{S_x^2}$$

La recta de regresión estimada por mínimos cuadrados es la que pasa por el vector de medias o centro de gravedad, (\bar{x}, \bar{Y}) , y tiene pendiente $\hat{\beta}_1 = \frac{S_{xy}}{S_x^2}$.

Residuos del modelo RLS

Supongamos el modelo $Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$, con $i = 1, \dots, n$ y denotemos por $\hat{\beta}_0$ y $\hat{\beta}_1$ los estimadores de los parámetros. Los errores de predicción son:

$$\hat{\varepsilon}_i = Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i, \quad i = 1, \dots, n$$

también denominados **residuos** de la regresión.

Para un conjunto de datos (Xi, Yi) se pueden construir infinitos modelos (rectas)

OBJETIVO:

- Hallar el modelo (recta) que mejor se "ajuste" a los datos
- Es decir: que tenga los residuos más pequeños posibles.

CONSULTA!!:

Desarrollar las ecuaciones de las estimaciones de β_0 y β_1 por mínimos cuadrados!

Estimación de σ^2

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{\varepsilon}_i^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

Values	
Beta0	141.703353881689
Beta1	0.294259414639639
n	10
s2x	7282.01
Sxy	2142.8
vecmed	num [1:2] 225 208
x	num [1:10] 175 189 344 88 114 338 2.
y	num [1:10] 195 215 243 162 185 231 .

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

$$\hat{y} = 141.70 + 0.29 X$$

Interpretación:

Proyección tiempo de uso de máquina = $141.70 + 0.29$

*(unidades producidas en cada corrida de producción)

- $\hat{\beta}_0 = 141.70$ min (tiempo promedio de uso de máquina, cuando no se produce ninguna unidad, $X=0$) - costos fijos
- $\hat{\beta}_1 = 0.29$ min/unidad (tasa marginal unitaria). Si se produce una unidad adicional, el tiempo de uso de máquina AUMENTA en 0.29 minutos (relación directa)

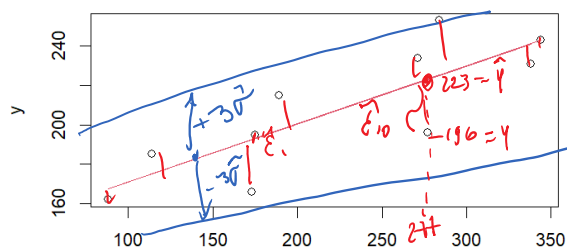
```
> Yproy = Beta0 + Beta1*x
> Yproy
[1] 193.1988 197.3184 242.9286 167.5982 175.2489 241.1630 221.4477 192.6102
[9] 225.2730 223.2132
```

```
> y
[1] 195 215 243 162 185 231 234 166 253 196
> resid = y - Yproy
> resid
[1] 1.80124856 17.68161675 0.07140748 -5.59818237 -9.75107285
[6] -10.16303603 12.55234475 -26.61023261 27.72697236 -27.21321174
```

$$\hat{\epsilon} = y - \hat{y}$$

$$\hat{\epsilon}_1 = 195 - 193.1988 = +1.80$$

$$\hat{\epsilon}_{10} = 196 - 223.21 = -27.213$$



```
> sum(resid) # practicamente CERO
[1] 2.842171e-14
```

$$\sum \hat{\epsilon} \approx 0 \rightarrow \frac{\sum \hat{\epsilon}}{n} \approx 0 \rightarrow \bar{\hat{\epsilon}} = 0$$

```
> SCR = sum(resid^2)
> SCR
[1] 2920.609
```

$$\sum \hat{\epsilon}^2 \leftarrow \min \left\{ \text{MCO} \right\}$$

```
> sigma2 = SCR / (n-2)
> sigma2
[1] 365.0762
```

$$\hat{\sigma}^2 = \frac{\sum \hat{\epsilon}^2}{n-2} = 365.0762$$

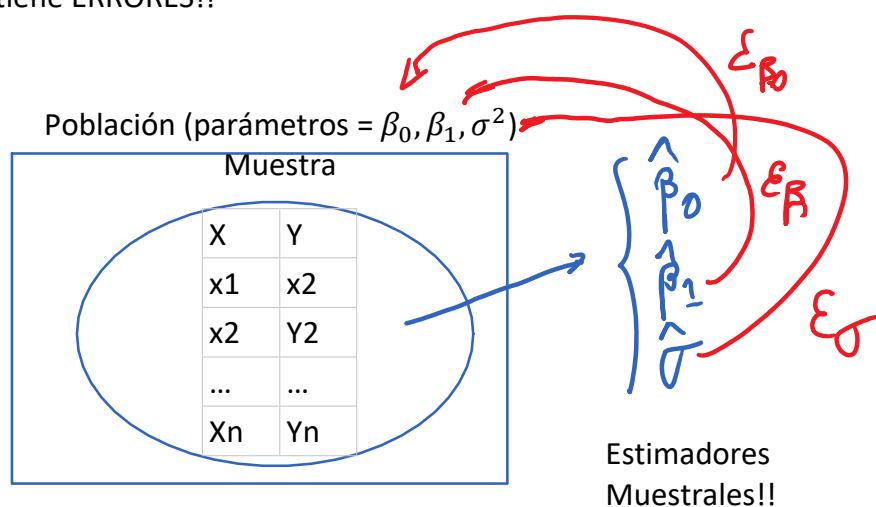
```
> sigma = sqrt(sigma2)
> sigma
[1] 19.10697
```

$$\hat{\sigma} = \sqrt{\hat{\sigma}^2} = 19.10$$

Inferencia del modelo: Pruebas T individuales

lunes, 12 de agosto de 2024 9:52

- Una cuestión es indicar cuánto valen las estimaciones del modelo, otra cosa es que esas estimaciones sean válidas. Porque estas estimaciones se obtuvieron a partir de una muestras, y toda muestra tiene ERRORES!!



$$\hat{\beta}_1 = \frac{S_{xy}}{S_x^2}$$

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

$$\hat{\sigma}^2 = \frac{\sum \epsilon^2}{n-2}$$

Parámetro	Estimador	Error muestral	Error típico	Estadístico Pivote
β_0	$\hat{\beta}_0$	$ \hat{\beta}_0 - \beta_0 $	$\hat{\sigma} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{nS_x^2}}$	$\frac{(\hat{\beta}_0 - \beta_0)}{\hat{\sigma} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{nS_x^2}}} = t_{(n-2)}$
β_1	$\hat{\beta}_1$	$ \hat{\beta}_1 - \beta_1 $	$\hat{\sigma} / (S_x \sqrt{n})$	$\frac{\hat{\beta}_1 - \beta_1}{\hat{\sigma} / (S_x \sqrt{n})} = t_{(n-2)}$
σ^2	$\hat{\sigma}^2$	$\hat{\sigma}^2 / \sigma^2$		$\chi^2 = (n-2) \hat{\sigma}^2 / \sigma^2$

Inferencia para el Beta Cero

```
> etB0 = sigma * sqrt((1/n) + ((mean(x)^2)/(n*S2x)))
> etB0 # estimacion del error tipico del beta cero
[1] 17.05838
```

```
> tB0 = Beta0 / etB0
> tB0 # estadístico de contraste de Beta 0
[1] 8.306967
```

$$t = \frac{\hat{\beta}_0 - \beta_0}{\hat{\sigma}_{\hat{\beta}_0}} = 8.30$$

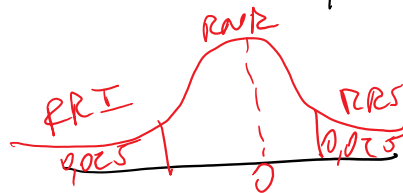
$H_0: \beta_0 = 0$
 $H_1: \beta_0 \neq 0$
 $gl(n-2) = 2$

P-Valor?

```
> PvalorB0 = 2*(1 - pt(tB0, n-2))
> PvalorB0 # es menor al 5% , H0 se rechaza
[1] 3.326772e-05
```

H0 se rechaza (Rechazamos que el Beta cero = 0)
 Si existen costos fijos en el modelo!!

Parámetro	Estimador	Error muestral	Error típico	Estadístico Pivote
β_0	$\hat{\beta}_0$	$ \hat{\beta}_0 - \beta_0 $	$\hat{\sigma} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{nS_x^2}}$	$\frac{(\hat{\beta}_0 - \beta_0)}{\hat{\sigma} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{nS_x^2}}} = t \quad (n-2)$



Inferencia para el Beta 1:

β_1	$\hat{\beta}_1$	$ \hat{\beta}_1 - \beta_1 $	$\hat{\sigma} / (S_x \sqrt{n})$	$\frac{\hat{\beta}_1 - \beta_1}{\hat{\sigma} / (S_x \sqrt{n})} = t \quad (n-2)$
		0.0708	0.0708	

```
> etB1 = sigma / sqrt(n*S2x)
> etB1
[1] 0.07080535
```

$$t = \frac{\hat{\beta}_1 - \beta_1}{\hat{\sigma}_{\hat{\beta}_1}} = \frac{0.2942}{0.0708} = 4.155$$

$H_0: \beta_1 = 0$
 $H_1: \beta_1 \neq 0$

P-Valor = 0.003 (< 5%)

```
[1] 4.155892
> PvalorB1 = 2*(1 - pt(tB1, n-2))
> PvalorB1 # es menor al 5% , H0 se rechaza
[1] 0.003183112
```

H0 se rechaza
 La tasa marginal unitaria de X es distinta a cero!

Prueba F Global

$H_0: Y = \beta_0 + \varepsilon$ para algún β_0

$H_a: Y = \beta_0 + \beta_1 X + \varepsilon$ para algún β_0 y algún β_1

Residuos de la regresión: $Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$

Residuos sin recta de regresión: $Y_i - \bar{Y}$

Descomposición:

$$Y_i - \bar{Y} = (Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) + (\hat{\beta}_0 + \hat{\beta}_1 x_i - \bar{Y})$$

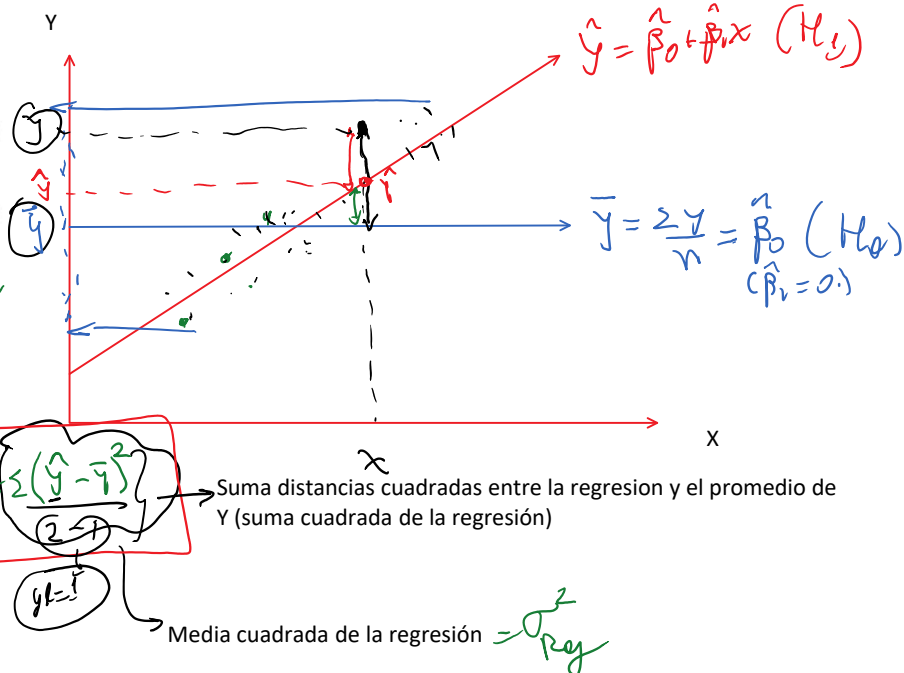
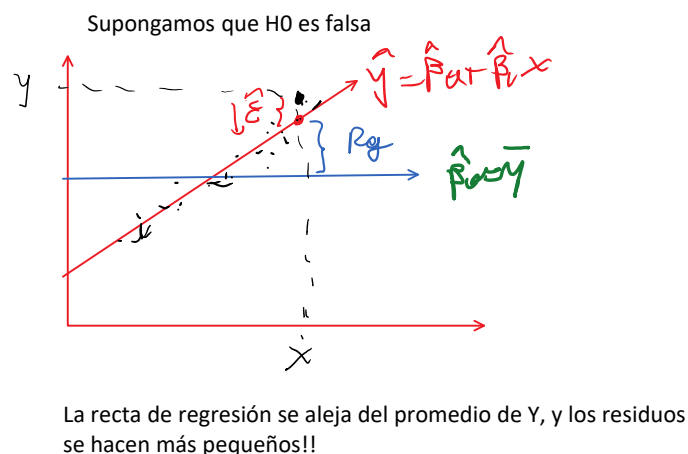
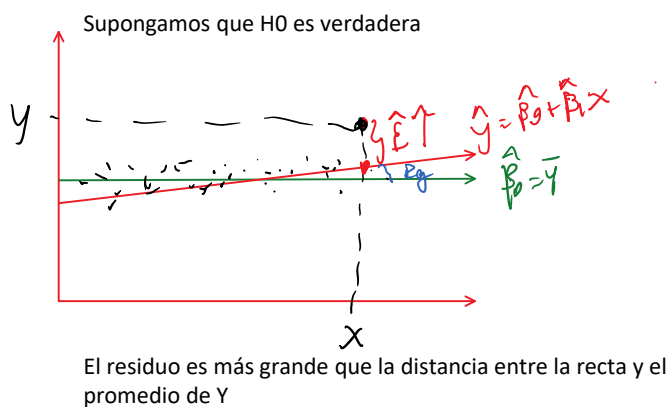


TABLA DE ANALISIS DE VARIANZAS DE LA REGRESION LINEAL

Fuente de variabilidad	Sumas cuadradas	Grados de libertad	Medias cuadradas	Estadístico F
Regresión	$\sum (\hat{Y} - \bar{Y})^2$	$2 - 1 = 1$	σ_{Reg}^2	
Error	$\sum (Y - \hat{Y})^2 = \sum \varepsilon^2$	$n - 2$	σ^2	
Total	$\sum (Y - \bar{Y})^2$	$n - 1$	S_y^2	



Luego, para que el modelo tenga sentido, se requiere que:

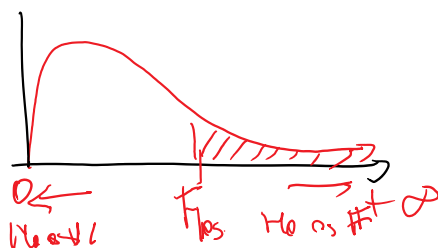
- Los residuos sea pequeños.
- La regresión se aleje del promedio de Y

A partir de esto, se construye un estadístico de contraste F:

- Si la H_0 es verdadera, F tiende a cero
- Si la H_0 es falsa, F tiende a + infinito

$$F = \frac{\sigma_{Reg}^2}{\sigma_{\varepsilon}^2} = \frac{\frac{\sum (\hat{Y} - \bar{Y})^2}{1}}{\frac{\sum \varepsilon^2}{n-2}}$$

$\leftarrow gl(numer)$
 $\leftarrow gl(denom)$



$$\left. \begin{array}{l} P\text{-Value} < \alpha \rightarrow H_0 \text{ es falsa} \\ P\text{-Value} \geq \alpha \rightarrow H_0 \text{ es verdadera} \end{array} \right\}$$

```
> SCT = sum((y - mean(y))^2) # suma cuadrada total/global
> SCT
[1] 9226
```

$$\rightarrow SCT = \sum (y - \bar{y})^2 = 9226$$

```
> SCReg = sum((Yproy - mean(y))^2) # suma cuadrada de la regresion
> SCReg
[1] 6305.391
```

$$\rightarrow SCReg = \sum (\hat{y} - \bar{y})^2 = 6305.391$$

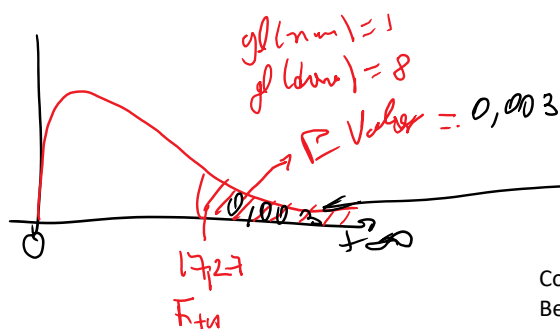
```
> SCR = sum((y - Yproy)^2) # suma cuadrada del residuo
> SCR
[1] 2920.609
```

$$\rightarrow SCE = \sum (y - \hat{y})^2 = \sum \hat{\epsilon}^2 = 2920.609$$

```
> SCReg + SCR #esta suma debe ser igual a la SCT
[1] 9226
```

```
> # PRUEBA GLOBAL F (contraste F)
> Ftest = (SCReg/1) / (SCR/(n-2)) # estadistico F
> Ftest
[1] 17.27144
```

$$F_{\text{test}} = \frac{\frac{SCReg}{1}}{\frac{SCR}{n-2}} = \frac{\frac{6305.391}{1}}{\frac{2920.609}{8}} = 17.27144$$



```
> PvalorF = 1 - pf(Ftest, 1, n-2)
> PvalorF # Este valor es menor al 5%, por tanto la H0 se rechaza
[1] 0.003183112
> # Beta1 es distinto de cero, luego Y si depende de X
```

Como el P-valor es menor al 5%, la H0 se rechaza!!
Beta 1 es distinto de cero (Y depende de X)