

Trabajo 1

3,9

Estudiantes

Brayan Alberto Patiño Alzate
Jenifer Tatiana Atehortua Duque
Samuel David Suarez Tovar
Felipe Andrés Villero Mejia
Equipo # 65

Docente

Javier Armando Lozano Rodriguez

Asignatura

Estadística II



UNIVERSIDAD
NACIONAL
DE COLOMBIA

Sede Medellín
5 de Octubre 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	4
1.5. Coeficiente de determinación múltiple R^2	5
2. Pregunta 2	5
2.1. Planteamiento pruebas de hipótesis y modelo reducido	5
2.2. Estadístico de prueba y conclusión	5
3. Pregunta 3	6
3.1. Prueba de hipótesis y prueba de hipótesis matricial	6
3.2. Estadístico de prueba	6
4. Pregunta 4	6
4.1. Supuestos del modelo	6
4.1.1. Normalidad de los residuales	7
4.1.2. Varianza constante	8
4.2. Verificación de las observaciones	9
4.2.1. Datos atípicos	9
4.2.2. Puntos de balanceo	10
4.2.3. Puntos influyentes	11
4.3. Conclusión	12

Índice de figuras

1.	Gráfico cuantil-cuantil y normalidad de residuales	7
2.	Gráfico residuales estudentizados vs valores ajustados	8
3.	Identificación de datos atípicos	9
4.	Identificación de puntos de balanceo	10
5.	Criterio distancias de Cook para puntos influenciales	11
6.	Criterio Dffits para puntos influenciales	12

Índice de tablas

1.	Tabla de valores coeficientes del modelo	3
2.	Tabla ANOVA para el modelo	4
3.	Resumen de los coeficientes	4
4.	Resumen tabla de todas las regresiones	5
5.	Resumen de diagnostico	10
6.	Resumen de diagnostico	12

1. Pregunta 1

20pt

Teniendo en cuenta la base de datos asignada, se realiza el ajuste del modelo de regresión lineal múltiple (RLM), explicando la eficacia sobre el control de infecciones hospitalarias en la cual hay 5 variables regresoras dadas por:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + \beta_4 X_{4i} + \beta_5 X_{5i} + \varepsilon_i, \varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2); 1 \leq i \leq 54$$

Se asignan las siguientes variables:

$$\begin{cases} Y : \text{Riesgo de Infección} \\ X1 : \text{Duración de estadia} \\ X2 : \text{Rutina de cultivos} \\ X3 : \text{Número de camas} \\ X4 : \text{Censo promedio diario} \\ X5 : \text{Número de enfermeras} \end{cases}$$

1.1. Modelo de regresión

Al ajustar el modelo, se obtienen los siguientes coeficientes:

Tabla 1: Tabla de valores coeficientes del modelo

	Valor del parámetro
β_0	-0.7722
β_1	0.1952
β_2	0.0255
β_3	0.0552
β_4	0.0076
β_5	0.0016

30pt

Por lo tanto, el modelo de regresión ajustado es:

$$\hat{Y}_i = -0.7722 + 0.1952X_{1i} + 0.0255X_{2i} + 0.0552X_{3i} + 0.0076X_{4i} + 0.0016X_{5i}$$

1.2. Significancia de la regresión

Para analizar la significancia de la regresión, se plantea la siguiente prueba de hipótesis:

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

Cuyo estadístico de prueba es:

$$F_0 = \frac{MSR}{MSE} \stackrel{H_0}{\sim} f_{5,48} \quad (1)$$

$$F_0 = \frac{10.82931}{1.01425} = 10.6772 \quad (2)$$

Ahora, se presenta la tabla Anova:

Tabla 2: Tabla ANOVA para el modelo

	Sumas de cuadrados	g.l.	Cuadrado medio	F_0	P-valor
Regresión	54.1466	5	10.82931	10.6772	6.18567e-07
Error	48.6838	48	1.01425		

5pt

De la tabla Anova, se observa un valor P aproximadamente igual a 0, por lo que se rechaza la hipótesis nula en la que $\beta_j = 0$ con $1 \leq j \leq 5$, aceptando la hipótesis alternativa en la que algún $\beta_j \neq 0$, en este caso el P-valor nos permite concluir que sí, el modelo de regresión es significativo, rechazando la hipótesis nula planteada anteriormente, esto quiere decir que alguna de las variables predictorias es significativa en el riesgo de adquirir infección en hospitales.

1.3. Significancia de los parámetros

En el siguiente cuadro respecto a los parametros individuales se presenta información de los mismos, lo cual permitirá determinar cuáles de ellos son significativos.

Tabla 3: Resumen de los coeficientes

	$\hat{\beta}_j$	$SE(\hat{\beta}_j)$	T_{0j}	P-valor
β_0	-0.7722	1.9008	-0.4063	0.6864
β_1	0.1952	0.0905	2.1575	0.0360
β_2	0.0255	0.0351	0.7247	0.4721
β_3	0.0552	0.0145	3.8224	0.0004
β_4	0.0076	0.0091	0.8416	0.4042
β_5	0.0016	0.0007	2.1533	0.0364

6pt

Los P-valores presentes en la tabla permiten concluir que con un nivel de significancia $\alpha = 0.05$, los parámetros β_1 β_3 y β_5 son significativos, Los P valores mostrados en la tabla me indican cuales de los parámetros está aportando significativamente al modelo de regresión, ya que sus P-valores 0.0360 y 0.0004 y 0.0364 respectivamente son menores al valor alfa.

1.4. Interpretación de los parámetros

3pt

Con lo anterior se puede determinar que $\hat{\beta}_1$: 0.0360, indica que por cada unidad que aumente el promedio de duración de la estadia en el hospital X_1 , el riesgo de infección aumenta 0.0360 unidades, cuando las demás variables predictorias se mantienen fijas.

Asi mismo $\hat{\beta}_3$: 0.0004 indica que por cada unidad que aumente el promedio de camas en el hospital X_3 , el riesgo de infección aumenta 0.0004 unidades, cuando las demás variables predictorias se mantienen fijas.

Y por ultimo $\hat{\beta}_5$: 0.0364 indica que por cada unidad que aumente el promedio de enfermeras en el hospital X_5 , el riesgo de infección aumenta 0.0364 unidades, cuando las demás variables predictorias se mantienen fijas.

1.5. Coeficiente de determinación múltiple R^2

3pt

$$R^2 = \frac{SSR}{SSR + SSE} \quad (3)$$

$$R^2 = \frac{54.1466}{54.1466 + 48.6838} = 0.526562 \quad (4)$$

El modelo tiene un coeficiente de determinación múltiple $R^2 = 0.526562$, lo que significa que aproximadamente el 52.66 % de la variabilidad total observada en los resultados de la prueba de riesgo es explicada por el modelo de regresión propuesto en el presente informe.

2. Pregunta 2

Opt

2.1. Planteamiento pruebas de hipótesis y modelo reducido

mas bajo, y estar 3

Las covariable con el P-valor más alto en el modelo fueron, rutina de cultivos X_2 (Valor $P = 0.4721$) y censo promedio diario X_4 (Valor $P = 0.4042$), ahora se procede a plantear la prueba de hipótesis para verificar la significancia en simultaneo del subconjunto con valores p más altos del modelo de regresión.

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

Tabla 4: Resumen tabla de todas las regresiones

	SSE	Covariables en el modelo
Modelo completo	48.684	X_1, X_2, X_3, X_4, X_5
Modelo reducido	83.101	X_2, X_4

 X_1, X_3, X_5

2.2. Estadístico de prueba y conclusión

Se construye el estadístico de prueba como:

$$\begin{aligned} F_0 &= \frac{(SSE(\beta_0, \beta_3, \beta_5) - SSE(\beta_0, \dots, \beta_5))/2}{MSE(\beta_0, \dots, \beta_5)} \stackrel{H_0}{\sim} f_{2,48} \\ &= \frac{(83.101 - 48.684/2)}{1.01425} \\ &= 17.4537 \end{aligned} \quad (5)$$

Opt

Ahora, comparando el F_0 con $F_{0.95,2,48} = 17.4537$, se puede ver que $F_0 < F_{0.95,2,48}$

Opt

Por consiguiente, con un nivel de significancia de 0.05 comparando F_0 con $F_{0.95,2,48} = 17.4537$, se puede ver que $F_0 < F_{0.95,2,48}$, como F_0 no esta dentro de la region de rechazo, no se rechaza H_0 y se concluye que las variables X_2 y X_4 no afectan de forma conjunta a la probabilidad de riesgo de infección.

No!!

3. Pregunta 3

7 p+

3.1. Prueba de hipótesis y prueba de hipótesis matricial

Se plantea la hipótesis de que la variable (X1) Duración de la estadia es igual de significativa a (X3) número de camas; además (X2) rutina de cultivos y (X4) censo promedio diario son igual de significativas. Por consiguiente se plantea la siguiente prueba de hipótesis:

$$\begin{cases} H_0 : \beta_1 = \beta_3; \beta_2 = \beta_4 \\ H_1 : \text{Alguna de las igualdades no se cumple} \end{cases}$$

reescribiendo matricialmente:

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

Con \mathbf{L} dada por

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

2 p+

El modelo reducido está dado por:

$$Y_i = \beta_0 + \beta_2 X_{2i}^* + \beta_3 X_{3i}^* + \beta_5 X_{5i} + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2); \quad 1 \leq i \leq 54$$

X 0 p+

Donde $X_{2i}^* = X_{2i} + X_{4i}$ y $X_{3i}^* = 3X_{1i} + X_{3i}$

X

3.2. Estadístico de prueba

El estadístico de prueba F_0 está dado por:

$$F_0 = \frac{(SSE(MR) - SSE(MF))/2}{MSE(MF)} \stackrel{H_0}{\sim} f_{2,48} \quad (6)$$

2 p+

$$F_0 = \frac{(SSE(MR) - 48.684/2)}{1.01425} \quad (7)$$

4. Pregunta 4

15 p+

4.1. Supuestos del modelo

Normalidad Tomando en cuenta el valor p no se rechaza H_0 siendo se podría decir que hay normalidad, sin embargo la gráfica nos muestra patrones irregulares, además las colas están alejadas de la línea de tendencia marcada en la gráfica lo que supone que el supuesto de normalidad no se este cumpliendo como lo indica el valor p, por consiguiente al tener mas peso el método grafico frente al valor se rechaza el supuesto de normalidad por medio del grafico de cuanti – cuantil.

Varianza En la gráfica observamos que la varianza no tiene una tendencia marcada hacia la linealidad, o hacia valores constantes, ya que los datos se muestran dispersos a lo largo de la gráfica, se puede decir que no se cumple el supuesto de varianza constante.

4.1.1. Normalidad de los residuales

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

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

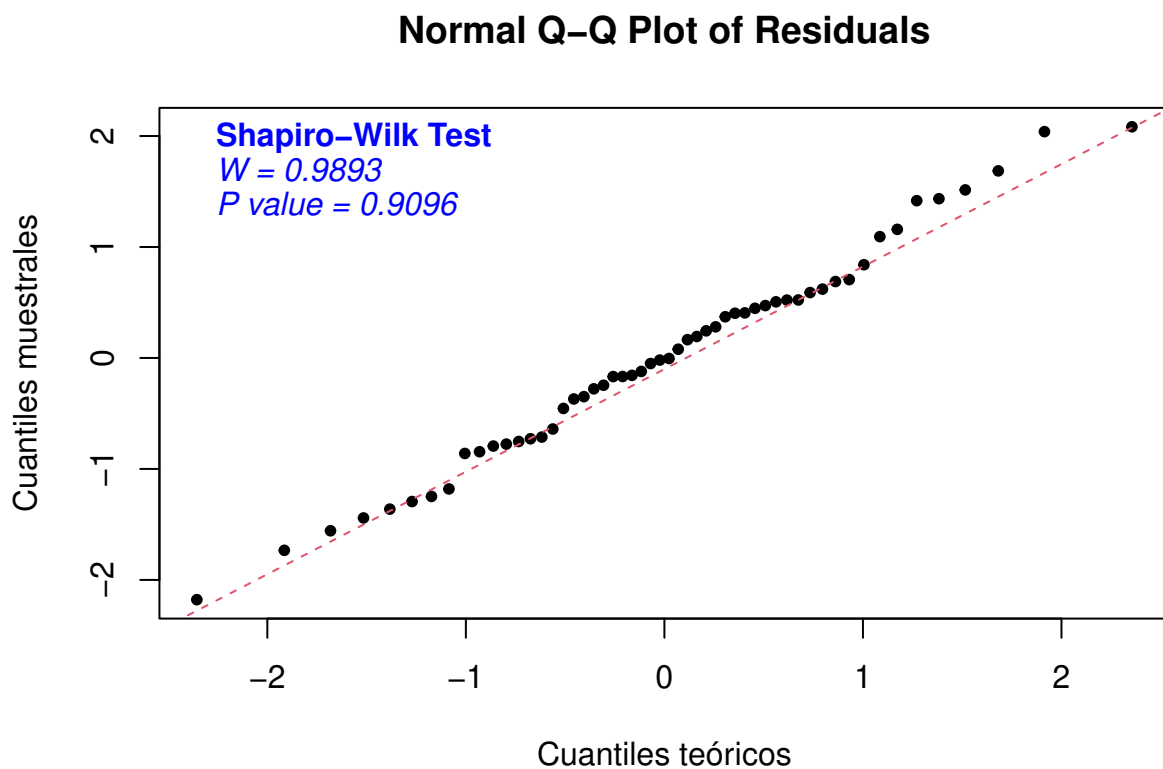


Figura 1: Gráfico cuantil-cuantil y normalidad de residuales

Tomando en cuenta el valor p no se rechaza H_0 , se puede decir que hay normalidad, sin embargo la gráfica nos muestra patrones irregulares en la parte superior de recta, se evidencia que la cola está un poco alejada de la línea de tendencia marcada en la gráfica lo que supone que el supuesto de normalidad no se está cumpliendo como lo indica el valor p. Sin embargo la gráfica de comparación de cuantiles permite ver patrones irregulares que posiblemente se deben a datos atípicos, pero son leves, así que, a partir del gráfico tampoco se rechaza la hipótesis nula.

4.1.2. Varianza constante

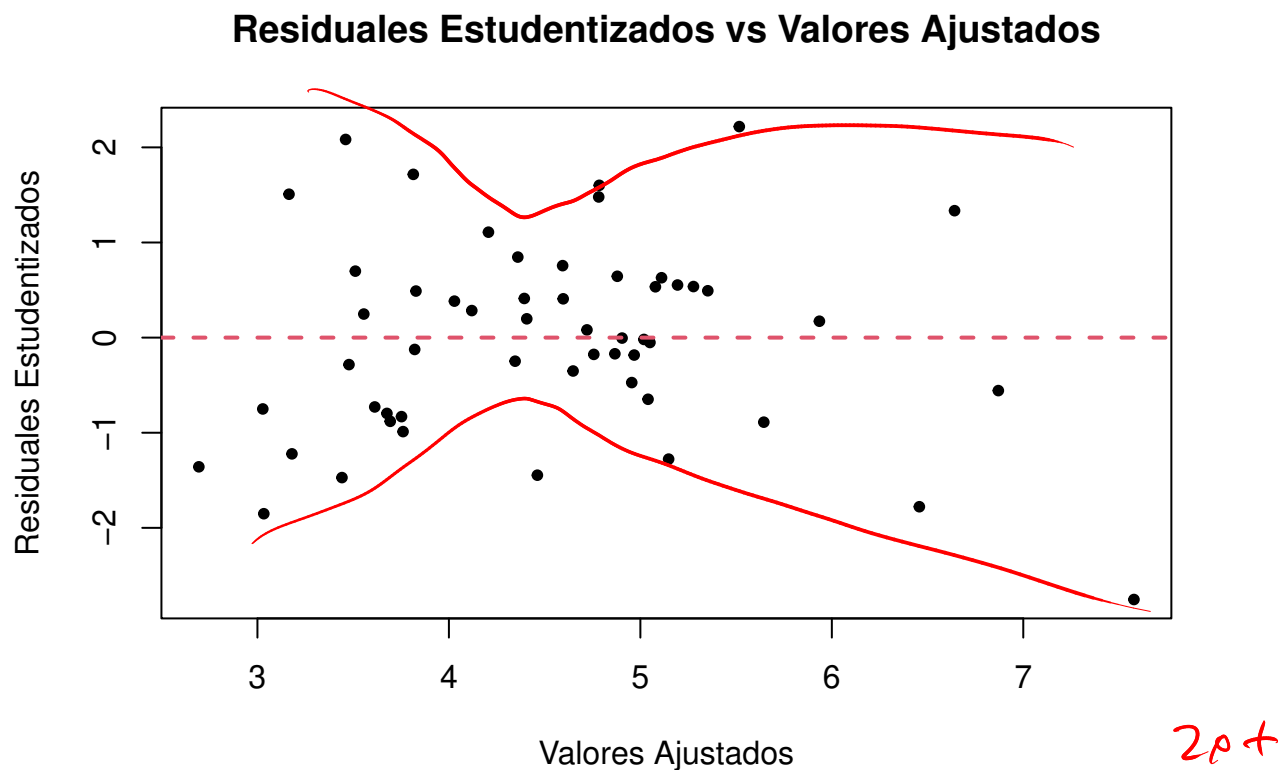


Figura 2: Gráfico residuales estudentizados vs valores ajustados

En el gráfico de residuales estudentizados vs valores ajustados se puede observar que no hay patrones en los que la varianza aumente, decrezca 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. Aunque La Varianza efectivamente no se cumple pero el análisis hecho, no explica el hecho expuesto lo suficientemente claro.

si hay patrón

4.2. Verificación de las observaciones

4.2.1. Datos atípicos

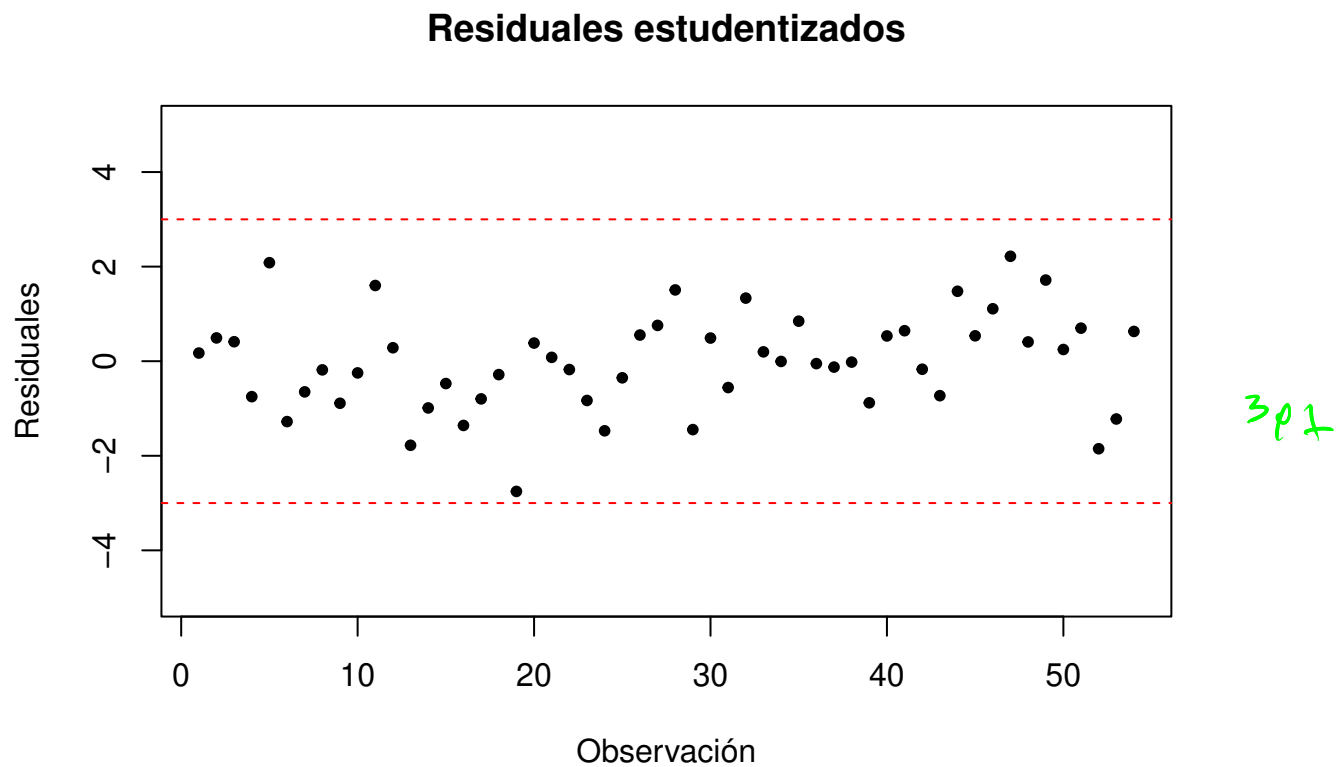
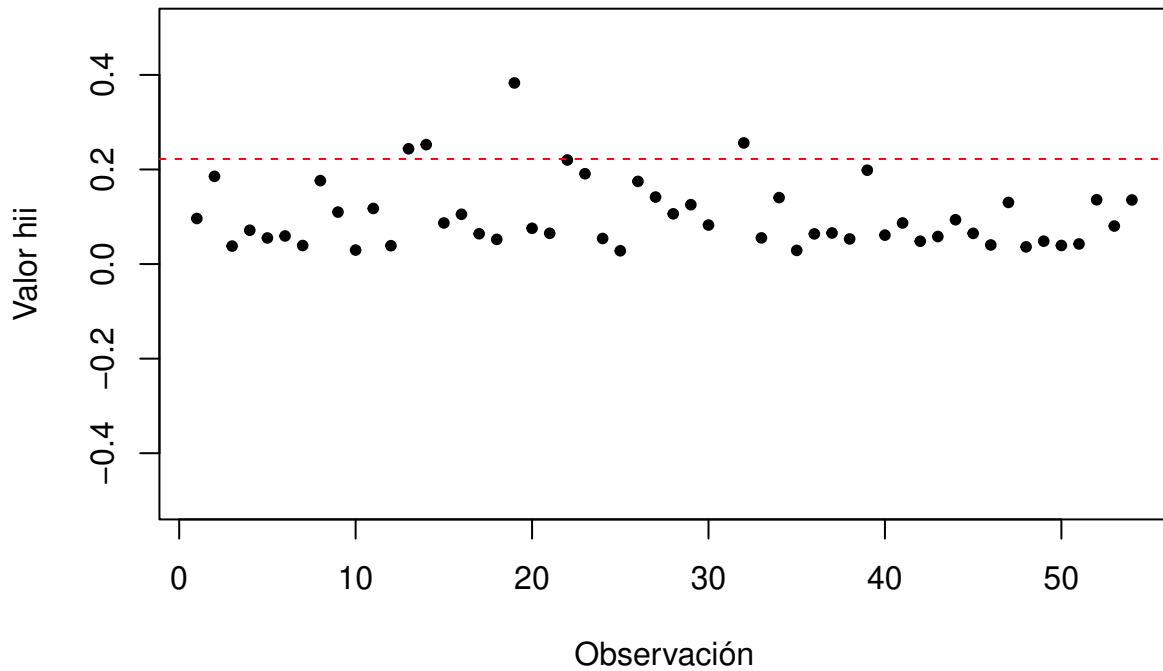


Figura 3: Identificación de datos atípicos

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

4.2.2. Puntos de balanceo

Gráfica de h_{ii} para las observaciones

causan...?

Figura 4: Identificación de puntos de balanceo

Tabla 5: Resumen de diagnostico

	<i>Res.stud</i>	<i>Cook.D_i</i>	<i>h_{ii}value</i>	<i>Df fits</i>
13	-1.7782	0.1699	0.2438	-1.0338
14	-0.9884	0.0550	0.2525	-0.5744
19	-2.7533	0.7844	0.3830	-2.3394
31	-0.5571	0.0674	0.5658	-0.6313
32	1.3343	0.1022	0.2561	0.7895

lpt

→ No se ve

son 5

Al observar la gráfica de observaciones vs valores h_{ii} , donde la línea punteada roja representa el valor $h_{ii} = 2 \frac{p}{n} = 0.22$, se puede apreciar que existen ~~5~~ datos del conjunto que son puntos de balanceo según el criterio bajo el cual $h_{ii} > 2 \frac{p}{n}$, los cuales son los presentados en la tabla. Bajo el criterio de puntos de balanceo y generan un aumento en la variabilidad en el modelo RLM.

4.2.3. Puntos influyentes

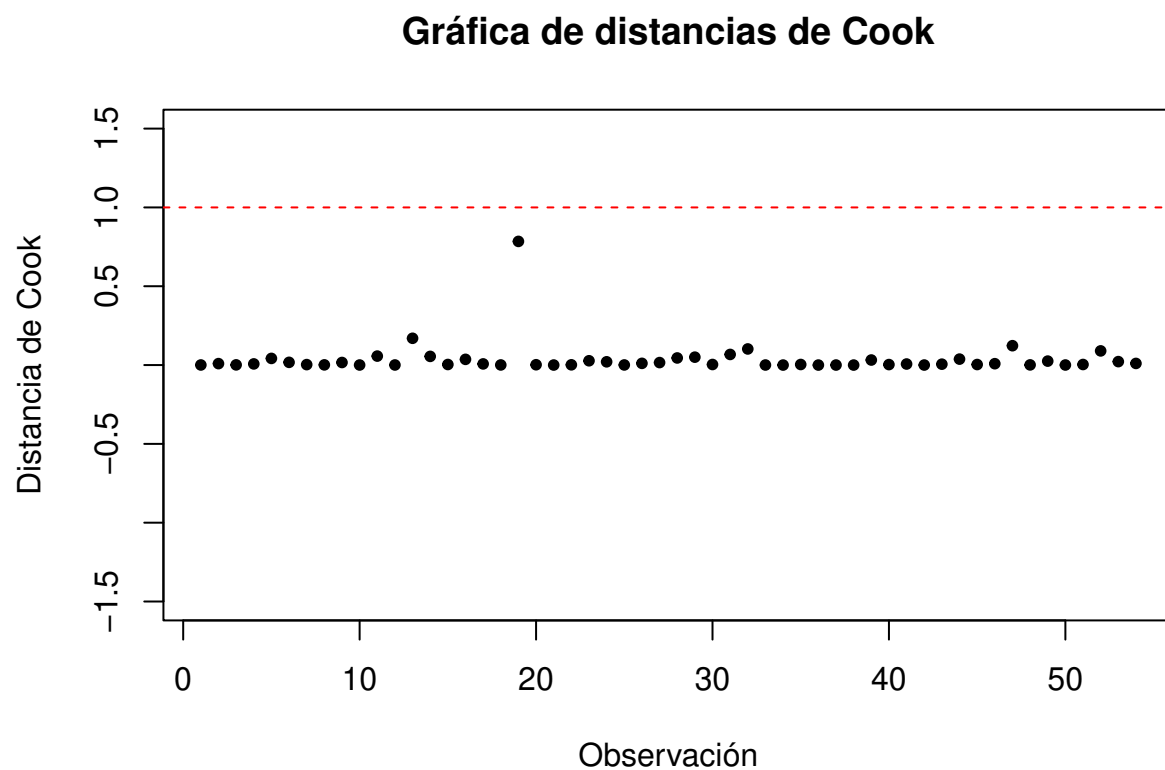


Figura 5: Criterio distancias de Cook para puntos influyentes

Como se puede visualizar en la gráfica de distancias de Cook, donde la línea punteada roja representa el valor 1, se puede inferir que no hay medidas de influencia a partir del criterio $D_i > 1$

Gráfica de observaciones vs Dffits

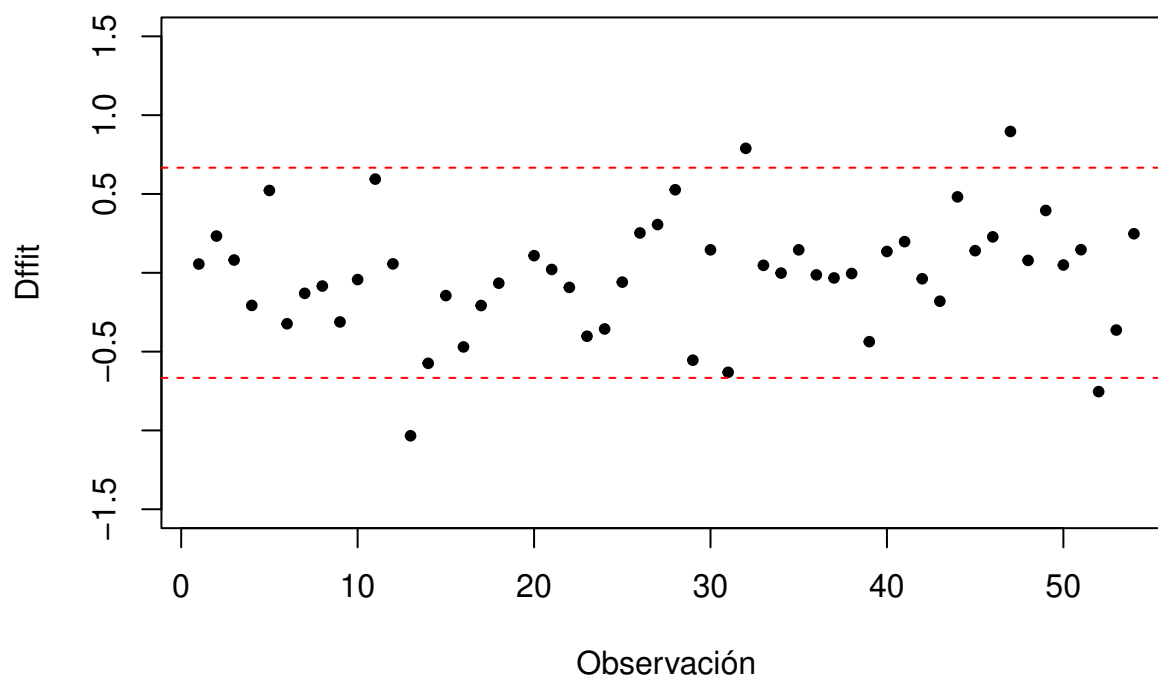


Figura 6: Criterio Dffits para puntos influyentes

Tabla 6: Resumen de diagnostico

	<i>Res.stud</i>	<i>Cook.D_i</i>	<i>h_{ii}value</i>	<i>Dffits</i>
13	-1.7782	0.1699	0.2438	-1.0338
19	-2.7533	0.7844	0.3830	-2.3394
32	1.3343	0.1022	0.2561	0.7895
47	2.2182	0.1228	0.1303	0.8967
52	-1.8514	0.0898	0.1359	-0.7539

2 pt

→ No se ve

(causan...?)

son 4

Como se puede ver, las observaciones 11, 14, 20 son puntos influyentes según el criterio de Dffits, el cual dice que para cualquier punto cuyo $|D_{ffit}| > 2\sqrt{\frac{p}{n}} = 0.42070311619$, es un punto influyente. Cabe destacar también que con el criterio de distancias de Cook, en el cual para cualquier punto cuya $D_i > 1$, es un punto influyente, ninguno de los datos cumple con serlo.

4.3. Conclusión

3 pt

El modelo es válido, ya que se cumplen los supuestos del error, también cumple con los supuestos extremos donde el análisis de este puede decir si pueden afectar en cierta medida a los mismo el modelo y por lo que se puede evidenciar no afecta significativamente en estos.