

4,15

## Eficacia en el Control de Infecciones Hospitalarias



Diego Mauricio Ballesteros Osorio

Valentin Jose Padilla Marimon

León Felipe Restrepo Perez

Santiago Molina Muñoz

Estadística II

Equipo 22

Escuela de Estadística, Universidad Nacional de Colombia

Medellín

2023

**Problema planteado:**

En un estudio a gran escala realizado en EE.UU sobre la eficacia en el control de infecciones hospitalarias se recogió información en 113 hospitales. Y se obtuvo una muestra aleatoria con 60 registros. Según la base de datos considere lo siguiente:

**Tabla 1. Descripción de variables regresoras y de respuesta.**

Variable	Descripción
<b>Y:</b> Riesgo de infección	Probabilidad promedio estimada de adquirir infección en el hospital (Porcentaje)
<b>X<sub>1</sub>:</b> Duración de la estadía	Duración promedio de la estadía de todos los pacientes en el hospital (Días)
<b>X<sub>2</sub>:</b> Rutina de cultivos	Razón del número de cultivos realizados en pacientes sin síntomas de infección hospitalaria, por cada 100
<b>X<sub>3</sub>:</b> Número de camas	Número promedio de camas en el hospital durante el periodo de estudio
<b>X<sub>4</sub>:</b> Censo promedio diario	Número promedio de pacientes en el hospital por día durante el periodo del estudio
<b>X<sub>5</sub>:</b> Número de enfermeras	Número promedio de enfermeras, equivalente a tiempo completo, durante el periodo de estudio

**Problemas a resolver:****Pregunta 1.**

18 pt

**1.1 Estimación del modelo de regresión lineal múltiple**

El modelo que se va a emplear para aproximar la relación entre la variable **Y** y las variables **X<sub>k</sub>** ; **k=1,2,3,4,5**. (Ver tabla 1). Está dado por:  $1 \leq i \leq 60$ .

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + \beta_4 X_{4i} + \beta_5 X_{5i} + \epsilon_i ; \epsilon_i \sim N(0, \sigma^2). \text{ iid}$$

Al hacer una estimación de los parámetros:  $\beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5$ . Se obtienen los parámetros estimados:  $B_j$ ;  $j = 0, 1, 2, 3, 4, 5$ . Expuestos en la siguiente tabla:

**Tabla 2. Coeficientes de regresión estimados de  $B_j$**

	Coeficientes de regresión estimados
$B_0$	-1.0186398
$B_1$	0.2313154
$B_2$	0.0226572
$B_3$	0.0643261
$B_4$	0.0053465
$B_5$	0.0018688

Con los datos dados por la **tabla 2**. La regresión estimada o modelo de regresión ajustado será de la forma:

$$\hat{Y}_i = -1.0186398 + 0.2313154 X_{1i} + 0.0226572 X_{2i} + 0.0643261 X_{3i} + 0.0053465 X_{4i} + 0.0018688 X_{5i}$$

## 1.2 Significancia de la regresión

Para plantear la significancia de la regresión se formula la siguiente prueba de hipótesis:

$$H_0: \beta_1 = \beta_2 = \beta_3 = \beta_4 = \beta_5 = 0, \quad H_1: \text{Algún } \beta_k \neq 0, \quad k = 1, 2, 3, 4, 5.$$

Para continuar con el análisis de la significancia de la regresión. Es necesario calcular la tabla de análisis de varianza que estará compuesta por distintos elementos que son: Fuente de variación, suma de cuadrados, grados de libertad, cuadrados medios y finalmente un estadístico de prueba con el que se probarán las hipótesis planteadas anteriormente. La forma de la tabla queda de la siguiente manera:

Tabla 3. Análisis de varianza

Fuente de variación	Suma de cuadrados	Grados de libertad	Cuadrados medios	Valor estadístico de prueba F
Regresión	67.1743	5	13.434854	$F_0 = 14.2182$
Error	51.0251	54	0.944909	
Total	118.1994	59		

Usando los datos dados por la tabla de análisis de varianza. (Ver tabla 3). El estadístico de prueba está dado por:

$$F_0 = \frac{\text{CUADRADO MEDIO DE LA REGRESIÓN}}{\text{CUADRADO MEDIO DEL ERROR}} = \frac{\left( \frac{\text{SUMA DE CUADRADOS DE LA REGRESIÓN (SCR)}}{\text{GRADOS DE LIBERTAD ASOCIADO A SCR}} \right)}{\left( \frac{\text{SUMA DE CUADRADOS DEL ERROR (SCE)}}{\text{GRADOS DE LIBERTAD ASOCIADO A SCE}} \right)} ; F_0 \sim f_{5,54}$$

Es decir,

$$F_0 = \frac{13.434854}{0.944909} = \frac{\frac{67.1743}{5}}{\frac{51.0251}{54}} = 14.2182$$

Se rechaza  $H_0$  con un nivel de significancia  $\alpha = 0.05$  si:  $|F_0| > f_{\alpha,5,54}$  o  $VP < \alpha$  donde  $VP = P(f_{5,54} > |F_0|)$ . Entonces se tiene que:  $(f_{\alpha,5,54} = 2.38607) \Rightarrow |F_0| = |14.2182| > 2.38607$  y  $VP = 7.10924e-09 < 0.05$ .

Con los resultados dados tanto por la región de rechazo como por el VP, concluimos el rechazo de la hipótesis nula  $H_0$  y decimos que podría haber al menos un  $\beta_k \neq 0$  en la regresión que la hace significativa.

### 1.3 Significancia individual de parámetros

Para probar la significancia individual de los parámetros y confirmar a su vez el rechazo de la hipótesis nula del subpunto 1.2 para los parámetros  $\beta_k$ . Debemos plantear nuevas pruebas de

hipótesis pero esta vez para cada parámetro de la regresión estimada o ajustada  $\beta_j$ , es decir, incluiremos el intercepto en este caso. Las pruebas de hipótesis tienen la forma:

$$H_0: \beta_j = 0, H_1: \beta_j \neq 0; j = 0, 1, 2, 3, 4, 5.$$

El estadístico de prueba adquiere la forma:

$T_{0,j}$

$$T_j = B_j / ee(B_j); T_j \sim t_{54}$$

ojo con la notación,  $\beta$ , no  $B$

Para calcular cada estadístico de prueba  $T_j$  según corresponda usamos las tablas 2 y 4.

**Tabla 4. Errores estándar (ee) y valores P de  $B_j$**

	Errores estándar	Valores P
$B_0$	1.4895450	0.497
$B_1$	0.0769191	0.004
$B_2$	0.0270241	0.405
$B_3$	0.0151851	8.91e-05
$B_4$	0.0074696	0.477
$B_5$	0.0007184	0.012

Así cada estadístico de prueba es:

$$T_0 = -0.684, T_1 = 3.007, T_2 = 0.838, T_3 = 4.236, T_4 = 0.716, T_5 = 2.602; T_j \sim t_{54}$$

Se rechaza  $H_0$  con un nivel de significancia  $\alpha = 0.05$  cuando: La región de rechazo de  $H_0$  está dada por  $RR = \{T_j | |T_j| > t_{\alpha/2, 54}\}$  ó  $VP < \alpha$  cuando  $VP = P(t_{54} > |T_j|) < \alpha$ . ( $t_{\alpha/2, 54} = 2.004879$ ).

Podemos observar que tanto por el criterio de valor P (Ver tabla 4) como de región de rechazo (RR), solo los parámetros  $B_1, B_3, B_5$  son significativos, es decir, solo en estos casos rechazamos la hipótesis  $H_0$  que nos dice que son iguales a 0, es decir, no significativos. Mientras que  $B_0, B_2$  y  $B_4$  no son significativos.

Sólo se interpretaban las significativas

6

#### 1.4 Interpretación de parámetros estimados *1pt*

- B*
- ~~$B_0$~~ : Este parámetro no tiene interpretación debido a que ninguna de las variables regresoras tomadas en cuenta son 0 simultáneamente en la muestra, es decir, el "punto" (0, 0, 0, 0, 0) no es considerado en la muestra.
  - ~~$B_1$~~ : Este parámetro nos dice que la probabilidad porcentual promedio estimada de adquirir infección en el hospital aumenta a razón de 0.2313154 unidades por cada unidad de aumento de duración de días promedio de la estadía de todos los pacientes en el hospital. *(cuando las demás se quedan constantes)*
  - ~~$B_2$~~ : Este parámetro nos dice que la probabilidad porcentual promedio estimada de adquirir infección en el hospital aumenta a razón de 0.0226572 unidades por cada unidad de aumento en la razón del número de cultivos realizados en pacientes sin síntomas de infección hospitalaria, por cada 100.
  - ~~$B_3$~~ : Este parámetro nos dice que la probabilidad porcentual promedio estimada de adquirir infección en el hospital aumenta a razón de 0.0643261 unidades por cada unidad de aumento en el número promedio de camas en el hospital durante el periodo de estudio.
  - ~~$B_4$~~ : Este parámetro nos dice que la probabilidad porcentual promedio estimada de adquirir infección en el hospital aumenta a razón de 0.0053465 unidades por cada unidad de aumento en el número promedio de pacientes en el hospital por día durante el periodo del estudio.
  - ~~$B_5$~~ : Este parámetro nos dice que la probabilidad porcentual promedio estimada de adquirir infección en el hospital aumenta a razón de 0.0018688 unidades por cada unidad de aumento en el número promedio de enfermeras, equivalente a tiempo completo, durante el periodo de estudio.

#### 1.5 Coeficiente de determinación múltiple e interpretación $R^2$

*3pt*

El coeficiente de determinación múltiple  $R^2$ , está definido de la forma:

$$R^2 = 1 - \frac{\text{SUMA DE CUADRADO DEL ERROR (SCE)}}{\text{SUMA TOTAL DE CUADRADOS CORREGIDOS (STCC)}} \quad \text{ò} \quad \frac{\text{SUMA DE CUADRADOS DE LA REGRESIÓN (SCR)}}{\text{SUMA TOTAL DE CUADRADOS CORREGIDOS (STCC)}}$$

Así:

$$R^2 = 1 - \frac{51.0251}{118.1994} \quad \text{ò} \quad \frac{67.1743}{118.1994} = 0.56831$$

$$R^2 = 0.56831(100) = 56.831\%$$

El coeficiente de determinación múltiple (en porcentaje) nos dice que el 56.831% de la variabilidad total en la probabilidad promedio estimada de adquirir infección en el hospital es explicada por el modelo planteado. Hay que recordar que este coeficiente estima su valor teniendo en cuenta todas las variables, independientemente de si sus parámetros que las

acompañan son o no significativos. Es por eso, que existe otra medida sin interpretación que reduce este error de estimación y es el  $R^2_{\text{ajustado}}$ .

## Pregunta 2.

5 p +

### 2.1 Planteamiento de pruebas de hipótesis y modelo reducido

Las 3 co-variables con el **Valor P** más alto fueron  $X_2$ ;  $X_4$ ;  $X_5$ , con el uso de la tabla de todas las regresiones posibles pretendemos hacer la siguiente prueba de hipótesis:

¿Por qué  
ese formato?

$$\begin{cases} H_0 : \beta_2 = \beta_4 = \beta_5 = 0 \\ H_1 : \text{Algun } \beta_j \text{ distinto de 0 para } j = 2, 4, 5 \end{cases}$$

**Tabla 5. Resumen tabla de todas las regresiones**

	SSE	Covariables
Modelo completo	51.025	$X_1, X_2, X_3, X_4, X_5$
Modelo reducido	58.098	$X_1, X_3$

Teniendo en cuenta los datos de la **Tabla 5**, podemos armar un modelo reducido para la prueba de significancia del subconjunto:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_3 X_{3i} + \epsilon_i ; \epsilon_i \sim N(0, \sigma^2). \text{ iid}$$

$i = 1, 2, \dots, 60$

### 2.2 Cálculo y análisis del estadístico de prueba

$$F_0 = \frac{\frac{SSE(\beta_0, \beta_1, \beta_3) - SSE(\beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5)}{3}}{MSE(\beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5)} ; F_0 \sim f_{3,54}$$

Por lo tanto

$$F_0 = \frac{\frac{58.098 - 51.0251}{3}}{0.944909} = 2,4952$$

3 p +

Ya teniendo este resultado, obtenemos un valor crítico de la distribución  $F_{3,54}$  con un nivel de significancia  $\alpha = 0.05$ , que es igual a  $F_{0.05,3,54} = 2.7758$ .

Como  $F_0 = 2.4952 < F_{0.05,3,54} = 2.7758$ , entonces se acepta  $H_0$  por cual el subconjunto no es significativo, quiere decir que es posible descartar las variables del subconjunto, se concluye que no hay suficiente evidencia para decir que al menos una de las variables  $X_2$ ,  $X_4$  y  $X_5$  tiene un efecto significativo en la variable de respuesta  $Y_i$ . En este caso, se debe utilizar el modelo reducido que solo incluye las variables  $X_1$ ,  $X_3$  y la constante para hacer predicciones.

### Pregunta 3.

¿El valor de la variable  $\beta_1$  es equivalente al valor de  $\beta_2$  y el valor de  $\beta_3$  es equivalente al valor de  $\beta_4$  en el modelo de RLM?

Para resolver, se tiene la siguiente prueba de hipótesis:

$$H_0: \begin{cases} \beta_1 = \beta_2 \\ \beta_3 = \beta_4 \end{cases} \quad \text{vs} \quad H_1: \begin{cases} \beta_1 \neq \beta_2 \\ \beta_3 \neq \beta_4 \end{cases}$$

o en forma matricial:

$$H_0: L\beta = \underline{0} \quad \text{vs.} \quad H_1: L\beta \neq \underline{0},$$

donde tenemos que:

$$L = \begin{bmatrix} 0 & 1 & -1 & 0 & 0 \\ 0 & 0 & 0 & 1 & -1 \end{bmatrix} \quad \times \quad I = \begin{bmatrix} 0 & 1 & -1 & 0 & 0 \\ 0 & 0 & 0 & 1 & -1 \end{bmatrix}$$

por lo tanto  $H_0$  se puede plantear como:

$$H_0: \begin{bmatrix} 0 & 1 & -1 & 0 & 0 \\ 0 & 0 & 0 & 1 & -1 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

congruentes con el error, sin embargo y debe tener todos los  $\beta$ 's del modelo, así como incluir a  $\beta_0$ , también

El modelo reducido está dado por:

$$Y_i = \beta_0 + \beta_1(X_{1i} + X_{2i}) + \beta_3(X_{3i} + X_{4i}) + \epsilon_i; \epsilon_i \sim N(0, \sigma^2). \text{ iid } 1 \leq i \leq 60$$

¿y  $\beta_5$ ? ¿no sacaron definitivamente del modelo o qué pasó ahí?



Para simplificar se tiene que:

$$Y_i = \beta_0 + \beta_1 X_{1i}^* + \beta_3 X_{3i}^* + \epsilon_i; X_{1i}^* = X_{1i} + X_{2i}, X_{3i}^* = X_{1i} + X_{2i}$$

la expresión del estadístico de prueba es:

$$F_0 = \frac{MSH}{MSE} = \frac{SSH/2}{MSE} = \frac{(SSE(RM) - SSE(FM))/2}{MSE} \sim F_{2,54}$$

$$F_0 = \frac{SSE(RM) - 25,51}{0.944909}$$

Pregunta 4.

4.1 Normalidad de los residuales

Para la validación de este supuesto, se planteará la siguiente prueba de hipótesis ~~Shapiro-Wilk~~, acompañada de un gráfico cuantil-cuantil:

$$\begin{cases} H_0 : \epsilon_i \sim \text{Normal} \\ H_1 : \epsilon_i \not\sim \text{Normal} \end{cases}$$

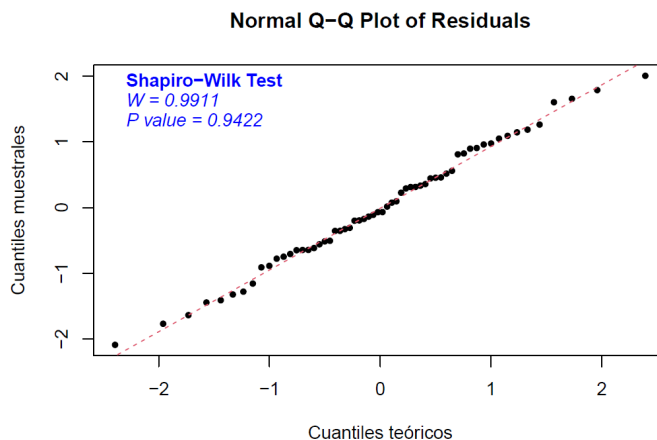


Figura 1: Gráfico cuantil-cuantil y normalidad de residuales

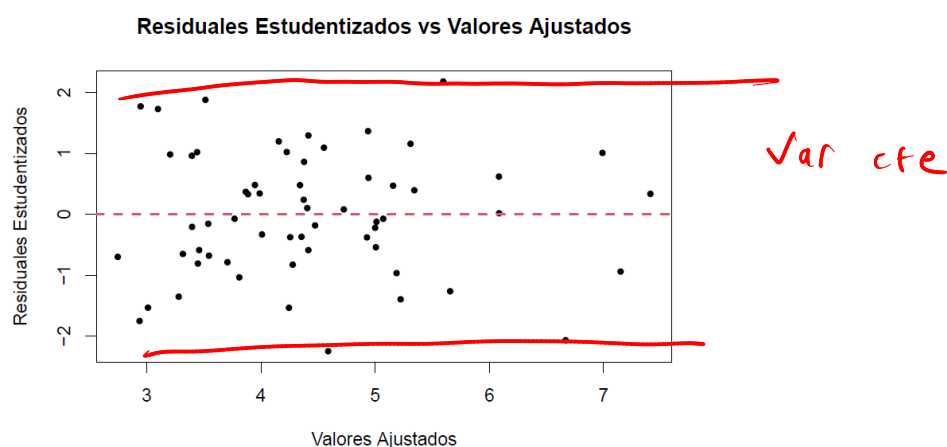
Al ser el **valor P** aproximadamente igual a **0.9422** y teniendo en cuenta que el nivel de significancia  $\alpha = 0.05$ , el p-valor es mucho mayor y por lo tanto, no se rechazaría la hipótesis

nula, es decir que los datos distribuyen normal con media 0 y varianza  $\sigma^2$  constante, esto es soportado además por el gráfico de comparación de cuantiles, se puede observar en la figura anterior que las colas no son tan pesadas con respecto a la diagonal trazada. Se determina por aceptar el cumplimiento de este supuesto. ✓

#### 4.2 Varianza constante

3 p +

Ahora se validará si la varianza cumple con el supuesto de ser constante. Realizaremos el análisis a través de un gráfico de residuales vs. valores ajustados en donde buscaremos patrones en la nube de puntos generada.

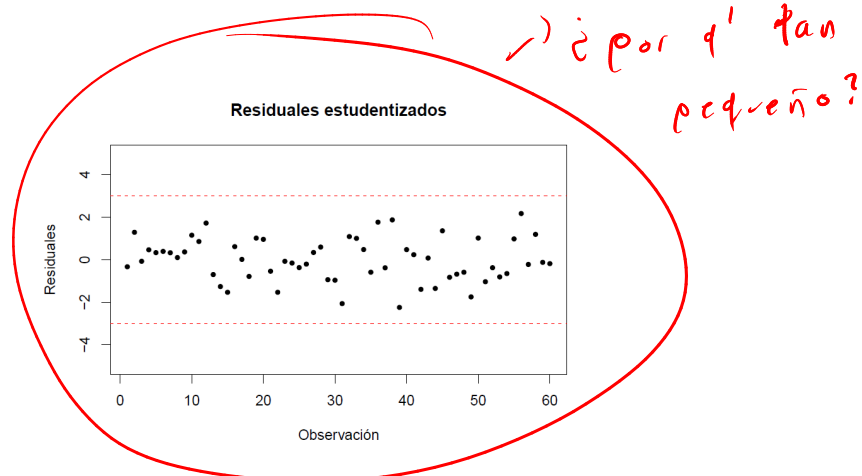


En el gráfico de Residuales Estudentizados vs Valores Ajustados se puede observar que no hay patrones en los que se note que la varianza aumenta o decrece, ni un comportamiento que permita descartar una varianza constante, al no haber evidencia suficiente en contra de este supuesto se acepta como cierto. Además es posible observar media 0. ✓

#### 4.3 Datos atípicos

3 p +

Se propone ahora un gráfico de residuales estudentizados que nos permita observar la existencia de valores atípicos.

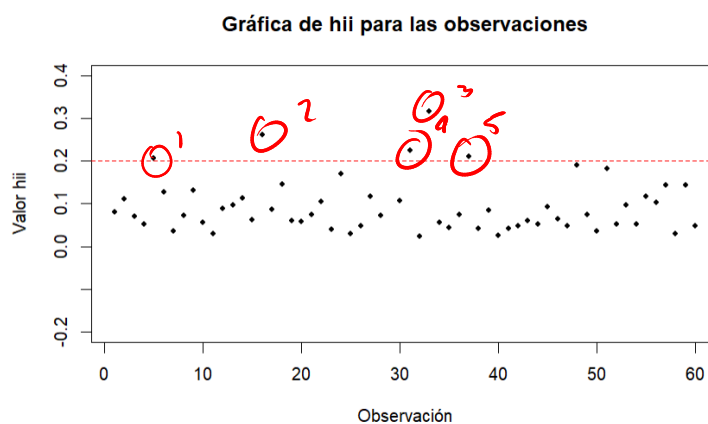


Como se puede observar en la gráfica anterior, no hay datos atípicos en el conjunto de datos pues ningún residual estudentizado sobrepasa el criterio de  $|r_{\text{estud}}| > 3$ . ✓

#### 4.4 Puntos de balanceo

2 pt

Para buscar los puntos de balanceo, que son aquellos tal que  $h_{ii} > 2 \frac{p}{n} = 0.2$ , nos apoyamos del gráfico de las observaciones versus los valores  $h_{ii}$



tienen 6 puntos de balanceo y en la gráfica sólo se ven 5,  $h_{ii}$  del dato 29 = 0.4959, debían ampliar límite superior a 0.5

Se identifica de la **tabla 6** y del gráfico anterior, que las observaciones **5, 16, 29, 31, 33 y 37** son puntos de balanceo. Estos son puntos alejados de los valores predictores, que pueden ser parte de la causa de nuestro pequeño coeficiente de determinación.

¿Qué más causan estos puntos de balanceo?

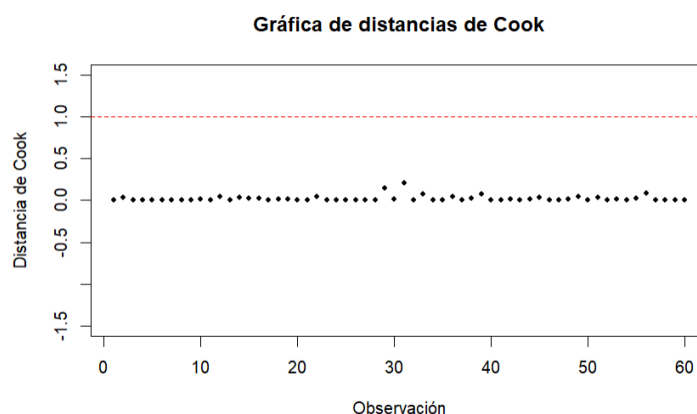
**Tabla 6. listado de puntos de balanceo**

	res.stud	Cooks.D	hii.value	Dffits
5	0,3346	0,0049	0,2064	0,1692
16	0,6176	0,0224	0,2609	0,3648
29	-0,9406	0,1450	0,4959	-0,9319
31	-2,0668	0,2072	0,2254	-1,1511
33	1.0066	0.0779	0.3156	0.6837
37	-0,3802	0,0064	0,2111	-0,1951

#### 4.5 Puntos influyentes

Para buscar puntos influyentes usaremos los criterios de las Distancias de Cook ( $D_i > 1$ ) y

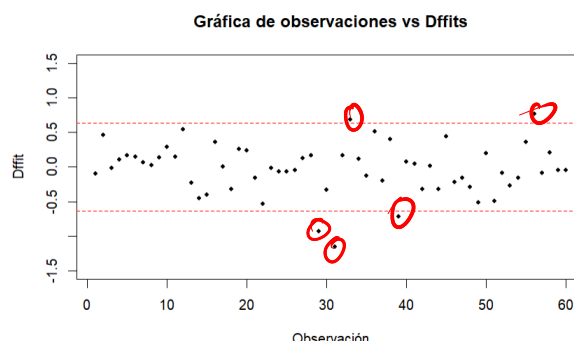
$$|DFITS_i| > 2\sqrt{\frac{p}{n}}.$$



Usando distancias de cook, en el gráfico anterior observamos que no hay  $D_i > 1$ , entonces por medio de este criterio no se encontraron puntos influyentes. ✓ 2 pr

Para el criterio de **DFITS**,  $|DFITS_i| > 2\sqrt{\frac{p}{n}} = 0.6324$ , el gráfico nos muestra la existencia de algunos puntos ~~influyentes~~:

*influyentes*



1,5 pt

explican bien que  
lo que causan es mover  
en vector

Son 5 puntos que corresponden a las observaciones **29, 31, 33, 39 y 56** que son puntos los cuales arrastran el modelo hacia su dirección. A continuación se muestra en la tabla 7 los valores influenciales hallados:

**Tabla 7. listado de puntos de influenciabiles (Criterio DFFITS)**

	res.stud	Cooks.D	hii.value	Dffits
<b>29</b>	-0,9406	0,145	0,4959	-0,9319
<b>31</b>	-2,0668	0,2072	0,2254	-1,1511
<b>33</b>	1,0066	0,0779	0,3156	0,6837
<b>39</b>	-2,247	0,0786	0,0854	-0,7145
<b>56</b>	2,1769	0,0909	0,1032	0,7658

#### 4.6 Conclusiones del modelo

2 pt

Del análisis de los supuestos podríamos concluir que el modelo, <sup>¿dado q. e...?</sup> dado que tiene un  $R^2$  (explica aproximadamente el 56% de la variabilidad del modelo), los errores cumplen ~~de una manera~~ ~~positiva~~ los supuestos, pero hay que tener en cuenta que los valores influenciales dado que arrastran el modelo hacia su dirección, podrían “desviar” el modelo de una mejor aproximación,

como también podría ser algo “natural” del experimento estos puntos. Esto, se podría observar comparando modelos con y sin los puntos influenciales.

¿Es válido o no?