

3,45

Trabajo 1

Estudiantes

**Rivera Bedoya Marhia Camila
Sanchez Vivas Harold Jhoan
Valderrama Posada Mateo
Tora Arroyave Oscar Julian**

Docente

Francisco Javier Rodriguez Cortes

Asignatura

Estadística II



UNIVERSIDAD
NACIONAL
DE COLOMBIA

Sede Medellín
25 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	9
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	11

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

15 pt

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 **Equipo45.txt**, las variables para el modelo son

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

X1 DE 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 ENF Número de enfermeras: Promedio de enfermeras, equivalentes a tiempo completo, durante el periodo del estudio.

El modelo que se propone es:

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

1.1. Modelo de regresión

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

Tabla 1: Tabla de valores de los coeficientes estimados

	Valor del parámetro
$\hat{\beta}_0$	-1.68786
$\hat{\beta}_1$	0.19317
$\hat{\beta}_2$	0.03430
$\hat{\beta}_3$	0.04258
$\hat{\beta}_4$	0.01994
$\hat{\beta}_5$	0.00069

2 pt

no va en ec. ajustada

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

$$\widehat{RI}_i = -1.68786 + 0.19317 DE_i + 0.0343 RC_i + 0.04258 NC_i + 0.01994 CPD_i + 6.9 \times 10^{-4} NENF_i + \varepsilon_i, \quad \varepsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

Donde las variables se mueven de acuerdo $1 \leq i \leq 60$

1.2. Significancia de la regresión

2,5 pt

Se plantea el siguiente Juego de Hipótesis

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

¿Es estadístico de prueba?

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	67.3216	5	13.46431	11.8112	9.69569e-08
Error	61.5578	54	1.13996		

Los resultados obtenidos de la Tabla Anova indican que la hipótesis nula debe ser rechazada de lo cual podemos concluir que el modelo al menos alguna de las variables es significativa

¿es significativo?

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	-1.6879	1.8239	-0.9254	0.3589
β_1	0.1932	0.0846	2.2832	0.0264
β_2	0.0343	0.0354	0.9698	0.3365
β_3	0.0426	0.0147	2.8922	0.0055
β_4	0.0199	0.0082	2.4384	0.0181
β_5	0.0007	0.0008	0.8141	0.4192

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:

- $\hat{\beta}_1 = 0.19317$: Si todas las demás variables predictoras se mantienen iguales, un aumento de un día en la Duración de la estancia en el hospital daría como resultado un aumento esperado en el promedio del Riesgo de infección en un porcentaje determinado por el valor de 0.19317%.
- $\hat{\beta}_3 = 0.04258$: 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 se incremente en un 0.04258%.

¿es más corrección?

constantes y no iguales

19,3117%

¿la probabilidad promedio

misma corrección
↑

- $\hat{\beta}_4 = 0.01994$: 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 aumenta en un 0.01994 %

1.5. Coeficiente de determinación múltiple R^2

3 pt

El modelo tiene un R^2 de 0.5224 lo cual significa que aproximadamente el 52.24 % 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 notar que las variables que contribuyen significativamente son la Duración de la estadía en el hospital, el Censo promedio diario de pacientes en el hospital y el número de camas. Esto se refleja en la importancia de los parámetros.

2. Pregunta 2

3 pt

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=0.3365, β_5 con VP= 0.4192, β_1 con VP= 0.02634. Por tanto, se plantea la siguiente prueba de hipótesis:

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

Hay 2 equipos con este mismo error tan particular.

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

$$\text{MR: } Rinf_i = \beta_0 + \beta_3 NCP_i + \beta_4 PD4_i + \varepsilon_i, \quad \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

Así se llama?

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	61.558	X1 X2 X3 X4 X5
Modelo reducido	73.462	X3 X4

✓

Así no se llaman sus variables

2.2. Estadístico de prueba y conclusiones

Se construye el estadístico de prueba como:

$$F_0 = \frac{(SSR(\beta_0, \beta_3, \beta_4) - SSR(\beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5) / 2)}{MSE(MF)} \stackrel{H_0}{\sim} F_{2, 54}$$

→ $\beta_0, \beta_3, \beta_4$

se está probando sobre $\beta_1, \beta_2, \beta_5$

0.3

1,5 pt

$$F_0 = \frac{(SSE(\beta_0, \beta_3, \beta_4) - SSE(\beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5))/2}{MSE(\beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5)} \stackrel{H_0}{\sim} f_{2,54}$$

$$= \frac{(73.462 - 61.558)/2}{61.558/54} = 5.2212$$

Cuando comparamos F_0 con $f_{0.05,2,54} = 3.168246$ a un nivel de significancia de $\alpha = 0.05$, y usamos un valor p de 0.0084527, vemos que el valor p es pequeño, lo que sugiere que debemos rechazar la hipótesis nula H_0 . Por lo tanto, llegamos a la conclusión de que no se puede descartar este subconjunto de datos del modelo.

Es esa no es la conclusión inicial, es una consecuencia

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:

$$\begin{cases} H_0 : \beta_4 = \beta_5, \beta_2 = \beta_3 \\ H_a : \text{Alguna de las desigualdades no se cumple} \end{cases}$$

Reescribiendo matricialmente:

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

Donde L está dada por:

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

Donde el modelo reducido está dado por:

$$RI = \beta_0 + \beta_1(DES_i) + \beta_2(RC_i + NCP_i) + \beta_4(CPD_i + ENF_i) + \varepsilon_i, \quad \varepsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

3.2. Estadístico de prueba

El estadístico de prueba F_0 está dado por:

$$F_0 = \frac{(SSE(MR) - SSE(MF))/2}{MSE(MF)} \stackrel{H_0}{\sim} f_{2,54}$$

Obteniendo esto podemos definir la region de rechazo de la hipotesis nula como $F_0 > F_{0.05,2,54} = 3.168246$ y con valor p: $P(F_{2,54} > |F_0|)$

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

2pt

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 están probando media de μ y var de σ^2

acompañado de un gráfico cuantil-cuantil:

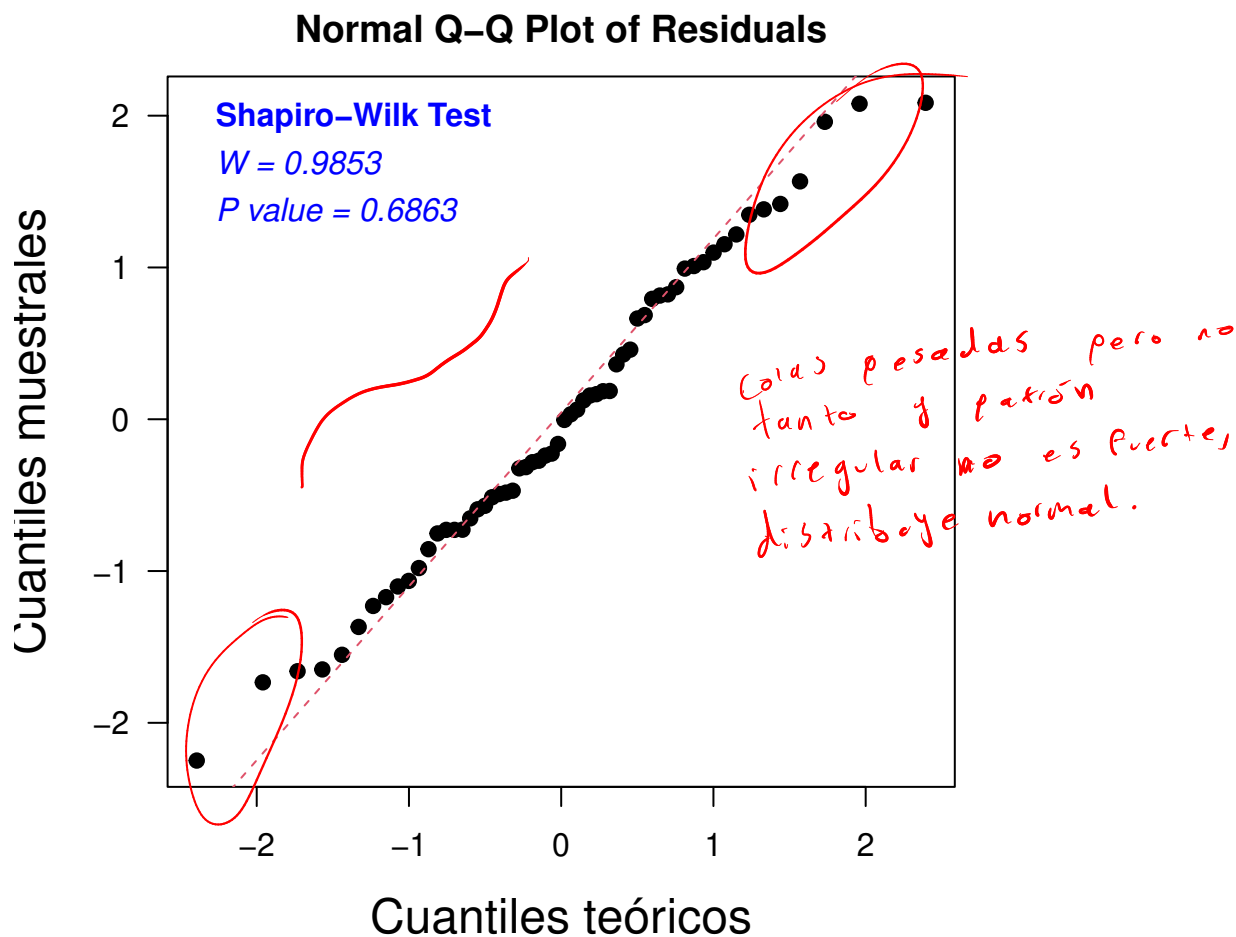


Figura 1: Gráfico cuantil-cuantil y normalidad de los residuales

Dado que el valor de p es alto, se puede concluir que no hay suficiente evidencia para rechazar la hipótesis nula H_0 . Por lo tanto, se puede inferir que el modelo es consistente con la suposición de que los residuales siguen una distribución normal.

No hacen análisis gráfico que es más importante errores

4.1.2. Media 0 y Varianza constante

3pt

En esta prueba se quiere probar

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



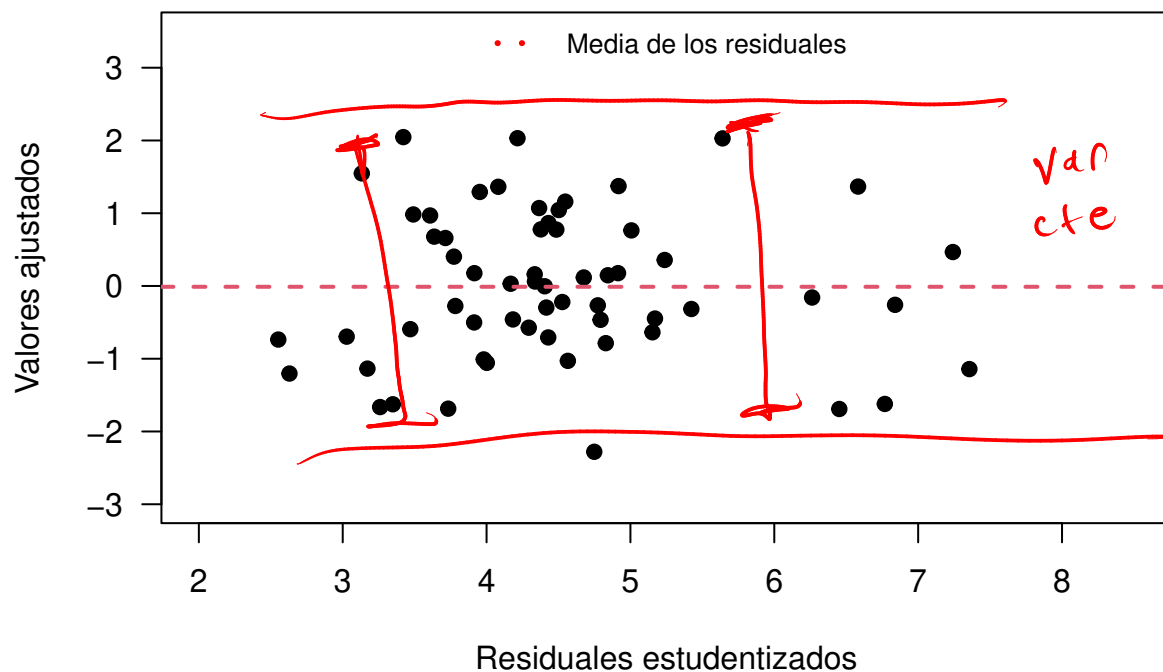


Figura 2: Gráfico residuales estudentizados vs valores ajustados

Se puede observar que la línea punteada roja, que representa la media de los errores, se encuentra cerca o en cero, lo que sugiere que los errores tienen una media cercana a cero. Además, al examinar los residuos, no se puede detectar ningún patrón y se ven uniformemente distribuidos, lo que indica que la varianza de los errores es constante a través de todo el rango de los valores observados. ✓

4.2. Observaciones extremas

4.2.1. Datos atípicos

1,5 pt

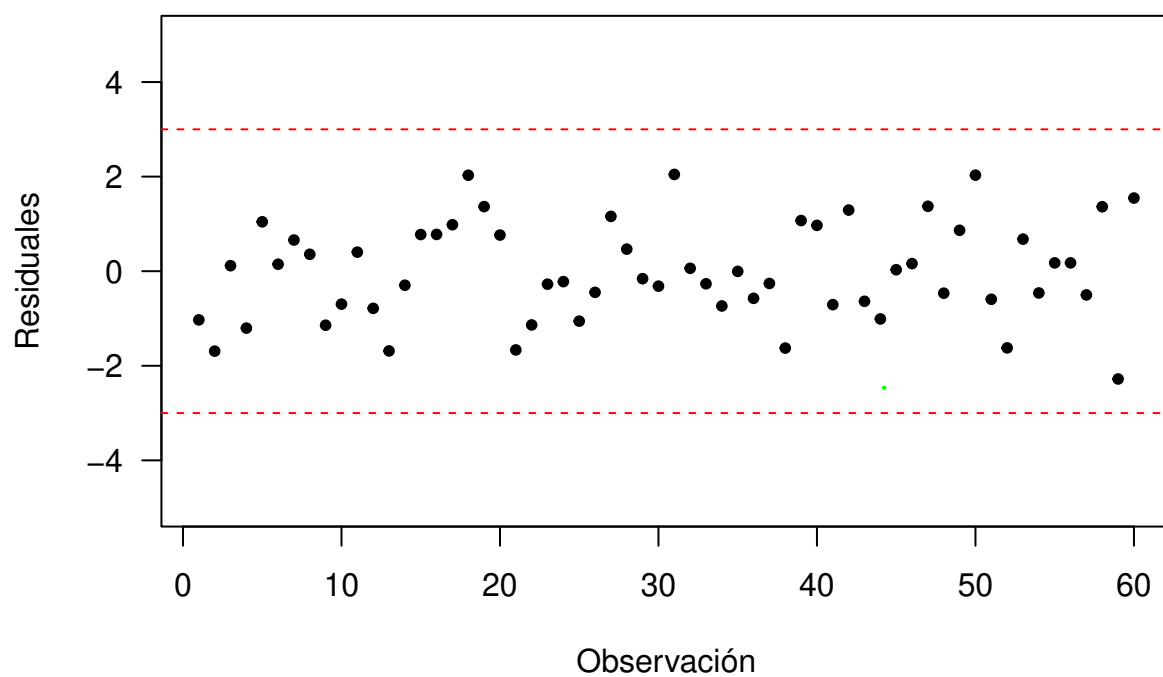


Figura 3: Identificación de datos atípicos

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

Segun cual criterio? qué dice?

4.2.2. Puntos de balanceo

2, 5 pt

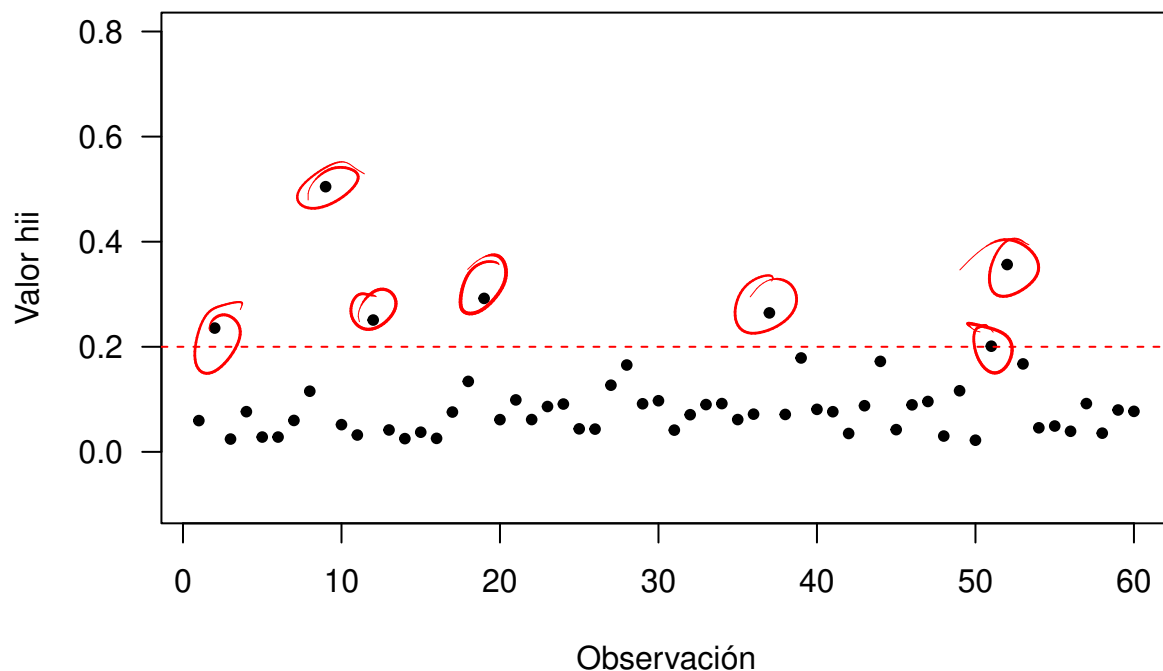


Figura 4: Identificación de puntos de balanceo

Tabla 5: Tabla de puntos de Balanceo

	Errores Estudentizados	D.Cook	Valor hii	DFFITS
2	-1.6907	0.1419	0.2356	-0.9385
9	-1.1425	0.2204	0.5047	-1.1533
12	-0.7858	0.0347	0.2509	-0.4547
19	1.3664	0.1264	0.2921	0.8777
37	-0.2592	0.0041	0.2645	-0.1555
51	-0.5936	0.0150	0.2011	-0.2978
52	-1.6214	0.2356	0.3565	-1.2067

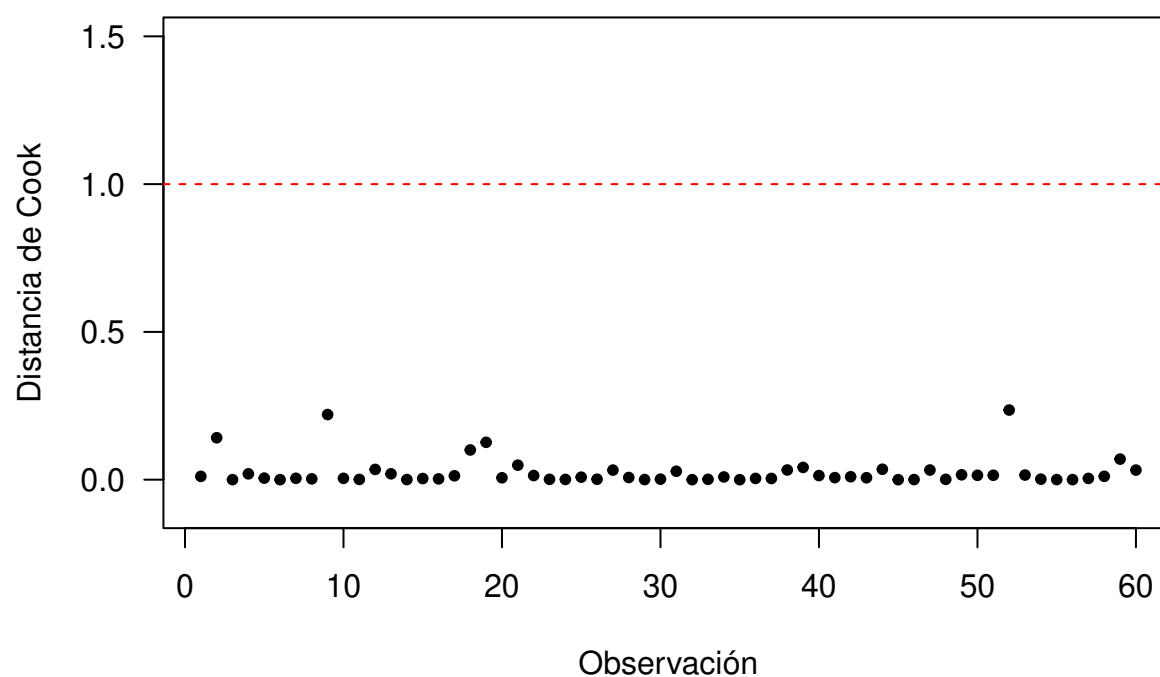
Es importante destacar que hay 7 datos que deben ser estudiados en términos de su impacto en el ajuste del modelo y sus propiedades. Estos datos corresponden a los puntos 2, 9, 12, 19, 37, 51 y 52, ya que son mayores al criterio $\frac{22}{n}$.

¿cuánto da?

¿cómo qué?
¿qué causan?

4.2.3. Puntos influenciales

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



redundantes

1,5 pt

Figura 5: Criterio distancias de Cook para puntos influenciales

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

¿Qué dice el criterio?

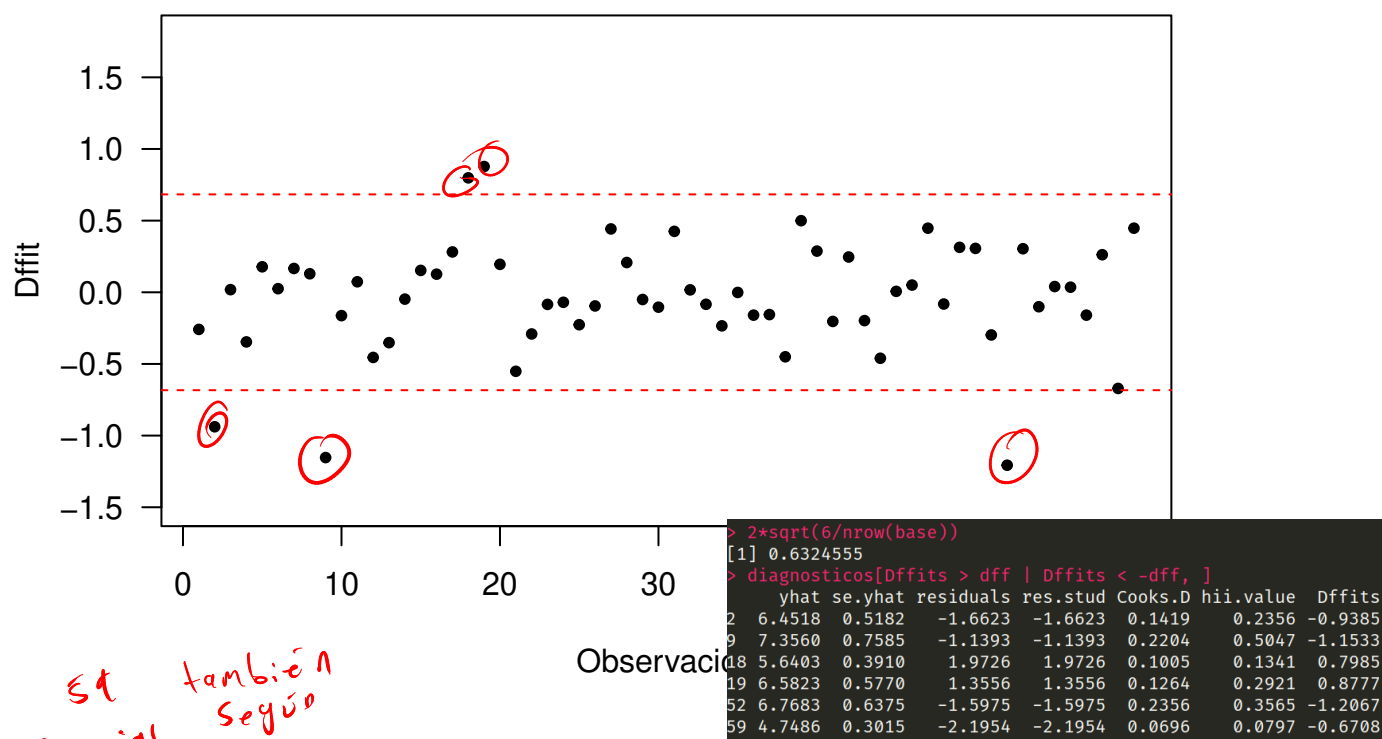


Figura 6: Criterio Dffits para puntos influyentes

Tabla 6: Tabla del criterio DFFITS para encontrar puntos influyentes

	Errores Estudentizados	D.Cook	Valor hii	DFFITS
2	-1.6907	0.1419	0.2356	-0.9385
9	-1.1425	0.2204	0.5047	-1.1533
18	2.0287	0.1005	0.1341	0.7985
19	1.3664	0.1264	0.2921	0.8777
52	-1.6214	0.2356	0.3565	-1.2067

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 2, 9, 18, 19 y 62 que pueden tener un impacto significativo en el modelo y deben ser investigadas con más detalle.

¿Qué causan?

¡No coincide con tabla

4.3. Conclusiones

El modelo cumple con los supuestos básicos de regresión lineal de que la media de los residuos es cercana a cero y la varianza es constante. Sin embargo, se observó una gran cantidad de datos de balanceo y puntos influyentes, lo que sugiere la necesidad de investigar si estos datos afectan significativamente el modelo y sus supuestos, incluida la normalidad de los residuos. En resumen, aunque el modelo cumple con los supuestos subyacentes, no podemos decir si el modelo es adecuado para hacer predicciones, y se requieren análisis

adicionales para evaluar el impacto de los datos de equilibrio e influencia para determinar si el modelo es un modelo preciso

Modelo válido o no?