

4,15

## **Trabajo 1**

Estudiantes

**Alejandra Diego Ramirez Agudelo**  
**Andres Felipe Ramirez Suarez**  
**Juan Diego Ramirez Agudelo**  
**Sara Paulina Aguirre Restrepo**  
Equipo # 37

Docente

**Fransisco Javier Rodriguez Cortes**

Asignatura

**Estadística II**



Sede Medellin  
30 de marzo de 2023

# Índice

<b>1. Pregunta 1</b>	<b>3</b>
1.1. Modelo de regresion . . . . .	3
1.2. Significancia de la regresion . . . . .	3
1.3. Significancia de los parametros . . . . .	3
1.4. interpretación de los parámetros . . . . .	4
1.5. Coeficiente de determinación $R^2$ . . . . .	4
<b>2. Pregunta 2</b>	<b>5</b>
2.1. Planteamiento prueba de hipótesis y modelo reducido . . . . .	5
2.2. Estadístico de prueba y conclusiones . . . . .	5
<b>3. Pregunta 3</b>	<b>5</b>
3.1. Prueba de hipótesis y prueba de hipótesis matricial . . . . .	5
3.2. Estadístico de prueba . . . . .	6
<b>4. Pregunta 4</b>	<b>6</b>
4.1. Supuestos del modelo . . . . .	6
4.1.1. Normalidad de los residuales . . . . .	6
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 influenciales . . . . .	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 de los coeficientes estimados . . . . .	3
2.	Tabla ANOVA significancia de la regresion . . . . .	3
3.	Resumen de los coeficientes . . . . .	4
4.	Resumen de todas las regresiones . . . . .	5
5.	Valor covariables (*) . . . . .	6
6.	TablaPuntos influenciales tomando en cuenta hii.value . . . . .	10
7.	Tabla Puntos influenciales tomando en cuenta Diffits . . . . .	12

## 1. Pregunta 1

Teniendo en cuenta la base de datos asignada. La cual es **equipo37.txt**. Las covariables son: duración de la estadía ( $X_1$ ), rutina de cultivos ( $X_2$ ), número de camas ( $X_3$ ), censo promedio diario ( $X_4$ ) y número de enfermeras ( $X_5$ ).

El modelo propuesto seria, entonces:

$$Y_j = \beta_0 + \beta_i X_{ij} + \varepsilon_j, \quad \varepsilon_j \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2); \quad 1 \leq i \leq 5; \quad 1 \leq j \leq 55$$

$$Y_j = \beta_0 + \beta_1 X_{1j} + \beta_2 X_{2j} + \beta_3 \dots$$

### 1.1. Modelo de regresion

Al ajustar el modelo, se obtienen los siguiente coeficientes:

Tabla 1: Tabla de valores de los coeficientes estimados

	Valor del parametro
$\hat{\beta}_0$	-2.8475
$\hat{\beta}_1$	0.2774
$\hat{\beta}_2$	0.0370
$\hat{\beta}_3$	0.0565
$\hat{\beta}_4$	0.0122
$\hat{\beta}_5$	0.0020

El modelo ajustado es:

$$\hat{Y}_j = -2.8475 + 0.2774X_{1j} + 0.037X_{2j} + 0.0565X_{3j} + 0.0122X_{4j} + 0.002X_{5j}$$

### 1.2. Significancia de la regresion

Para la interpretación se utilizará la tabla de análisis de varianza.

La cual da como resultado lo siguiente:

Tabla 2: Tabla ANOVA significancia de la regresion

	Suma de cuadrados	Grados de libertad	Cuadrado medio	$F_0$	Valor-P
Modelo de regresion	84.8460	5	16.96920	16.518	1.55928e-09
Error	50.3384	49	1.02731		

Tomando a  $\alpha = 0.05$ .

Se tiene que  $1.55928 \times 10^{-9} < 0.05$ , por lo cual se rechaza la  $H_0$  (el modelo no es significativo), concluyendo que el modelo de RLM propuesto es significativo. Eso significa que el riesgo de infección hospitalarias dependen de forma significativa de por lo menos una de las variables predictoras.

### 1.3. Significancia de los parametros

En el siguiente cuadro se presenta información de los parámetros. Lo cual permitirá calcular la significancia de las variables predictoras.

los parámetros

Tabla 3: Resumen de los coeficientes

	Estimacion $\beta_j$	$se(\hat{\beta}_j)$	$T_{0j}$	Valor-P
$\beta_0$	-2.8475	1.7659	-1.6125	0.1133
$\beta_1$	0.2774	0.1225	2.2637	0.0281
$\beta_2$	0.0370	0.0333	1.1089	0.2729
$\beta_3$	0.0565	0.0166	3.4056	0.0013
$\beta_4$	0.0122	0.0086	1.4280	0.1596
$\beta_5$	0.0020	0.0007	2.9134	0.0054

Los valores P permiten concluir con una significancia  $\alpha = 0.05$  que  $\hat{\beta}_1$ ,  $\hat{\beta}_3$  y  $\hat{\beta}_5$  son los únicos ~~valores~~ <sup>Parámetros</sup> significativos. Por otro lado, en el caso  $\hat{\beta}_0$  el valor 0 no se encuentra incluido en el ajuste. Por ende, no es interpretable y como su valor es mayor al alfa, tampoco es significativo.

#### 1.4. interpretación de los parámetros

- $\hat{\beta}_1$  su valor  $P = 0.0281 < 0.05$ , por lo que es un parámetro significativo en presencia de los demás en el modelo. Se tiene que su valor estimado es  $\hat{\beta}_1 = 0.2774$ , además  $X_1$  representa duración de la estadía de cada paciente en el hospital. Se podría interpretar que el aumento en 1 unidad de este número representa un aumento de 0.2774 en el riesgo de infección hospitalarias, mientras las otras variables permanecen constantes. ✓ *la probabilidad promedio.*
- $\hat{\beta}_3$  su valor  $P = 0.0013 < 0.05$ , por lo que es un parámetro significativo en presencia de los demás en el modelo. Se tiene que su valor estimado es  $\hat{\beta}_3 = 0.0565$ , además  $X_3$  representa el numero promedio de camas en el hospital durante el periodo de estudio. Se podría interpretar que el aumento en 1 unidad de este número representa un aumento de 0.0565 en el riesgo de infección hospitalarias, mientras las otras variables permanecen constantes. ✓
- $\hat{\beta}_5$  su valor  $P = 0.0054 < 0.05$ , por lo que es un parámetro significativo en presencia de los demás en el modelo. Se tiene que su valor estimado es  $\hat{\beta}_5 = 0.0020$ , además  $X_5$  representa el numero promedio de enfermeras equivalentes a tiempo completo. Se podría interpretar que el aumento en 1 unidad de este número representa un aumento de 0.0020 en el riesgo de infección hospitalarias, mientras las otras variables permanecen constantes. ✓

#### 1.5. Coeficiente de determinación $R^2$

Basado en la tabla **ANOVA** se tiene que:

$$R^2 = \frac{84.8460}{84.8460 + 50.3384} = 0.6276316$$

El  $R^2$  mide la proporción de la variabilidad total observada en la respuesta que es explicada por el modelo propuesto. El modelo tiene un  $R^2 = 0.6276316$  lo cual significa que las variables independientes explican aproximadamente el 62.76 % de la variabilidad de  $Y$  (Riesgos de infección hospitalarias). ✓

## 2. Pregunta 2

3pt

### 2.1. Planteamiento prueba de hipótesis y modelo reducido

pedían las 3 covariables.

Los parámetros cuyos valores P fueron los más altos corresponden a  $\beta_2, \beta_4, \beta_0$ . Por lo tanto, se plantea la siguiente prueba de hipótesis:

$$\begin{cases} H_0 : \beta_2 = \beta_4 = \beta_0 = 0 \\ H_a : \text{Algun } \beta_i \text{ distinto a } 0 \text{ con } i = 2, 4, 0 \end{cases}$$

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

$$\text{MR: } Y_j = \beta_1 X_{1j} + \beta_3 X_{3j} + \beta_5 X_{5j} + \varepsilon_j, \quad \varepsilon_j \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$$

$i = 1, 2, \dots$

Tabla 4: Resumen de todas las regresiones

	SSE	Covariables en el modelo
Modelo completo	50.338	X1 X2 X3 X4 X5
Modelo reducido	54.036	X1 X3 X5

esto fue ajustado con  $\beta_0$

### 2.2. Estadístico de prueba y conclusiones

Se construye el estadístico de prueba

$$F_0 = \frac{(SSE(\beta_1, \beta_3, \beta_5) - SSE(\beta_0, \beta_1, \dots, \beta_5))/3}{MSE(\beta_0, \beta_1, \dots, \beta_5)} \stackrel{H_0}{\sim} f_{3,49}$$

1pt

$$F_0 = \frac{(35.623 - 34.832)/3}{34.832/49} = 1.199902$$

Al menos son congruentes con el error

Ahora, comparando con un nivel de significancia  $\alpha = 0.05$ ,  $F_0$  con  $f_{0.05, 3, 49} = 2.7939489$  y valor  $P = 0.3196538$ .

Como  $F_0$  es MENOR que 2.7939489. Entonces, no es posible rechazar  $H_0$ , por lo que el subconjunto de las variables no es significativo. Teniendo en cuenta lo anterior, es posible descartar las variables. Debido a que son insignificantes para el modelo. Es decir, su valores son iguales a 0.

2pt

estadísticamente iguales a 0

## 3. Pregunta 3

5pt

### 3.1. Prueba de hipótesis y prueba de hipótesis matricial

Se plantea la siguiente prueba de hipótesis:

$$\begin{cases} H_0 : \beta_1 = \beta_3 - \beta_5, \beta_4 = \beta_2 + \beta_5 \\ H_a : \text{Alguna de las desigualdades no se cumple} \end{cases}$$

Reescribiendo matricialmente:

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

Donde  $\mathbf{L}$  está dada por:

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

Con  $r = 2$  filas linealmente independientes

Donde el modelo reducido está dado por:

$$\begin{aligned} \text{MR: } Y_j &= \beta_0 + (\beta_3 - \beta_5)X_{1j} + \beta_2 X_{2j} + \beta_3 X_{3j} + (\beta_2 + \beta_5)X_{4j} + \beta_5 X_{5j} + \varepsilon_j \\ &= \beta_0 + \beta_2(X_{2j} + X_{4j}) + \beta_3(X_{1j} + X_{3j}) + \beta_5(X_{4j} + X_{5j} - X_{1j}) + \varepsilon_j \\ &= \beta_0 + \beta_2 X_{2j}^* + \beta_3 X_{3j}^* + \beta_5 X_{5j}^* + \varepsilon_j \end{aligned}$$

Suponiendo que:  $\varepsilon_j \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$

En donde los valores de las covariables (\*) son:

Tabla 5: Valor covariables (\*)

	$X_{2j}^*$	$X_{3j}^*$	$X_{5j}^*$
Valor	$X_{2j} + X_{4j}$	$X_{1j} + X_{3j}$	$X_{4j} + X_{5j} - X_{1j}$

### 3.2. Estadístico de prueba

El estadístico de prueba  $F_0$  está dado por:

$$F_0 = \frac{(SSE(MR) - SSE(MF))/2}{MSE(MF)} \stackrel{H_0}{\sim} f_{2,49}$$

$$F_0 = \frac{(SSE(MR) - 50.338/2)}{50.338/49} \stackrel{H_0}{\sim} f_{2,49}$$

## 4. Pregunta 4

### 4.1. Supuestos del modelo

#### 4.1.1. Normalidad de los residuales

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

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

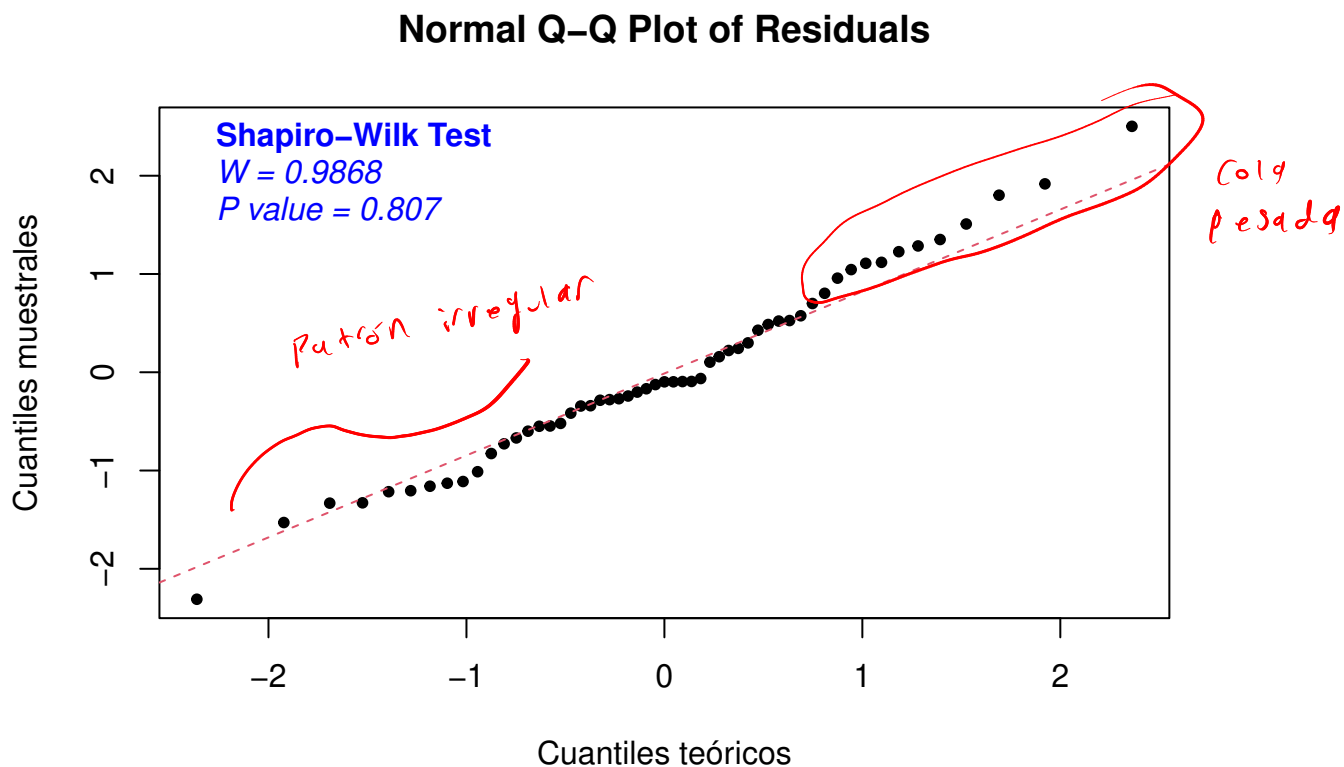


Figura 1: Gráfico cuantil-cuantil y normalidad de residuales

Al ser el P-valor aproximadamente igual a 0.807 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  $\mu$  y varianza  $\sigma^2$ , sin embargo la gráfica de comparación de cuantiles permite ver patrones irregulares, al tener más poder el análisis gráfico, se termina por rechazar el cumplimiento de este supuesto. ✓

*Cola pesada*

*No están probando media constante  $\mu$  ni var de  $\sigma^2$*



## 4.1.2. Varianza constante

3p 1

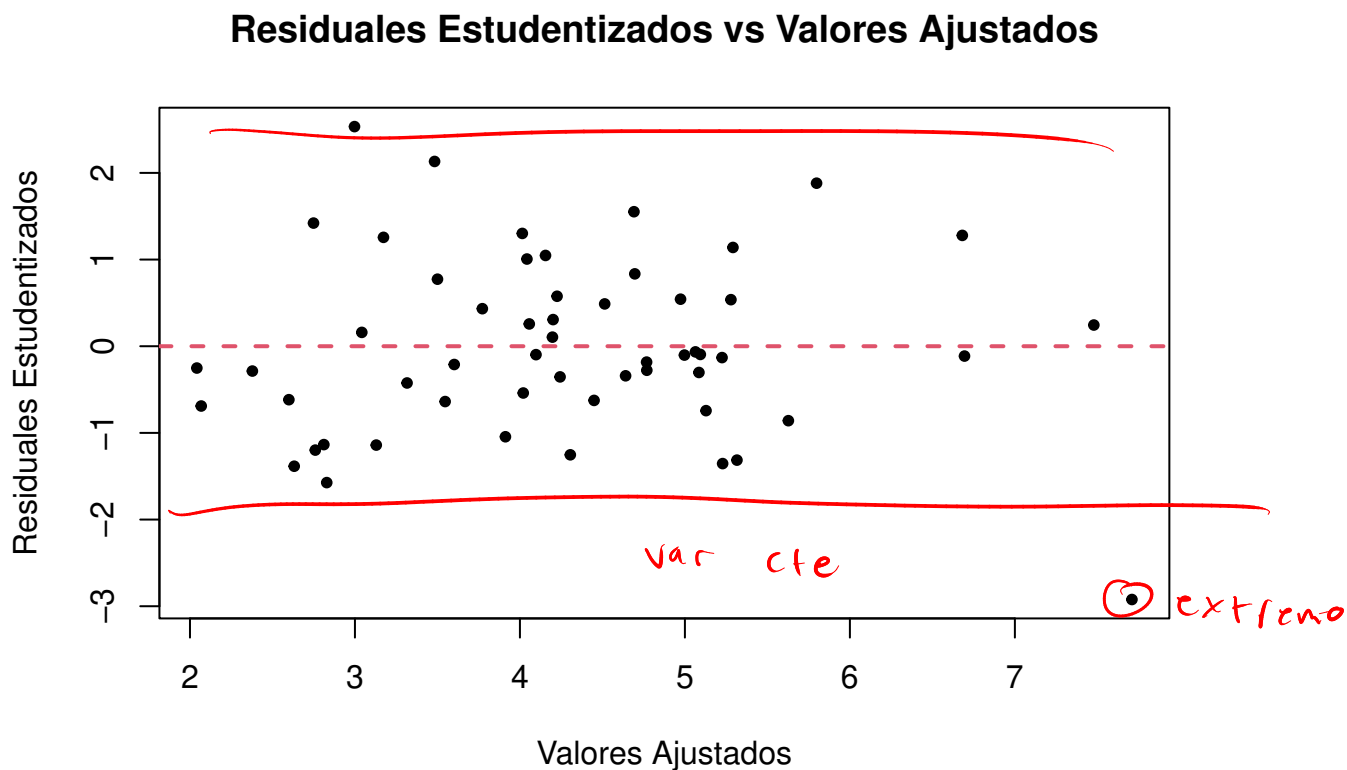


Figura 2: Gráfico residuales estudentizados vs valores ajustados

El modelo cumple con media 0 y además podemos ver en el gráfico de residuales estudentizados vs valores ajustados se puede observar que no hay patrones en los que la varianza muestre un comportamiento que permita descartar una varianza constante, al no haber evidencia suficiente en contra de este supuesto se acepta como cierto

no cualquier  
patrón es adjudicado a eso

El análisis pudo haber sido hecho mucho mejor.

## 4.2. Verificación de las observaciones

### 4.2.1. Datos atípicos

3 p +

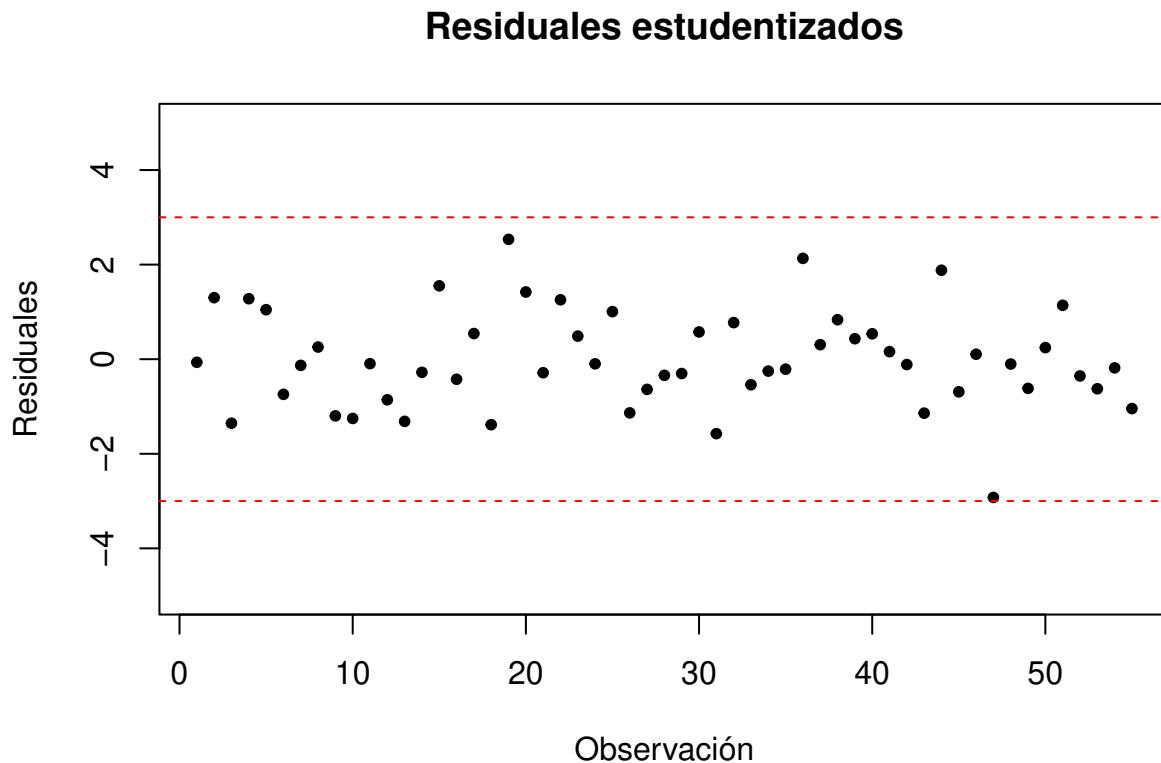


Figura 3: Identificación de datos atípicos

Como se puede observar en la gráfica anterior, no hay datos atípicos en el conjunto de datos pues ningún residual estudentizado sobrepasa el criterio de  $|r_{estud}| > 3$ , solo hay dato que se acerca mucho a este criterio, el cual sería el dato Numero 47 teniendo un residual estudentizado de -2.9225. ✓

## 4.2.2. Puntos de balanceo

2 pt

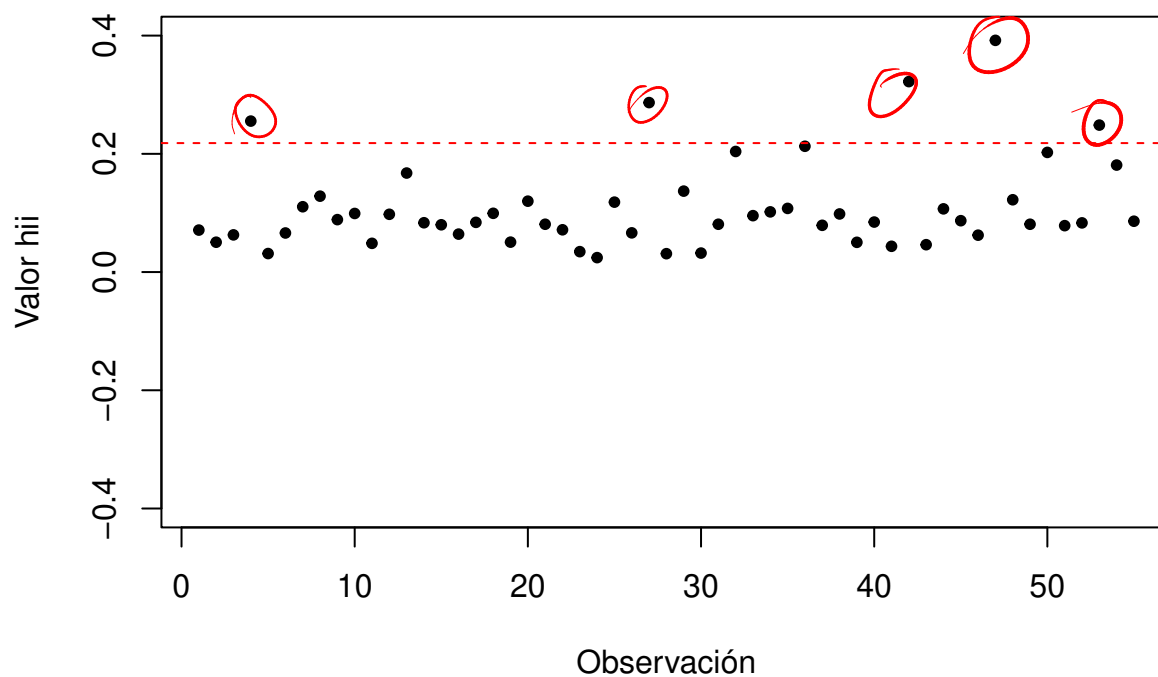
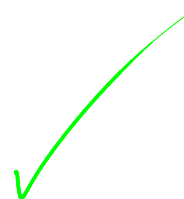
Gráfica de  $h_{ii}$  para las observaciones

Figura 4: Identificación de puntos de balanceo

Tabla 6: TablaPuntos ~~influenciados~~ tomando en cuenta  $h_{ii}.value$ 

	$h_{ii}.value$
punto 4	0.2554
punto 27	0.2867
punto 42	0.3222
punto 47	0.3920
punto 53	0.2486



Al observar la gráfica de observaciones vs valores  $h_{ii}$ , donde la línea punteada roja representa el valor  $h_{ii} = 2\frac{p}{n} = 0.2181$ , se puede apreciar que existen 5 datos del conjunto que son puntos de balanceo según el criterio bajo el cual  $h_{ii} > 2\frac{p}{n} = 0.2181$ , los cuales son los presentados en la tabla.

¿Qué causas?

## 4.2.3. Puntos influyentes

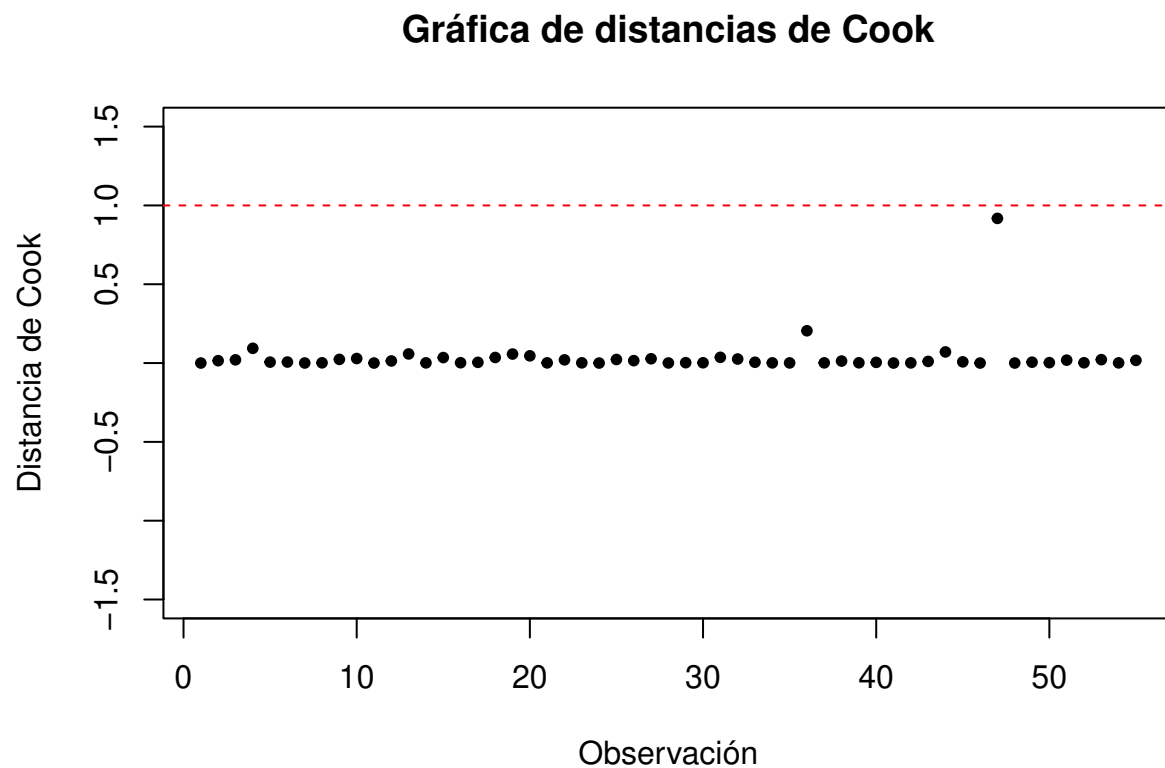
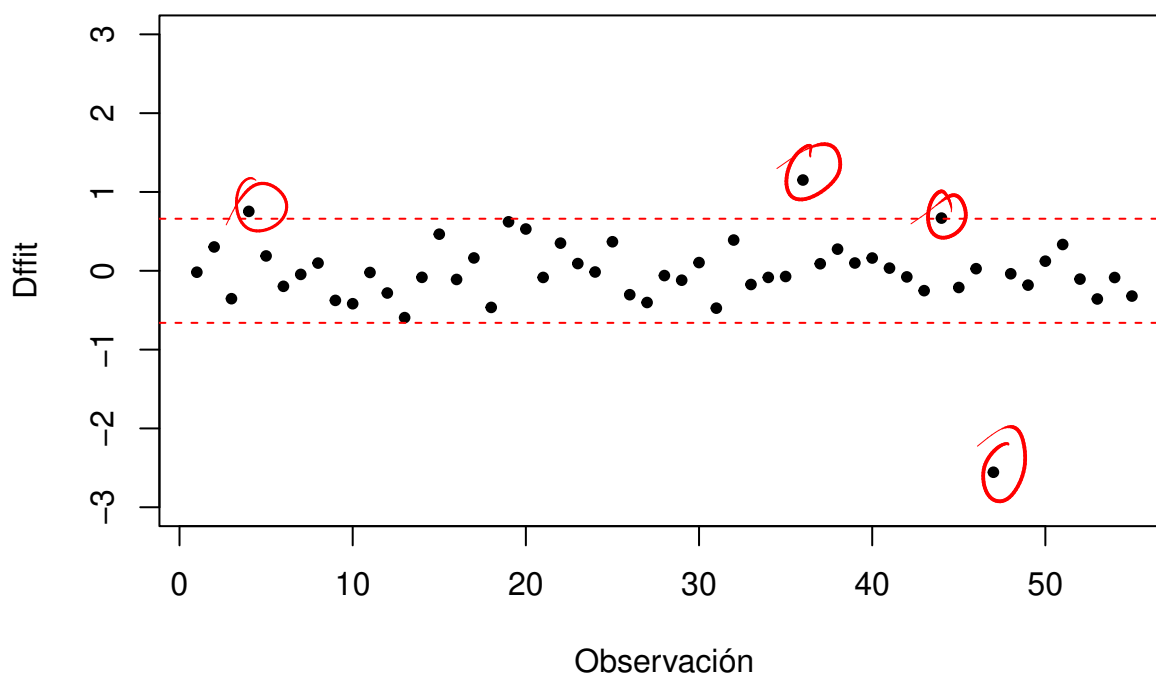


Figura 5: Criterio distancias de Cook para puntos influyentes

Como se observa en la gráfica ~~no~~ tenemos ningún punto influyente tomando en cuenta el criterio de Cook, con lo cual se podría decir que los datos no están siendo afectados de manera significativa por ningún punto ~~particular~~ en la gráfica. influyente

2 pt

### Gráfica de observaciones vs Dffits



1,5 pt

Figura 6: Criterio Dffits para puntos influenciales

Tabla 7: Tabla Puntos influenciales tomando en cuenta Dffits

	Dffits
punto 4	0.7540
punto 36	1.1518
punto 44	0.6686
punto 47	-2.5559

¿cuanto dice?

Como se puede ver, hay 4 puntos influenciales según el criterio de Dffits, los cuales son 4, 36, 44 y 47, respectivamente, este criterio dice que para cualquier punto cuyo  $|D_{ffit}| > 2\sqrt{\frac{p}{n}}$ , 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, y como se dijo anteriormente, ninguno de los datos cumple con serlo.

¿qué causan?

#### 4.3. Conclusión

1,5 pt

1. El modelo no sigue una distribución normal, ya que se observan patrones irregulares en la gráfica de comparación de cuantiles.
2. El modelo cumple con la suposición de que la media es 0 y que la varianza es constante, ya que en la gráfica de residuales estudentizados vs valores ajustados no se observan patrones que sugieran lo contrario.

3. No hay valores atípicos en el conjunto de datos, ya que ningún residual estudentizado sobrepasa el criterio de  $|\text{restud}| > 3$ .
4. Basándonos en la tabla de análisis de varianza se rechaza la  $H_0$  (el modelo no es significativo), concluyendo que el modelo de RLM propuesto es significativo.

En general, estos resultados sugieren que el modelo es adecuado para analizar los datos, aunque es importante tener en cuenta que la falta de normalidad en la distribución de los datos puede afectar la precisión y la validez de los resultados del análisis estadístico. *X*

*falso, su modelo no es válido por no cumplir todos los supuestos*