

Residual estandarizado

También conocidos como residuales estudentizados de manera interna, se define como:

$$d_i = \frac{e_i}{\hat{\sigma} \sqrt{1 - h_{ii}}}$$

La media de un residual estandarizado es 0 y su varianza es 1. Además, asumiendo normalidad, valores por encima de 2, o por debajo de -2 son considerados outliers.

```
modelo |> resid() -> r
summary(modelo)$sigma -> s
modelo |> model.matrix() -> X
X %*% solve(t(X) %*% X) %*% t(X) -> H ← matriz hat
(r/(s*sqrt(1-diag(H)))) -> res_stand
res_stand |> round(2)
```

1	2	3	4	5	6	7	8	9	10	11	12
-0.68	-0.42	-0.09	-0.58	-0.38	-0.30	-0.12	3.11	-0.40	-0.46	-0.08	0.10

Nota: se redondea solo con la finalidad de mostrar los residuales en la diapositiva.

Residual estandarizado

También conocidos como residuales estudentizados de manera interna, se define como:

$$d_i = \frac{e_i}{\hat{\sigma} \sqrt{1 - h_{ii}}}$$

Residual estudentizado

Sin embargo, de haber outliers, los residuales estandarizados podrían no ser normales. Ante ello, se propone redefinir la expresión para la estimación del CME, retirando la i-ésima observación. Así, el residual estudentizado (externamente) sería:

$$t_i = \frac{e_i}{\hat{\sigma}_{(i)} \sqrt{1 - h_{ii}}}$$

↪ (i) : al retirar la i-ésima observación

Leverage

Sabemos que:

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

h_{2T} (n × n) n × 1 n × 1

$$\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}_{n \times 1} \quad \hat{\mathbf{y}} = \begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{pmatrix}_{n \times 1}$$

Es decir:

$$\hat{y}_i = h_{i1}y_1 + h_{i2}y_2 + \dots + h_{ii}y_i + \dots + h_{in}y_n$$

$$\begin{bmatrix} h_{11} & h_{12} & \dots & h_{1n} \\ h_{21} & h_{22} & & h_{2n} \\ \vdots & \vdots & & \vdots \\ h_{n1} & h_{n2} & & h_{nn} \end{bmatrix} \times \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} h_{11}y_1 + h_{12}y_2 + \dots + h_{1n}y_n \\ \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \vdots \end{bmatrix}$$

\downarrow

$> \frac{2k}{n} ?$

- Si \Rightarrow Es leverage
- No \Rightarrow No es leverage



Distancia de Cook

Mide la influencia de cada observación en los coeficientes de regresión estimados

$$D_i = \frac{\underbrace{(\hat{\beta}_{(i)} - \hat{\beta})'}_{1 \times k} \underbrace{\mathbf{X}'\mathbf{X}}_{k \times k} \underbrace{(\hat{\beta}_{(i)} - \hat{\beta})}_{k \times 1}}{k\hat{\sigma}^2} = \frac{d_i^2 h_{ii}}{k(1 - h_{ii})}$$

$$\begin{array}{ll} \hat{\beta}_{k \times 1} & \hat{\beta}_{(i)}_{k \times 1} \\ X_{n \times k} & X'_{k \times n} \end{array}$$

Obs	D_i	> ?
1	D_1	
2	D_2	
3	D_3	
\vdots	\vdots	
n	D_n	

$$D_i = \frac{(\hat{\beta}_{(i)} - \hat{\beta})' X' X (\hat{\beta}_{(i)} - \hat{\beta})}{k \hat{\sigma}^2} = \frac{d_i^2 h_{ii}}{k(1 - h_{ii})}$$

```
# ¿observación 1 es influyente?
modelo = lm(y~x)
modelo_sin = lm(y[-1]~x[-1])
modelo |> coef() -> beta
modelo_sin |> coef() -> beta_sin
modelo |> model.matrix() -> X
(modelo |> sigma())**2 -> cme
2 -> k
(t(beta-beta_sin)%*%t(X)%*%X)%*(beta-beta_sin))/(k*cme)
```

```
[1,] 0.04252736 = D1
```

```
modelo |> rstandard() -> rsta
modelo |> hatvalues() -> h
(rsta[1]**2*h[1])/(k*(1-h[1]))
```

```
1
0.04252736 = D1
```

```
modelo |> cooks.distance() |> round(2) |> as.vector()
```

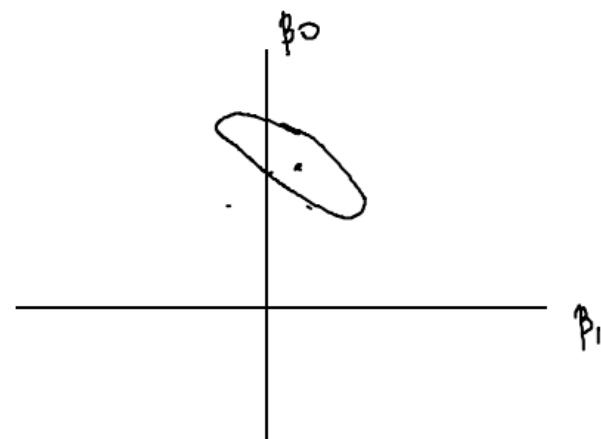
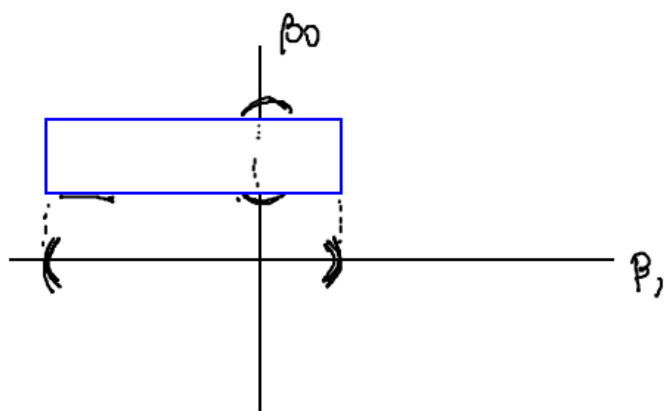
```
[1] 0.04 0.01 0.00 0.03 0.01 0.00 0.00 0.50 0.01 0.59 0.00 0.00
```

¿Entre qué valores fluctúa F? Interprete los valores de la siguiente tabla:

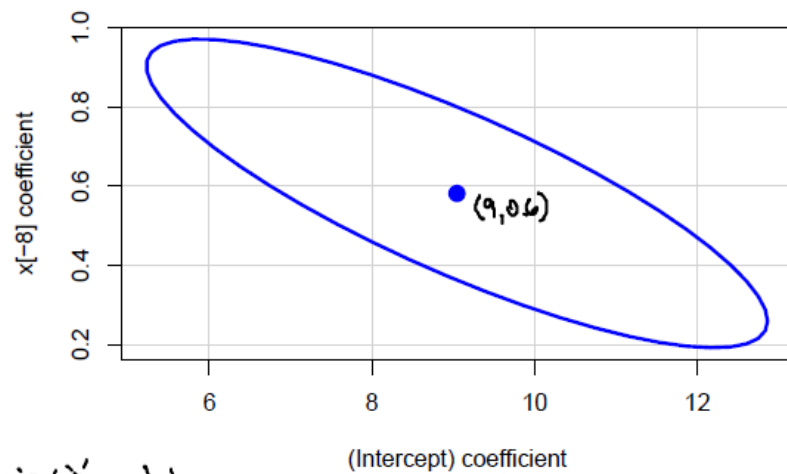
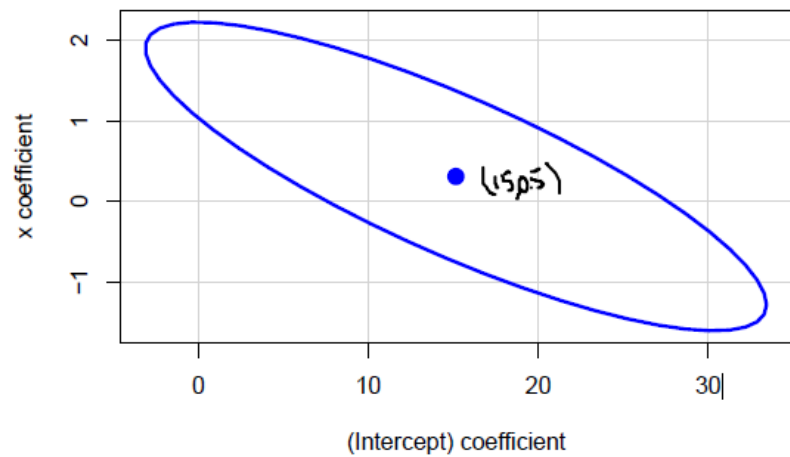
```
alfas = c(0.01,0.05,0.1,0.25,0.5,0.66,0.75,0.90,0.95,0.99)
data.frame(alfas,valorF=qf(alfas,k,n-k))
```

	alfas	valorF	
1	0.01	0.01006044	Si $D_i > 0.01006 \Rightarrow$ al retirar la i -ésima obs, $\hat{\beta}$ se mueve 1%
2	0.05	0.05155730	
3	0.10	0.10647844	
4	0.25	0.29611921	
5	0.50	0.74349177	equilibrio
6	0.66	1.20403474	más flexibles
7	0.75	1.59753955	
8	0.90	2.92446596	
9	0.95	4.10282102	
10	0.99	7.55943216	

$(\quad) \beta_1$



$$\hat{y} = \underset{\uparrow}{\hat{\beta}_0} + \underset{\downarrow}{\hat{\beta}_1} x_1$$



Variaçión del centro y la variabilidad

