

9,0

Trabajo 1

Estudiantes

Alejandro Diaz López
Juan José Flórez Ospina
Juan Diego Giraldo Jaramillo
Mariam Saavedra Navaja

Equipo 6

Docente:

Julieth Veronica Guarín Escudero

Asignatura:

Estadística II



UNIVERSIDAD
NACIONAL
DE COLOMBIA

Sede Medellín
30 de marzo de 2023

Índice

1. Pregunta 1	3
1.1. Modelo de regresión	3
1.2. Significancia de la regresión	3
1.3. Significancia de los parámetros	4
1.4. Interpretación de los parámetros	5
1.5. Coeficiente de determinación múltiple R^2	5
2. Pregunta 2	5
2.1. Planteamiento pruebas de hipótesis	5
2.2. Estadístico de prueba y conclusión	6
3. Pregunta 3	6
3.1. Prueba de hipótesis y prueba de hipótesis matricial	6
3.2. Estadístico de prueba	7
4. Pregunta 4	7
4.1. Supuestos del modelo	7
4.1.1. Normalidad de los residuales	7
4.1.2. Varianza constante	9
4.2. Verificación de las observaciones	10
4.2.1. Datos atípicos	10
4.2.2. Puntos de balanceo	11
4.2.3. Puntos influenciales	12
4.3. Conclusión	13

Índice de figuras

1.	Gráfico cuantil-cuantil y normalidad de residuales	8
2.	Gráfico residuales estudentizados vs valores ajustados	9
3.	Identificación de datos atípicos	10
4.	Identificación de puntos de balanceo	11
5.	Criterio distancias de Cook para puntos influenciales	12
6.	Criterio Dffits para puntos influenciales	13

Índice de cuadros

1.	Tabla de valores coeficientes del modelo	3
2.	Tabla ANOVA para el modelo	4
3.	Resumen de los coeficientes	4
4.	Tabla de todas las regresiones resumida	5

1. Pregunta 1 18pt

Teniendo en cuenta la base de datos brindada, en la cual hay 5 variables regresoras dadas por:

Y : Riesgo de infección ✓

X_1 : Duración de la estadía ✓

X_2 : Rutina de vultivos ✓

X_3 : Número de camas ✓

X_4 : Censo promedio diario ✓

X_5 : Número de enfermeras ✓

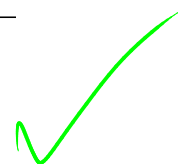
$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + \beta_4 X_{i4} + \beta_5 X_{i5} + \varepsilon_i, \varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2); 1 \leq i \leq 60 \quad \checkmark$$

1.1. Modelo de regresión 2pt

Al cargar y ajustar el modelo lineal, se obtienen los siguientes coeficientes estimados del MRLM:

Cuadro 1: Tabla de valores coeficientes del modelo

Valor del parámetro	
β_0	0.0882
β_1	0.1873
β_2	-0.0100
β_3	0.0605
β_4	0.0203
β_5	0.0008



ec ajustada no
leva supuestos

Por lo tanto, el modelo de regresión ajustado es:

$$\hat{Y}_i = 0.0882 + 0.1873X_{1i} - 0.01X_{2i} + 0.0605X_{3i} + 0.0203X_{4i} + 0.0008X_{5i} + \varepsilon_i, \varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2); 1 \leq i \leq 60$$



1.2. Significancia de la regresión 4pt

Para realizar la prueba de significancia del modelo de regresión planteamos el siguiente juego de hipótesis

$$\begin{cases} H_0 : \beta_1 = \beta_2 = \beta_3 = \beta_4 = \beta_5 = 0 \\ H_1 : \text{algún } \beta_j \neq 0, j=1, 2, 3, 4, 5 \end{cases}$$



Cuyo estadístico de prueba es:

$$F_0 = \frac{MSR}{MSE} \stackrel{H_0}{\sim} f_{5,54} \quad \checkmark \quad (1)$$

Ahora, se presenta la tabla Anova:

Cuadro 2: Tabla ANOVA para el modelo

	Sumas de cuadrados	Grados de libertad.	Cuadrado medio	F_0	valor P
Regresión	91.8115	5	18.362302	21.6734	7.80158e-12
Error	45.7503	54	0.847228		

De la tabla Anova, se compara el valor P que es de 7.80158e-12, con un nivel de significancia $\alpha = 0.05$ permitiendo que se rechaze la hipótesis nula en la que $\beta_1 = \beta_2 = \beta_3 = \beta_4 = \beta_5 = 0$, probando ~~que existe una relación en la regresión.~~

esa no es la conclusión bajo la p.H.

1.3. Significancia de los parámetros 6pt

En la siguiente tabla se aprecia la información de los parámetros que permite determinar la significancia de cada uno de estos.

Cuadro 3: Resumen de los coeficientes

	Valorestimado	$SE(\hat{\beta}_j)$	T_{0j}	P-valor
β_0	0.0882	1.5010	0.0588	0.9533
β_1	0.1873	0.0690	2.7145	0.0089
β_2	-0.0100	0.0285	-0.3514	0.7267
β_3	0.0605	0.0154	3.9187	0.0003
β_4	0.0203	0.0072	2.8405	0.0063
β_5	0.0008	0.0008	1.0852	0.2826

Para analizar la significancia de los coeficientes en el modelo lineal emplearemos un nivel de significancia = 0.05 para comparar con el valor P arrojado por el resumen de los coeficientes.

Los valores P presentes en la tabla permiten concluir que con un nivel de significancia $\alpha = 0.05$ que los parámetros β_1 , β_3 y β_4 son significativos, pues sus valores P son menores a α .

1.4. Interpretación de los parámetros 3 pt

$\hat{\beta}_1$: Por cada día que aumenta la estadía del paciente en el hospital, la probabilidad promedio del riesgo de infección aumenta en 0.187 cuando el resto de variables predictoras se mantienen fijas. ✓

$\hat{\beta}_3$: Por cada incremento en 1 unidad de camas en el hospital, la probabilidad promedio del riesgo de infección aumenta en 0.0605 cuando el resto de variables predictoras se mantienen fijas. ✓

$\hat{\beta}_4$: Por cada aumento en la cantidad promedio de enfermeras en el hospital, la probabilidad promedio del riesgo de infección aumenta en 0.0008 cuando el resto de variables predictoras se mantienen fijas. ✓

1.5. Coeficiente de determinación múltiple R^2 3 pt

Para hallar el coeficiente de determinación múltiple R^2 empleamos el SSR y SSE dados por la tabla ANOVA

$$R^2 = \frac{SSR}{SST} = \frac{SSR}{SSR+SSE} = \frac{91.8115}{137.5618} = 0.667420 \quad \checkmark$$

Este coeficiente de determinación nos dice que aproximadamente el 66.74 % de la variabilidad observada en la respuesta es explicada por el modelo de regresión propuesto en el presente informe. ✓

2. Pregunta 2 4 pt

2.1. Planteamiento pruebas de hipótesis

Las variables con el valor P más alto en el modelo fueron X_1, X_2, X_5

para probar la significancia simultánea de estos tres coeficientes de la regresión planteamos las hipótesis:

$$\begin{cases} H_0 : \beta_1 = \beta_2 = \beta_5 = 0 \\ H_1 : \text{algún } \beta_j \neq 0, j=1, 2, 5 \end{cases} \quad \checkmark$$

Cuadro 4: Tabla de todas las regresiones resumida

	SSE	Variables en el modelo				
Modelo completo	45.750	X1	X2	X3	X4	X5
Modelo reducido	57.368			X3	X4	

El modelo completo es el visto en el inicio de la pregunta 1.

El modelo reducido es de la forma:

$$Y_i = \beta_0 + \beta_3 X_{3i} + \beta_4 X_{4i} + \varepsilon_i; \varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2); 1 \leq i \leq 60$$

2.2. Estadístico de prueba y conclusión

Se calcula el estadístico de la prueba de la forma:

$$\begin{aligned} F_0 &= \frac{SSR(\beta_1, \beta_2, \beta_5 | \beta_0, \beta_3, \beta_4) / 3}{SSE(\beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5) / 54} \\ &= \frac{(SSE(\beta_0, \beta_3, \beta_4) - SSE(\beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5)) / 3}{SSE(\beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5) / 54} \\ &= \frac{11.618 / 3}{47.750 / 54} = 4.349557 \stackrel{H_0}{\sim} f_{3,54} \end{aligned} \quad (2)$$

Para el criterio de decisión calculamos el valor crítico a un nivel de significancia $\alpha = 0.05$ de una distribución $f_{0.05, 3, 54} = 4.349557$

Como $F_0 > f_{0.05, 3, 54}$, se rechaza la hipótesis nula, lo que quiere decir que al menos una de las variables regresoras asociadas a la Duración de la estadía, Rutina de cultivos y Número de enfermeras (X_1, X_2, X_5), es significativa en presencia del resto de variables y por lo tanto hace a este conjunto un conjunto significativo y no podemos descartarlo. subconjunto.

3. Pregunta 3

3.1. Prueba de hipótesis y prueba de hipótesis matricial

Se quiere estudiar si el efecto de la duración de estadía de los pacientes en el hospital es igual a la rutina de cultivos realizados en los pacientes sin síntoma de infección hospitalaria, por cada 100 pacientes. Además deseamos estudiar si el efecto promedio de camas en el hospital durante el periodo del estudio es igual al efecto del número promedio de pacientes en el hospital por día durante el periodo de estudio.

Para responder a la pregunta se plantea la siguiente prueba de hipótesis:

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

O equivalentemente,

$$H_0 : \beta_1 - \beta_2 = 0; \beta_3 - \beta_4 = 0$$

Además, se puede representar matricialmente de la siguiente forma:

$$\begin{cases} H_0 : \mathbf{L}\underline{\beta} = \underline{\mathbf{0}} \\ H_1 : \mathbf{L}\underline{\beta} \neq \underline{\mathbf{0}} \end{cases}$$

Con \mathbf{L} dada por

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

El modelo reducido es:

$$\begin{aligned} Y &= \beta_o + \beta_1(X_1 + X_2) + \beta_3(X_3 + X_4) + \beta_5 X_5 + \varepsilon \\ &= \beta_o + \beta_1 X_{1,2} + \beta_3 X_{3,4} + \beta_5 X_5 \end{aligned} \quad (3)$$

Donde $X_{1,2} = X_1 + X_2$ y $X_{3,4} = X_3 + X_4$

3.2. Estadístico de prueba

Se tiene que el estadístico de prueba F_0 está dado por:

$$\begin{aligned} F_0 &= \frac{SSH/gl.ssh}{MSE(MF)} = \frac{(SSE(MR) - SSE(MF))/2}{MSE(MF)} \\ &= \frac{(SSE(MR) - 45.7503)/2}{0.847228} \stackrel{H_0}{\sim} f_{2,54} \end{aligned} \quad (4)$$

donde $SSE(RM)$ corresponde al error estándar del modelo reducido, $SSE(RM)$ corresponde al error estándar del modelo full, r corresponde al número de filas linealmente independientes en la matriz L que son 2, y el $MSE(FM)$ a la media de errores estándar del modelo full.

Si $F_0 > f_{0.05,2,54}$ entonces se rechaza la hipótesis nula y al menos una de los dos supuestos que queremos probar no se cumple con una significancia del 95 %.

4. Pregunta 4

4.1. Supuestos del modelo

4.1.1. Normalidad de los residuales

Para la validación de la normalidad de los residuales, se plantea la siguiente prueba de hipótesis ~~shapiro-wilk~~, acompañada de un gráfico cuantil-cuantil:

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

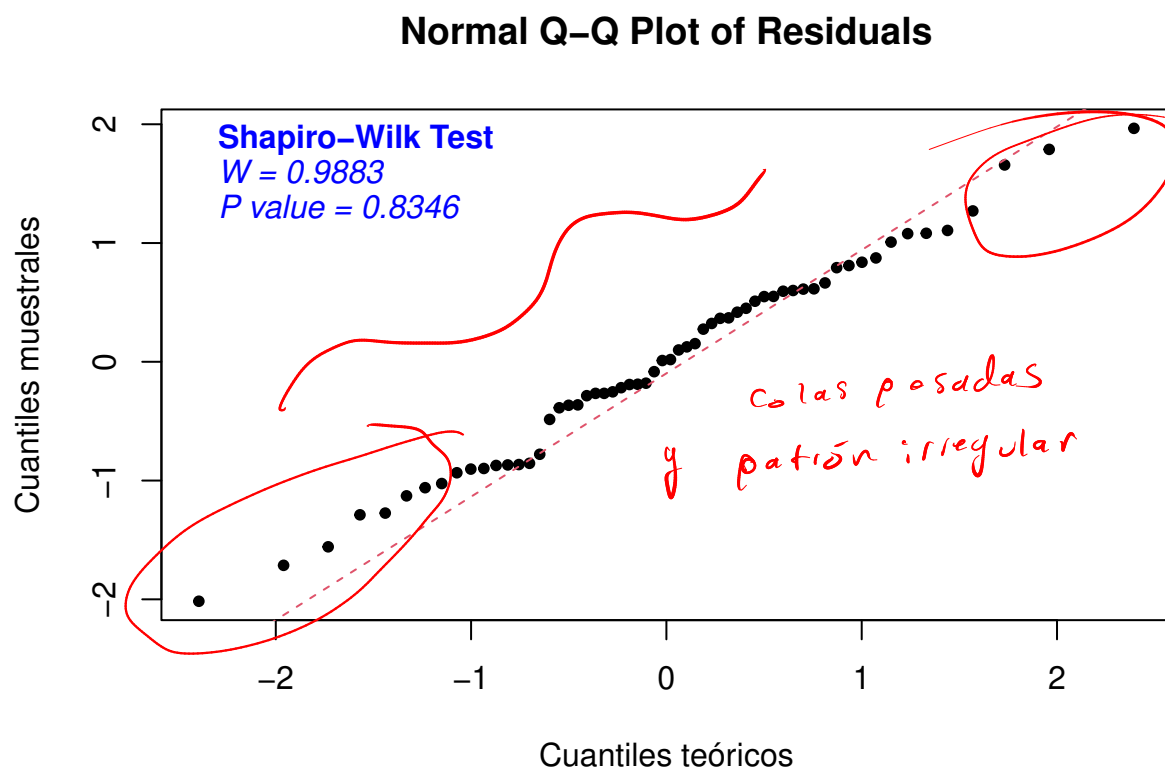


Figura 1: Gráfico cuantil-cuantil y normalidad de residuales

En la prueba teórica de normalidad Shapiro-Wilk se ^{puede} notar que el valor P es 0.8346, al ser mayor que $\alpha = 0.05$, podemos afirmar que los errores cumplen el supuesto de distribución normal. En esta gráfica también se pueden observar algunos datos aparentemente “outliers” que se comprobarán con los criterios que siguen a continuación. Con la gráfica y el valor- P , podemos aceptar la normalidad de los errores en nuestro modelo.

Muy mal análisis gráfico, no distribuye normal

3 p 1

4.1.2. Varianza constante

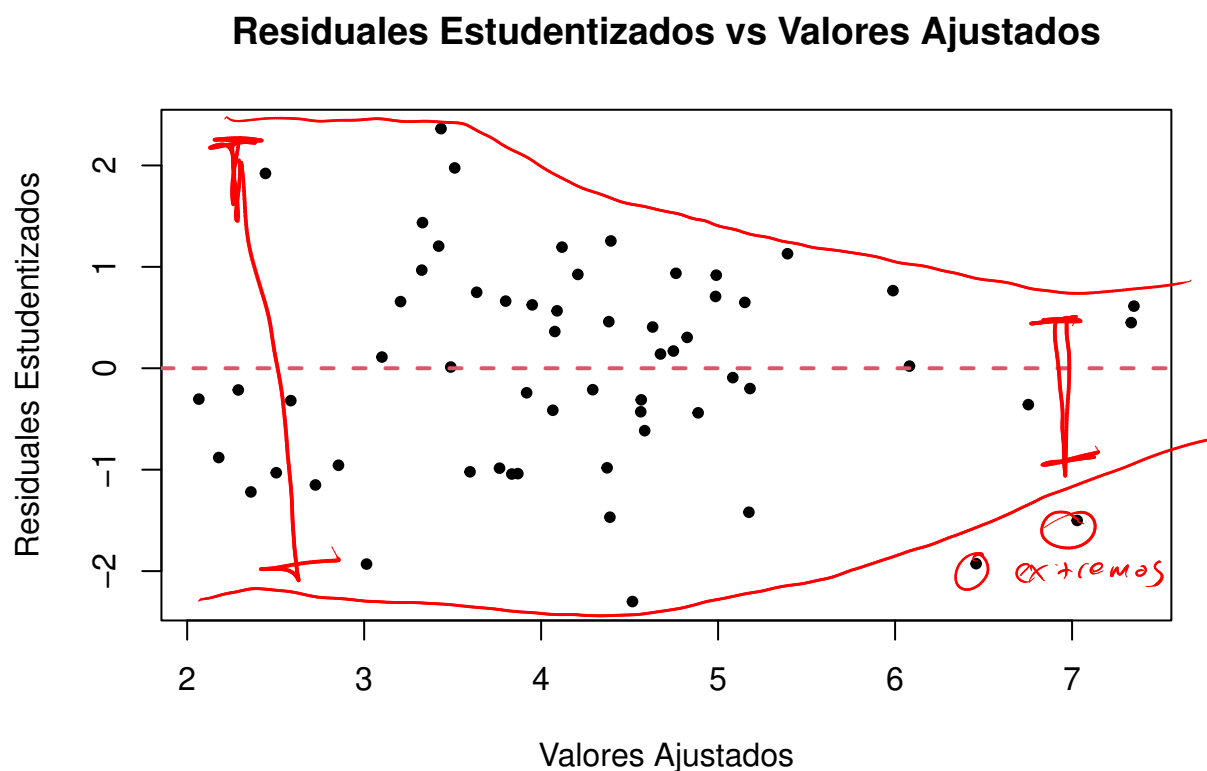


Figura 2: Gráfico residuales estudentizados vs valores ajustados

Para analizar si la varianza de los residuales es costante, vemos esta gráfica donde podemos notar algunas tendencias de varianza no constante, es decir, se puede percibir una concavidad en los cuantiles positivos y una convexidad en los cuantiles negativos, sin embargo, esta evidencia gráfica no es lo suficientemente significativa como para rechazar el supuesto de varianza constante.

sí lo es, sin embargo por el análisis les vaigo el punto

su model no es válido por normalidad ni por var cte.

4.2. Verificación de las observaciones

4.2.1. Datos atípicos

3 p +

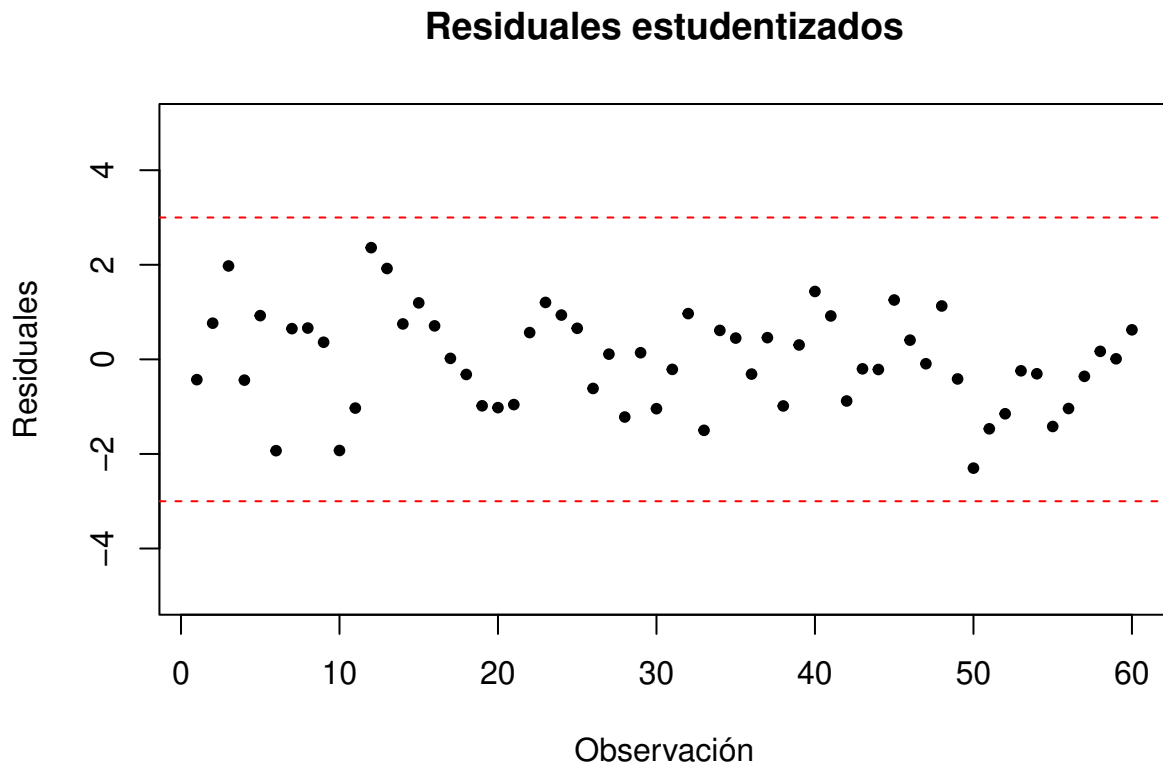


Figura 3: Identificación de datos atípicos

Como se puede observar en la gráfica anterior, no hay datos atípicos en el conjunto de datos pues ningún residual estudentizado cumple el criterio de $|r_i| > 3$. Esto significa que nuestro modelo no tiene “outliers”.

✓

4.2.2. Puntos de balanceo 2,5 pt

Gráfica de hii para las observaciones

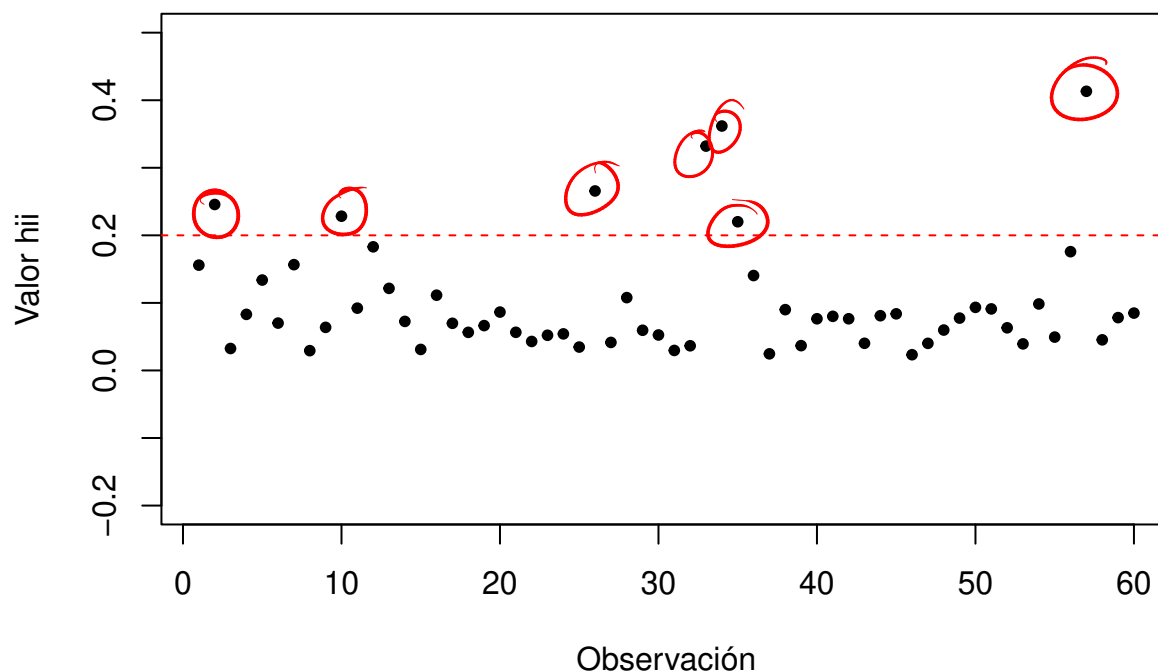


Figura 4: Identificación de puntos de balanceo

##	res.stud	Cooks.D	hii.value	Dffits
## 2	0.7655	0.0318	0.2457	0.4352
## 10	-1.9267	0.1830	0.2283	-1.0758
## 26	-0.6152	0.0228	0.2657	-0.3679
## 33	-1.5005	0.1866	0.3321	-1.0707
## 34	-0.6126	0.0355	0.3618	0.4586
## 35	0.4500	0.0095	0.2201	0.2373
## 57	-0.3587	0.0151	0.4132	-0.2985

→ No se presentan salidas de R en un reporte estadístico como este.

En esta ~~tabla~~ se puede observar que los puntos 2, 10, 26, 33, 34, 35 y 57 cumplen con el criterio $h_{ii} > 2\frac{p}{n}$. Por lo tanto estos son puntos de balanceo, lo cual nos dice que son puntos que se encuentran muy alejados respecto a los otros puntos en el eje X. Veremos con la prueba de distancias de Cook y Diagnóstico DFFITS si estos son puntos influyentes o no. ~

, cuánto da $2\frac{p}{n}$?

¿Qué causa esto en el modelo?

4.2.3. Puntos influenciales

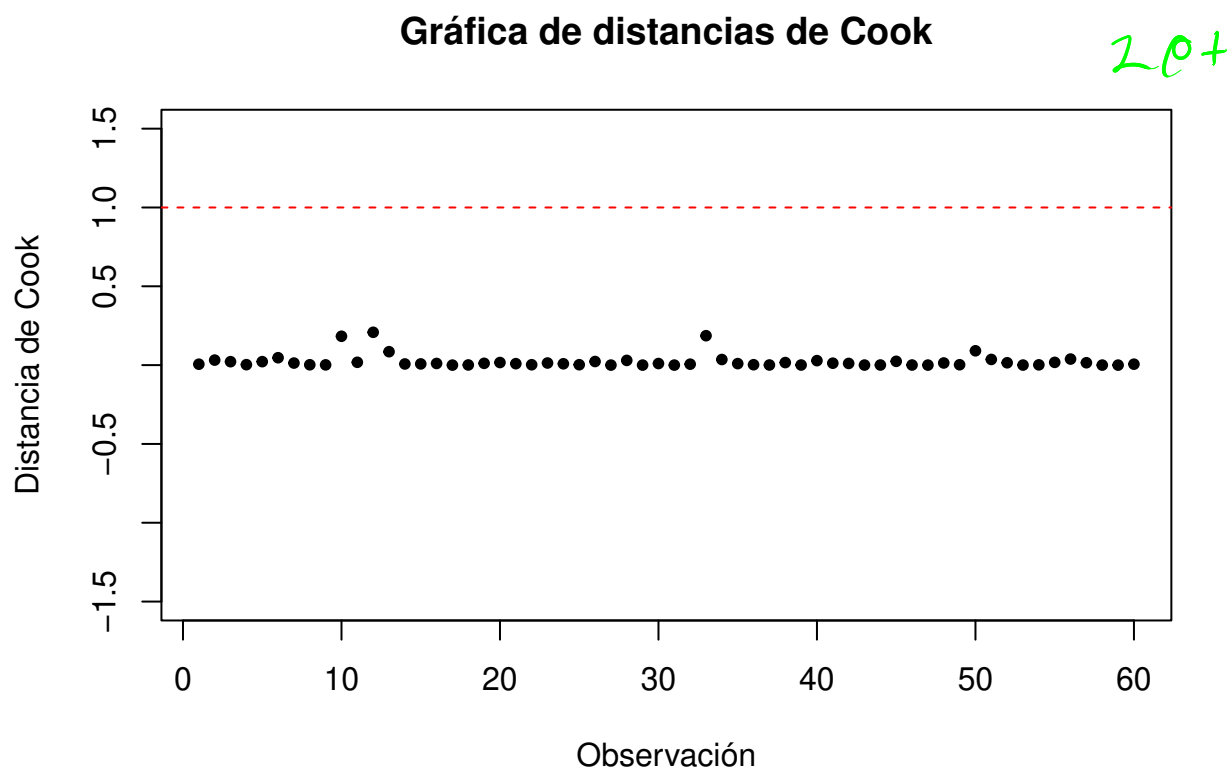


Figura 5: Criterio distancias de Cook para puntos influenciales

Según este criterio de $D_i > 1$, no hay ningún punto que cumpla las condiciones necesarias para llamarlos puntos influenciales, pero usaremos también el método DFFITS para determinar si alguno de los ~~residuales~~ ^{datos} del modelo son puntos influenciales.

~~residuales~~
datos

~~método~~
criterio

Gráfica de observaciones vs Dffits

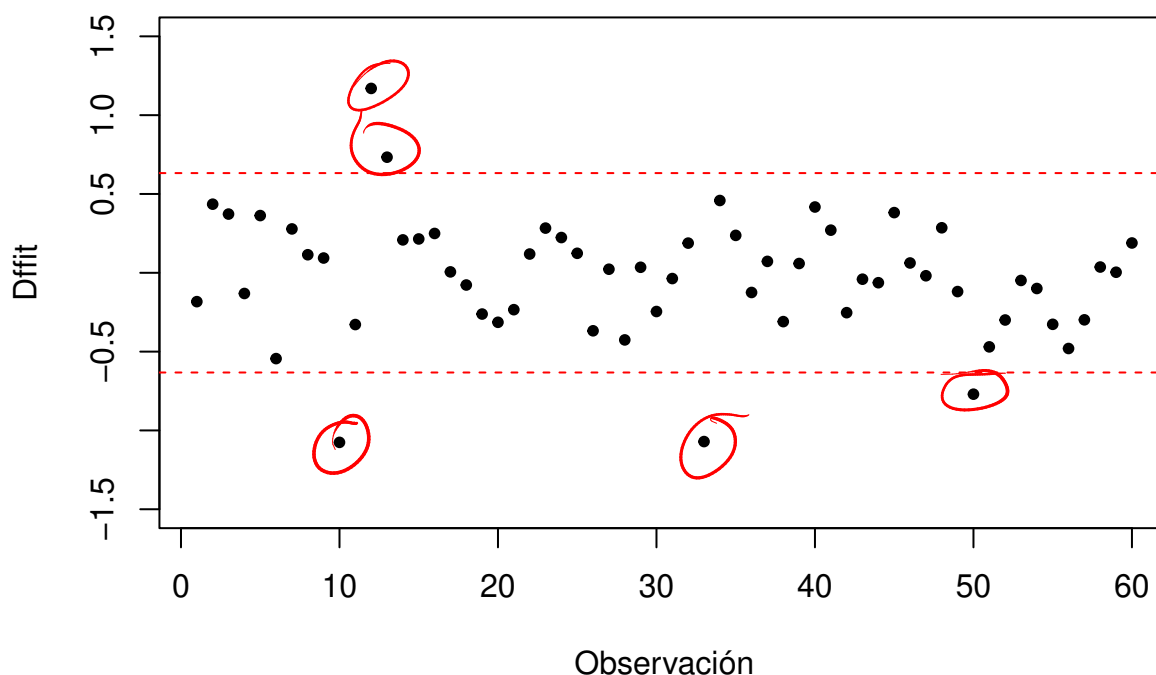


Figura 6: Criterio Dffits para puntos influyentes

##	res.stud	Cooks.D	hii.value	Dffits
## 10	-1.9267	0.1830	0.2283	-1.0758
## 12	2.3624	0.2083	0.1830	1.1697
## 13	1.9211	0.0851	0.1215	0.7334
## 33	-1.5005	0.1866	0.3321	-1.0707
## 50	-2.3003	0.0908	0.0934	-0.7700

No dicen qué causan los puntos influyentes

los datos

En esta gráfica se puede notar que hay 5 puntos influyentes, las muestras 10, 12, 13, 33 y 50. Estos puntos cumplen con el criterio $|DFITS_i| > 2\sqrt{\frac{p}{n}}$.

Todos los puntos influyentes deberían ser investigados, pero en especial, comparando este criterio con el de puntos de balanceo podemos decir entonces que los puntos 10 y 33 son puntos de balanceo e influyentes que pueden alterar nuestros análisis, por lo que sería recomendable tener especial atención a la hora de aplicar el modelo al incluir estos puntos, pues son los datos más atípicos.

4.3. Conclusión

Podemos concluir que nuestro modelo funciona, cumple con el supuesto de normalidad en los errores, y aunque hay una breve tendencia de no normalidad en la varianza, no es sufi-

¿están probando normalidad en varianza?

No hay bueno o malo.

¿Normalidad de varianzas o el de homocedasticidad?

14

cientemente contundente la evidencia gráfica para rechazar este supuesto. Se puede concluir también que hay pocos puntos de balanceo y que en general la muestra elegida se comporta bastante bien con el modelo planteado.

¡, hay 7 !!!

No dicen si el modelo es válido.