

3,85

Trabajo 1

Estudiantes

**Jalan Howard Hudgson
Juan José Lopez Taborda
Diego Quintero Gaítan**

Docente

Francisco Javier Rodríguez Cortes

Asignatura

Estadística II



UNIVERSIDAD
NACIONAL
DE COLOMBIA

Sede Medellín
27 de Marzo de 2023

Índice

1. Pregunta 1	2
1.1. Modelo de regresión	2
1.2. Significancia de la regresión	2
1.3. Significancia de los parámetros	3
1.4. Interpretación de los parámetros	3
1.5. Coeficiente de determinación múltiple R^2	4
1.6. Comentarios	4
2. Pregunta 2	4
2.1. Planteamiento prueba de hipótesis y modelo reducido	4
2.2. Estadístico de prueba y conclusiones	4
3. Pregunta 3	5
3.1. Prueba de hipótesis y prueba de hipótesis matricial	5
3.2. Estadístico de prueba	5
4. Pregunta 4	5
4.1. Supuestos del modelo	6
4.1.1. Normalidad de los residuales	6
4.1.2. Media 0 y Varianza constante	6
4.2. Observaciones extremas	8
4.2.1. Datos atípicos	8
4.2.2. Puntos de balanceo	9
4.2.3. Puntos influyentes	10
4.3. Conclusiones	11

Índice de figuras

1. Gráfico cuantil-cuantil y normalidad de los residuales	6
2. Gráfico residuales estudentizados vs valores ajustados	7
3. Identificación de datos atípicos	8
4. Identificación de puntos de balanceo	9
5. Criterio distancias de Cook para puntos influyentes	10
6. Criterio Dffits para puntos influyentes	10

Índice de tablas

1.	Tabla de valores de los coeficientes estimados	2
2.	Tabla anova significancia de la regresión	3
3.	Resumen de los coeficientes	3
4.	Resumen de todas las regresiones	4
5.	Tabla de puntos de Balanceo	9
6.	Tabla del criterio DFFITS para encontrar puntos influenciales	11

1. Pregunta 1

17 p+

Estime un modelo de regresión lineal múltiple que explique el riesgo de infección en términos de las variables restantes (actuando como predictoras) Analice la significancia de la regresión y de los parámetros individuales. Interprete los parámetros estimados. Calcule e interprete el coeficiente de determinación múltiple R^2

Teniendo en cuenta la base de datos asignada a nuestro equipo, la cual es **Equipo42.txt**, las variables para el modelo son

Y RI Riesgo de infección en porcentaje: Probabilidad promedio estimada de adquirir infección en el hospital.

X1 DEHOS Duración de la estadía en días: Duración promedio de la estadía de todos los pacientes en el hospital.

X2 RC Rutina de cultivos: Razón del número de cultivos realizados en pacientes sin síntomas de infección hospitalaria, por cada 100 pacientes.

X3 NCP Número de camas: Promedio de camas en el hospital durante el periodo del estudio.

X4 CPD Censo promedio diario: Número promedio de pacientes en el hospital por día durante el periodo del estudio.

X5 ENFP Número de enfermeras: Promedio de enfermeras, equivalentes a tiempo completo, durante el periodo del estudio.

El modelo que se propone es:

$$RI_i = \beta_0 + \beta_1 DEHOS_i + \beta_2 RC_i + \beta_3 NCP_i + \beta_4 CPD_i + \beta_5 ENFP_i + \varepsilon_i, \quad \varepsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

1.1. Modelo de regresión

$i = 1, 2, \dots, 65$
20+

Al ajustar el modelo de regresión para el riesgo de infección en un hospital, se obtienen los siguientes coeficientes:

Tabla 1: Tabla de valores de los coeficientes estimados

	Valor del parámetro
$\hat{\beta}_0$	-1.84471
$\hat{\beta}_1$	0.18354
$\hat{\beta}_2$	0.02896
$\hat{\beta}_3$	0.04178
$\hat{\beta}_4$	0.02269
$\hat{\beta}_5$	0.00121

Por lo que el modelo con los respectivos valores de los parámetros es:

$$\widehat{RI}_i = -1.84471 + 0.18354 DEHOS_i + 0.02896 RC_i + 0.04178 NCP_i + 0.02269 CPD_i + 0.00121 ENFP_i + \varepsilon_i, \quad \varepsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

Donde las variables se mueven de acuerdo $1 \leq i \leq 65$

no va en
ec ajustada

1.2. Significancia de la regresión

4 p+

Se plantea el siguiente Juego de Hipótesis

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

Se utilizará la siguiente tabla ANOVA para evaluar la significancia de la regresión:

Tabla 2: Tabla anova significancia de la regresión

	Sumas de cuadrados	g.l	Cuadrado medio	F_0	Valor-P
Modelo de regresión	70.1988	5	14.039762	14.3115	3.61266e-09
Error	57.8797	59	0.981011		

Los resultados obtenidos de la Tabla Anova indican que la hipótesis nula debe ser rechazada. Esto nos lleva a concluir que la regresión es significativa según la evidencia muestral analizada.

1.3. Significancia de los parámetros

$$\begin{cases} H_0 : \beta_j = 0 \\ H_a : \text{Algún } \beta_j \neq 0 \text{ para } j = 1, 2, 3, 4, 5 \end{cases}$$

La tabla a continuación muestra los criterios utilizados para evaluar la significancia de los parámetros de forma individual:

Tabla 3: Resumen de los coeficientes

	Estimación β_j	$se(\hat{\beta}_j)$	T_{0j}	Valor-P
β_0	-1.8447	1.6298	-1.1319	0.2623
β_1	0.1835	0.0750	2.4472	0.0174
β_2	0.0290	0.0294	0.9844	0.3289
β_3	0.0418	0.0129	3.2407	0.0020
β_4	0.0227	0.0080	2.8505	0.0060
β_5	0.0012	0.0007	1.7369	0.0876

Los resultados de las pruebas: valor del estadístico de prueba y el valor p para la prueba se obtiene en las dos últimas columnas de la tabla de los parámetros estimados.

Con un nivel de significancia $\alpha = 0.05$ se concluye que los parámetros individuales $\beta_1, \beta_3, \beta_4$ son significativos cada uno en presencia de los demás parámetros. Por el contrario los parámetros $\beta_0, \beta_2, \beta_5$ individualmente no son significativos en presencia de los demás parámetros.

1.4. Interpretación de los parámetros

A continuación se hará la interpretación de los parámetros que son significativos, ya que los otros parámetros no tienen interpretación y no aportan al modelo.

- $\hat{\beta}_1 = 0.18354$: Si se mantiene constante el efecto de las demás variables predictoras, un incremento de un día en la Duración de la estancia en el hospital (medida en días) resultaría en un aumento esperado del promedio del Riesgo de infección en un ~~0.18354%~~ ^{18,354%}, según los resultados del análisis de regresión.
- $\hat{\beta}_3 = 0.04178$: Si el número promedio de camas en el hospital durante el periodo de estudio aumenta en una unidad, manteniendo constantes las demás variables predictoras, se espera que el promedio del Riesgo de infección aumente en un ~~0.04178%~~ ^{4,178%}.

- $\hat{\beta}_4 = 0.02269$: si el número censo del promedio Diario del paciente en el hospital durante el periodo de estudio se incrementa en una unidad, cuando las demás variables se mantienen constantes, se espera que el promedio del Riesgo de infección se incrementa en un ~~0.02269%~~ 2,269%

1.5. Coeficiente de determinación múltiple R^2

El modelo tiene un R^2 de 0.5481 lo cual significa que aproximadamente el 54.81 % de la variabilidad total en el porcentaje de Riesgo de infección es explicado por el modelo RLM ✓

¿Cómo se calcula?

1.6. Comentarios

En el modelo de regresión, se puede observar que las variables que tienen un aporte significativo son la Duración de la estadía en el hospital (DE), el Censo promedio diario de pacientes en el hospital (CDP) y el número de camas, lo cual se ve reflejado en la significancia de los parámetros. ✓

2. Pregunta 2

Use la tabla de todas las regresiones posibles, para probar la significancia simultánea del subconjunto de tres variables con los valores p más grandes del punto anterior. Según el resultado de la prueba es posible descartar del modelo las variables del subconjunto? Explique su respuesta.

2.1. Planteamiento prueba de hipótesis y modelo reducido

Los parámetros cuyos valores P fueron los más altos corresponden a β_2 con VP=0.3284, β_5 con VP= 0.087, β_1 con VP= 0.0174. Por tanto, se plantea la siguiente prueba de hipótesis:

$$\begin{cases} H_0 : \beta_1 = \beta_2 = \beta_5 = 0 \\ H_a : \text{Algún } \beta_j \neq 0 \text{ para } j = 1, 2, 5 \end{cases}$$

El modelo completo es el definido en la sección 1.1, y el modelo reducido es:

$$\text{MR: } \text{Rinf}_i = \beta_0 + \beta_3 \text{NCP}_i + \beta_4 \text{PD}_i + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$$

Se presenta la siguiente tabla con el resumen de todas las regresiones para plantear el estadístico de prueba:

Tabla 4: Resumen de todas las regresiones

	SSE	Covariables en el modelo
Modelo completo	57.880	X1 X2 X3 X4 X5
Modelo reducido	75.342	X3 X4

Así no se llaman sus var's

2.2. Estadístico de prueba y conclusiones

Se construye el estadístico de prueba como:

$$F_0 = \frac{(SSR(\beta_0, \beta_3, \beta_4 | \beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5) / 2)}{MSE(MF)} \stackrel{H_0}{\sim} f_{2, 9}$$

$SSR(\beta_1, \beta_2, \beta_5 | \beta_0, \beta_3, \beta_4) / 3$

MSE(MF)

3 X

$$F_0 = \frac{(SSE(\beta_0, \beta_3, \beta_4) - SSE(\beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5)) / 2}{MSE(\beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5)} \stackrel{H_0}{\sim} f_{2,59}$$

$$= \frac{(75.342 - 57.880) / 2}{57.880 / 59} = 8.899948$$

Al comparar a un nivel de significancia $\alpha = 0.05$, F_0 con $f_{0.05, 2, 59} = 3.153123$. Con valor $P = 4.1878942 \times 10^{-4}$ con un nivel de significancia del 5%, y el valor p obtenido es pequeño. Por lo tanto, la evidencia sugiere que se debe rechazar la hipótesis nula H_0 . Por tanto concluimos que no se puede descartar este subconjunto de datos del modelo.

La conclusión de la PH es que son significativas y luego en sí

3. Pregunta 3

Plantee una pregunta donde su solución implique el uso exclusivo de una prueba de hipótesis lineal general de la forma $H_0: L\beta = 0$ (solo se puede usar este procedimiento y no SSextra). Especifique claramente la matriz L , el modelo reducido y la expresión para el estadístico de prueba (no hay que calcularlo).

3.1. Prueba de hipótesis y prueba de hipótesis matricial

Se plantea la siguiente prueba de hipótesis:

$$\begin{cases} H_0: \beta_1 = \beta_5, \beta_2 = \beta_3 \\ H_a: \text{Alguna de las desigualdades no se cumple} \end{cases}$$

Reescribiendo matricialmente:

$$\begin{cases} H_0: L\beta = 0 \\ H_a: L\beta \neq 0 \end{cases}$$

Donde L está dada por:

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

Donde el modelo reducido está dado por:

$$RI = \beta_0 + \beta_1(DES_i + ENFP_i) + \beta_2(RC_i + NCP_i) + \beta_4CPD_i + \varepsilon_i, \quad \varepsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

3.2. Estadístico de prueba

El estadístico de prueba F_0 está dado por:

$$F_0 = \frac{(SSE(MR) - SSE(MF)) / 2}{MSE(MF)} \stackrel{H_0}{\sim} f_{2, 59}$$

Obteniendo esto podemos definir la región de rechazo de la hipótesis nula como $F_0 > F_{0.05, 2, 59} = 3.153123$ y con valor $p: P(F_{2, 59} > |F_0|)$

4. Pregunta 4

Realice una validación de los supuestos en los errores y examine si hay valores atípicos, de balanceo e influencias. Qué puede decir acerca de la validez de éste modelo?. Argumente su respuesta.

4.1. Supuestos del modelo

4.1.1. Normalidad de los residuales

2 pt

Para la validación de este supuesto, se plantea la siguiente prueba de hipótesis (shapiro-wilk)

$$\begin{cases} H_0 : \varepsilon_i \sim N(\mu, \sigma^2) \\ H_a : \varepsilon_i \not\sim N(\mu, \sigma^2) \end{cases}$$

no se está probando media cte μ y var cte σ^2

acompañado de un gráfico cuantil-cuantil:

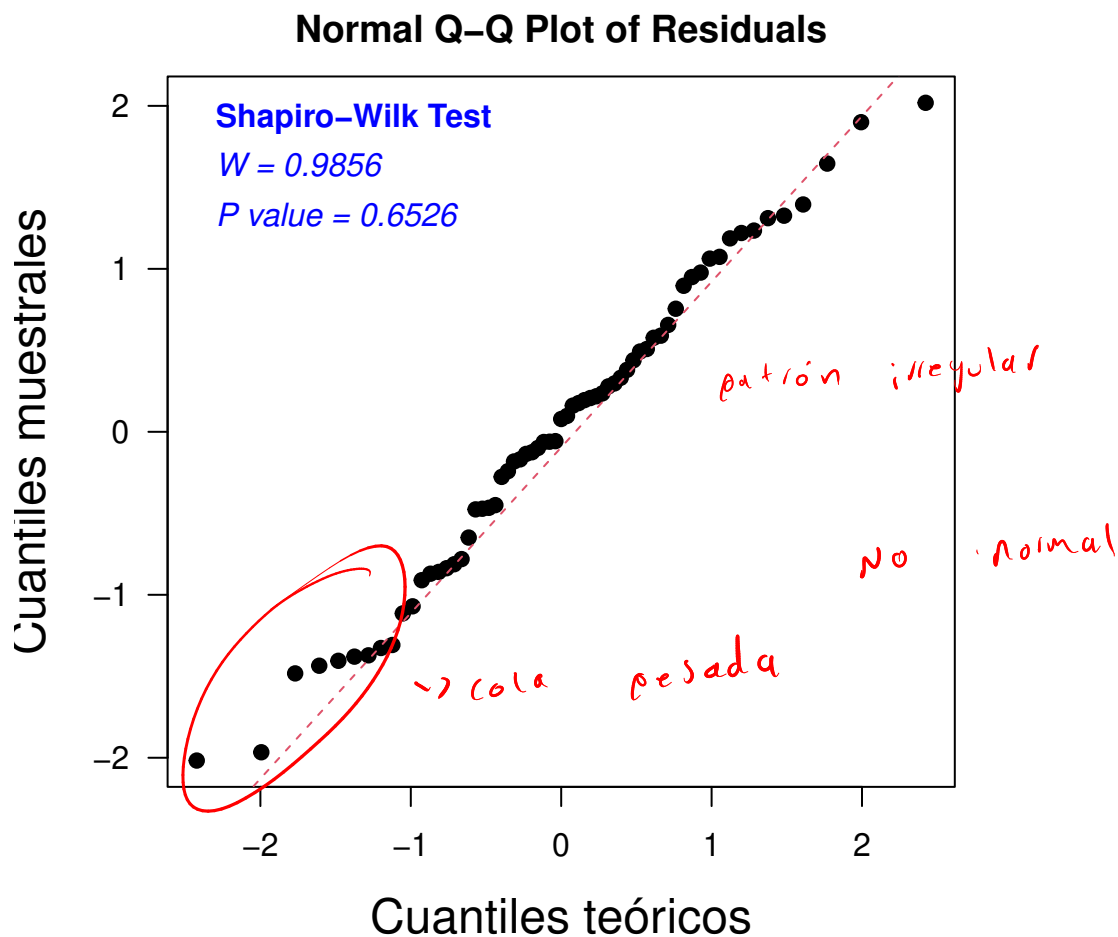


Figura 1: Gráfico cuantil-cuantil y normalidad de los residuales

Ya que el valor de P es alto, se puede inferir que no hay suficiente evidencia para rechazar la hipótesis nula H_0 . Por lo tanto, se puede concluir que el modelo es congruente con la asunción de que los residuos se distribuyen normalmente.

Es más importante el análisis gráfico y no lo hicieron

4.1.2. Media 0 y Varianza constante

2 pt

En esta prueba se quiere probar

$$H_0 : V[\varepsilon_i] = \sigma^2 \quad \text{vs} \quad V[\varepsilon_i] \neq \sigma^2$$



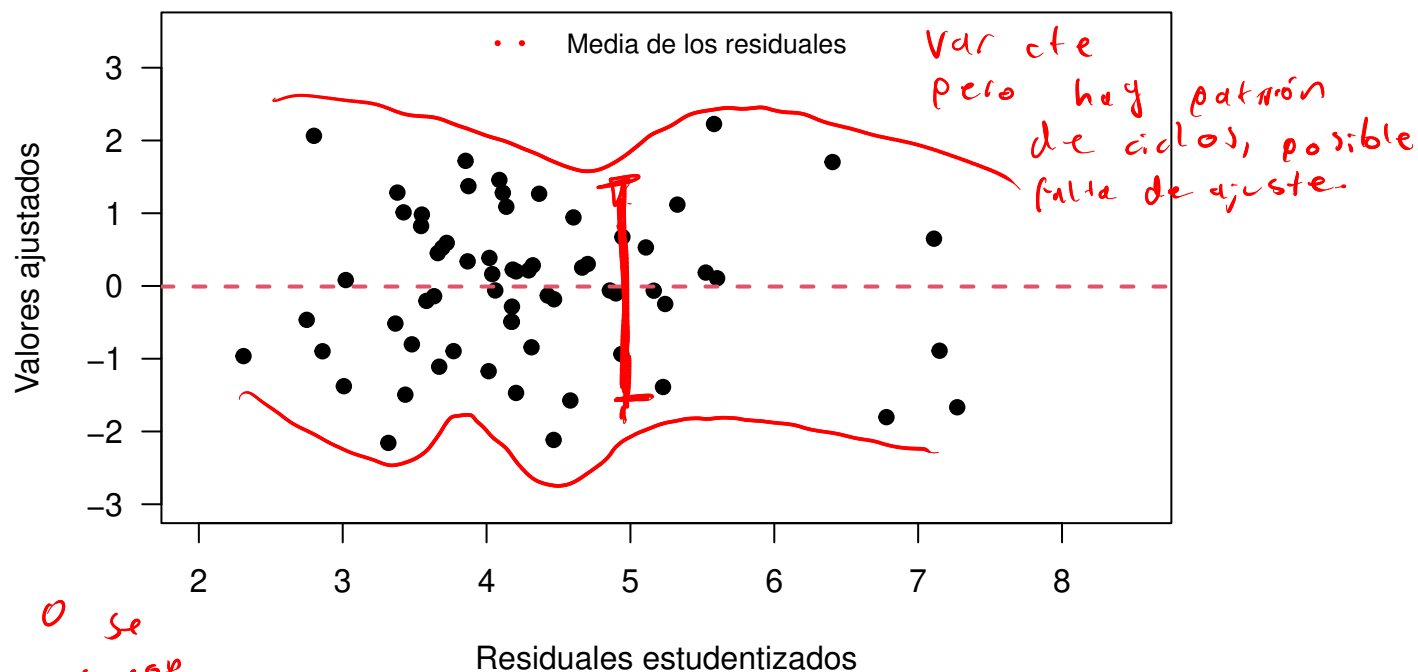


Figura 2: Gráfico residuales estudentizados vs valores ajustados

Podemos notar que la línea punteada roja, que muestra la media de los errores, está en cero o cerca de cero. Esto sugiere que los errores tienen una media cercana a cero. Al analizar los residuos, no se observa ningún patrón discernible, lo que indica que la varianza de los errores es constante en todo el rango de los valores observados. En resumen, la media cercana a cero de los errores y la constancia de la varianza de los residuos sugieren que el modelo se ajusta bien a los datos y cumple con los supuestos básicos de la regresión lineal.

→ No cualquier patrón es por var no cte

mejor o se prueba es con residuales crudos

4.2. Observaciones extremas

4.2.1. Datos atípicos

2,5 p+

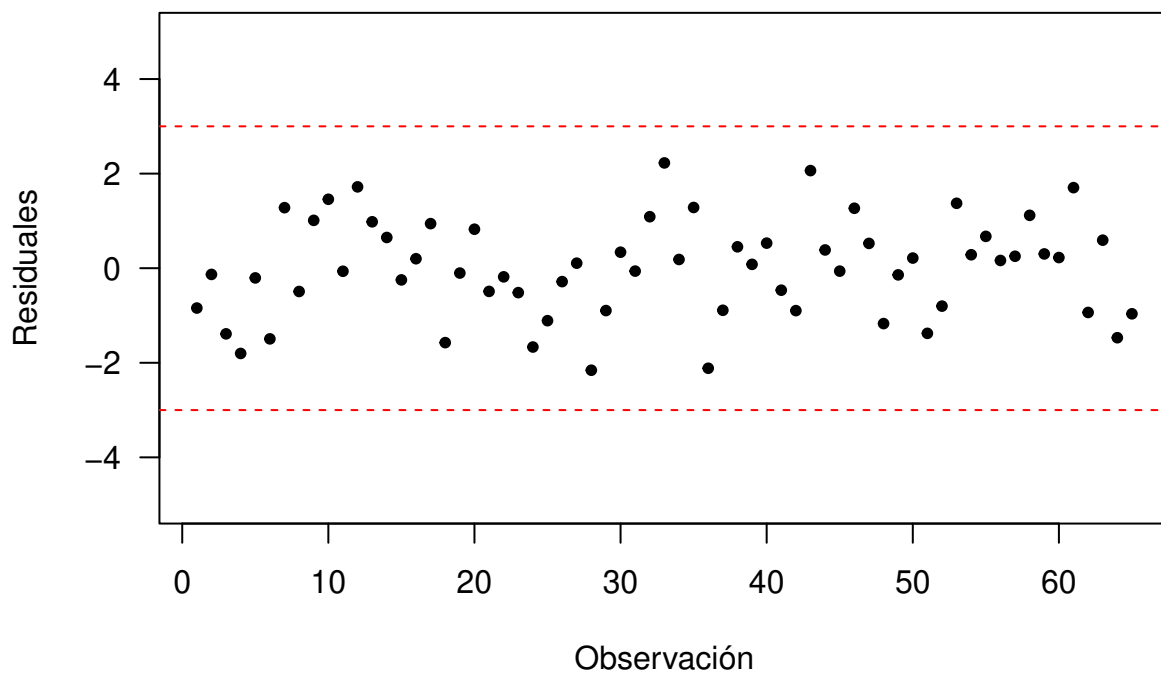


Figura 3: Identificación de datos atípicos

Notese que segun este criterio no existen puntos atipicos que deban ser investigados

¿Qué dice el criterio?

4.2.2. Puntos de balanceo

2,5 pr

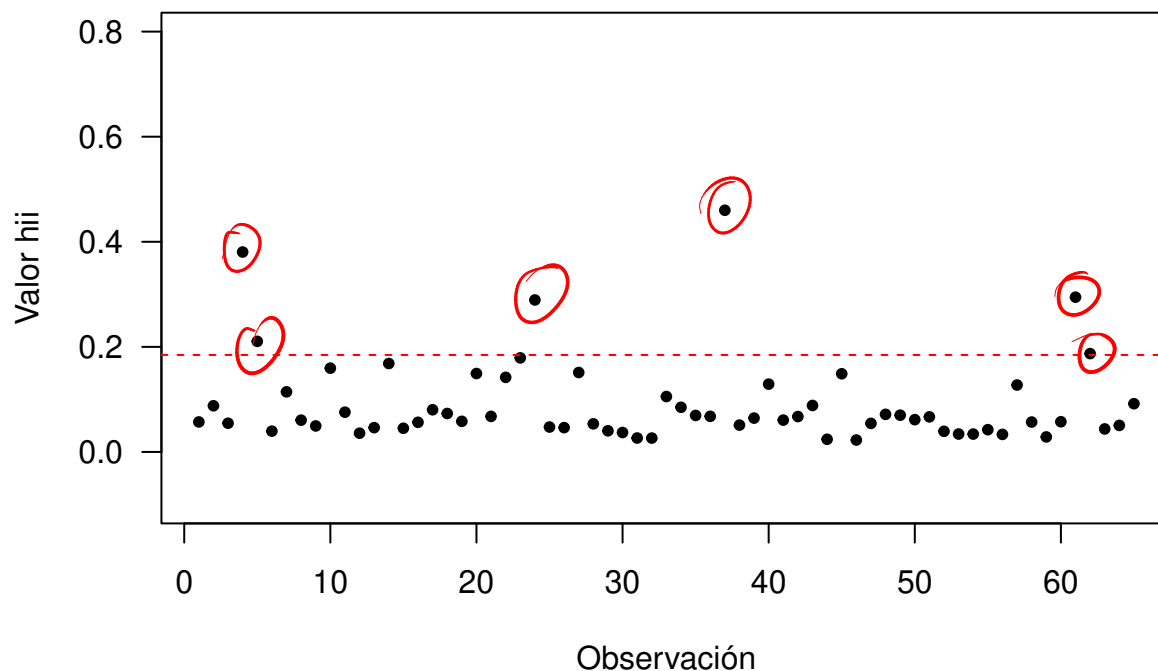


Figura 4: Identificación de puntos de balanceo

Tabla 5: Tabla de puntos de Balanceo

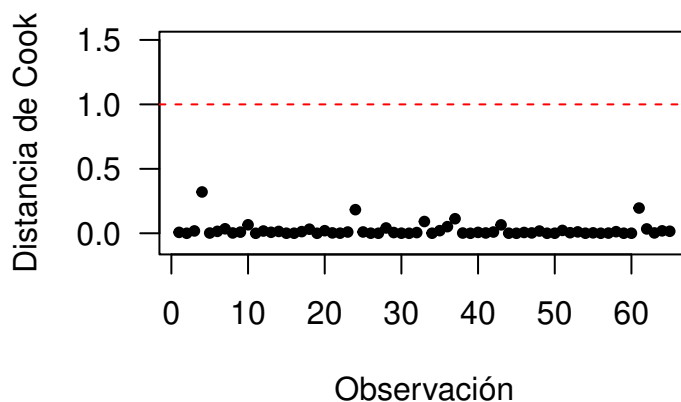
	Errores Estudentizados	D.Cook	Valor hii	DFFITS
4	-1.8031	0.3206	0.3805	-1.4131
5	-0.2050	0.0019	0.2105	-0.1058
24	-1.6672	0.1831	0.2893	-1.0638
37	-0.8901	0.1129	0.4600	-0.8215
61	1.7038	0.1958	0.2946	1.1012
62	-0.9362	0.0337	0.1872	-0.4493

Es importante destacar que hay seis datos que deben ser cuidadosamente analizados en términos de su impacto en el ajuste del modelo y sus propiedades. Estos datos corresponden a los puntos 4, 5, 24, 37, 61 y 62, ya que son mayores que el valor crítico $\frac{2p}{n}$. Estos puntos de balanceo pueden tener una gran influencia en el modelo y, por lo tanto, es crucial llevar a cabo un análisis detallado de su impacto antes de sacar conclusiones definitivas.

¿En qué influyen específicamente?

4.2.3. Puntos influyentes

Bajo el criterio de Cook, se hace la siguiente gráfica:



2 pr

Figura 5: Criterio distancias de Cook para puntos influyentes

Bajo el criterio de cook, se obtuvo la anterior gráfica. A partir de la gráfica podemos concluir que no existen puntos influyentes bajo este criterio

↓
redundantes

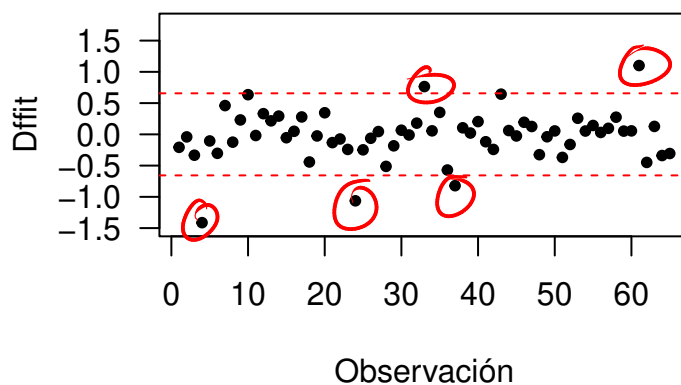


Figura 6: Criterio Dffits para puntos influyentes

Tabla 6: Tabla del criterio DFFITS para encontrar puntos influyentes

	Errores Estudentizados	D.Cook	Valor hii	DFFITS
4	-1.8031	0.3206	0.3805	-1.4131
24	-1.6672	0.1831	0.2893	-1.0638
33	2.2262	0.0913	0.1055	0.7644
37	-0.8901	0.1129	0.4600	-0.8215
61	1.7038	0.1958	0.2946	1.1012

Usando el criterio de Dffits, se ha generado el gráfico anterior, el cual sugiere que hay varios valores influyentes en el modelo. Específicamente, las observaciones 4, 24, 33, 37 y 61 parecen tener un impacto significativo en el modelo y deben ser investigadas con más detalle. Es necesario realizar un análisis adicional para determinar si estos valores influyentes deben ser eliminados o ajustados de alguna manera y evaluar su impacto en el modelo de regresión.

4.3. Conclusiones

El modelo cumple con los supuestos básicos de la regresión lineal, es decir, que la media de los residuos es cercana a cero y la varianza es constante. Sin embargo, se observa un gran número de datos de balanceo e influenciadores en el modelo, lo que indica la necesidad de investigar si estos datos están afectando significativamente el modelo y sus supuestos, incluyendo la normalidad de los residuos. En conclusión, aunque el modelo cumple con los supuestos básicos, no es adecuado para hacer predicciones y se debe llevar a cabo un análisis adicional para evaluar la influencia de los datos de balanceo e influyentes y determinar si el modelo es una buena estimación.

¿Por qué no si cumple los supuestos? No significa decirlo si era válido según lo que encontraron.