

3,65

Trabajo 1

Estudiantes

Doris Alejandra Ramos Figueroa
Jhon Sebastian Chidiak Olaya
Laura Valentina Rincon Guataquira

Docente

Julieth Veronica Guarin Escudero

Asignatura

Estadística II



UNIVERSIDAD
NACIONAL
DE COLOMBIA

Sede Medellín
30 de marzo de 2023

Índice

1. pregunta 1	3
1.1. Modelo de regresion	3
1.2. Significancia de la regresión	4
1.3. Significancia de los parámetros	4
1.4. Interpretación de los parámetros	5
2. Pregunta 2	5
2.1. Planteamiento prueba de hipotesis y modelo reducido	5
2.2. Estadístico de prueba y conclusiones	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 de los coeficientes estimados	3
2.	Tabla ANOVA significancia de la regresión	4
3.	Resumen de los coeficientes	4
4.	Resumen de todas las regresiones	6
5.	Tabla de observaciones de Balanceo	11
6.	Punto inflencial	13

1. pregunta 1 18 pt

Teniendo en cuenta la base de datos equipo 13, nuestro modelo de RLM que explica el porcentaje de la eficacia en el control de infecciones hospitalarias (Y: Riesgo de infección) en términos de las variables:

- Duración de la estadia (X_1).
- Rutina de cultivos (X_2).
- Número de camas (X_3).
- Censo promedio diario (X_4).
- número de enfermeras (X_5).



El modelo se propone a continuación:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + \beta_4 X_{4i} + \beta_5 X_{5i} + \varepsilon_i; \quad \varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2); \quad 1 \leq i \leq 55$$

1.1. Modelo de regresión 3 pt

Al ajustar el modelo el modelo que explica el porcentaje de la eficacia en el control de infecciones hospitalarias, se obtienen los siguientes coeficientes:

Cuadro 1: Tabla de valores de los coeficientes estimados

	V. Parámetro
β_0	-0.0058
β_1	0.1675
β_2	-0.0009
β_3	0.0619
β_4	0.0154
β_5	0.0017



Por lo tanto, el modelo de regresión ajustado es:

$$\hat{Y}_i = -0.0058 + 0.1675X_{1i} - 9 \times 10^{-4}X_{2i} + 0.0619X_{3i} + 0.0154X_{4i} + 0.0017X_{5i}$$



Donde $1 \leq i \leq 55$



1.2. Significancia de la regresión 4 pt

Se plantea el siguiente Juego de Hipotesis, para analizar la significancia de la regresion:

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

Cuyo estadístico de prueba es:

$$F_0 = \frac{\cancel{MST}}{MSE} \stackrel{H_0}{\sim} \cancel{f_{5, \text{error}(55)-6}} \quad \checkmark \quad F_0 = \frac{MSR}{MSE} \neq \frac{MST}{MSE} \quad \sim f_{5, 99} \quad (1)$$

Para la significancia de la regresión tomaremos los resultados de la tabla Anova presentada a continuacion:

Cuadro 2: Tabla ANOVA significancia de la regresión

	Suma cuadratica	Grados de libertad	Cuadrado medio	F_0	Valor-P
Modelo LRM	69.4276	5	13.885513	19.6988	1.03971e-10
Error	34.5397	49	0.704892		

De acuerdo con la tabla anova, con un nivel de significancia del $\alpha = 0.05$, por el criterio del valor-p encontrado en la anterior tabla anova (tabla 2). Podemos determinar que nuestro modelo rechaza H_0 debido a que $V_p < \alpha = 0.05$ yes significativo ya que existe algun parametro distinto de cero. . Al rechazarla probamos que existe una relación de regresión, sin embargo, esto no garantiza que el modelo resulte útil para hacer predicciones. ✓

1.3. Significancia de los parámetros 6 pt

En el siguiente cuadro se presenta información de los parámetros, la cual permitirá determinar cuáles de ellos son significativos.

Cuadro 3: Resumen de los coeficientes

	$\hat{\beta}_j$	$SE(\hat{\beta}_j)$	T_{0j}	P-valor
β_0	-0.0058	1.6503	-0.0035	0.9972
β_1	0.1675	0.0735	2.2799	0.0270
β_2	-0.0009	0.0309	-0.0304	0.9759
β_3	0.0619	0.0141	4.3929	0.0001
β_4	0.0154	0.0071	2.1749	0.0345
β_5	0.0017	0.0006	2.6766	0.0101

Observando los resultados obtenidos, podemos determinar que con un nivel de significancia $\alpha = 0.05$ los parámetros $\beta_1, \beta_3, \beta_4$ y β_5 son significativos porque sus valores p son menores a $\alpha = 0.05$. Y los parámetros β_0 y β_2 no son significativos debido a que su p es mayor a $\alpha = 0.05$. ✓

1.4. Interpretación de los parámetros 2,5 pt

¿grasa corporal! ojo con el plagio.

A continuación se hará la interpretación de los parámetros que son significativos, ya que los otros parámetros pueden ser eliminados si no aportan el modelo de la grasa corporal.

- $\hat{\beta}_1 = 0.1675$: Indica que por cada unidad que aumenta la duración de la estancia, la probabilidad promedio estimada de adquirir una infección en el hospital aumenta en un 0.1675 unidades cuando las demás variables predictoras se mantienen fijas. ✓
- $\hat{\beta}_3 = 0.0619$: Indica que por cada unidad que aumenta el número de camas, la probabilidad promedio estimada de adquirir una infección en el hospital aumenta en un 0.0619 unidades cuando las demás variables predictoras se mantienen fijas. ✓
- $\hat{\beta}_4 = 0.0154$: Indica que por cada unidad que aumenta el censo promedio diario, la probabilidad promedio estimada de adquirir una infección en el hospital aumenta en un 0.0154 unidades cuando las demás variables predictoras se mantienen fijas. ✓
- $\hat{\beta}_5 = 0.0017$: Indica que por cada unidad que aumenta el número de enfermeras, la probabilidad promedio estimada de adquirir una infección en el hospital aumenta en un 0.0017 unidades cuando las demás variables predictoras se mantienen fijas. ✓

Coeficiente de determinación múltiple R^2 2,5 pt

El coeficiente de determinación del modelo es $R^2 = 0.6678$ lo que indica que el modelo explica 66.78 % de la variabilidad total en el riesgo de infección en el hospital y es explicada por el modelo propuesto de regresión lineal múltiple.

¿cómo se calcula!

2. Pregunta 2 2,5 pt

2.1. Planteamiento prueba de hipótesis y modelo reducido

Los tres parámetros cuyos valores P fueron los más altos corresponden a $\beta_1, \beta_2, \beta_4$. Debido a esto, se propone la siguiente prueba de hipótesis:

$$\begin{cases} H_0 : \beta_1 = \beta_2 = \beta_4 = 0 \\ H_1 : \text{Algún } \beta_j \text{ distinto de 0 para } j = 1, 2, 4 \end{cases}$$

✓

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

Cuadro 4: Resumen de todas las regresiones

	Suma cuadratica del error	Predictores del modelo
M.Completo	34.540	X1 X2 X3 X4 X5
M.Reducido	45.963	X3 X5

El modelo completo es el definido en la sección 1.1, y el modelo reducido para la prueba de significancia del subconjunto es:

$$Y_i = \beta_0 + \beta_3 X_{3i} + \beta_5 X_{5i} + \varepsilon_i; \varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2); 1 \leq i \leq 55$$

2.2. Estadístico de prueba y conclusiones

2,5 pt

Se construye el estadístico de prueba como:

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

$$F_0 = \frac{(45.963 - 34.540)/3}{0.9380204081632} \stackrel{H_0}{\sim} f_{3,49}$$

$$F_0 = 4.06 \stackrel{H_0}{\sim} f_{3,49}$$

Comparamos F_0 con $f_{0.95,3,49} = 2.7939$, se evidencia que $F_0 = 4.09 > f_{0.95,3,49}$, por lo cual a partir de la prueba de H_0 realizada se concluye que es mayor, entonces se rechaza la hipótesis nula planteada y también nos indica que el subconjunto es significativo. Por lo tanto, es posible descartar las variables del M.Reducido.

0 pt

¡Falso! Todo lo contrario

3. Pregunta 3

5 pt

3.1. Prueba de hipótesis y prueba de hipótesis matricial

Se hace la pregunta si... por consiguiente se plantea la siguiente prueba de hipótesis:

¿Plagio?
¿copy paste!

$$\begin{cases} H_0 : \beta_1 = \beta_5; \beta_2 = \beta_3 \\ H_1 : \text{Alguna de las igualdades no se cumple} \end{cases}$$

reescribiendo matricialmente:

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

2 pt

Donde Nuestra matriz \mathbf{L} esta dada por:

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

✓

El modelo reducido está dado por:

$$Y_i = \beta_0 + \beta_1 X_{1i}^* + \beta_2 X_{2i}^* + \beta_4 X_{4i} + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2); \quad 1 \leq i \leq 55$$

✓

1 pt

Donde $X_{1i}^* = X_{1i} + X_{5i}$ y $X_{2i}^* = X_{2i} + X_{3i}$

3.2. Estadístico de prueba

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

✓

2 pt

$$F_0 = \frac{(SSE(MR) - 34.540)/2}{0.705} \stackrel{H_0}{\sim} f_{2,49}$$

✓

4. Pregunta 4

11 pt

4.1. Supuestos del modelo

4.1.1. Normalidad de los residuales

3 pt

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

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

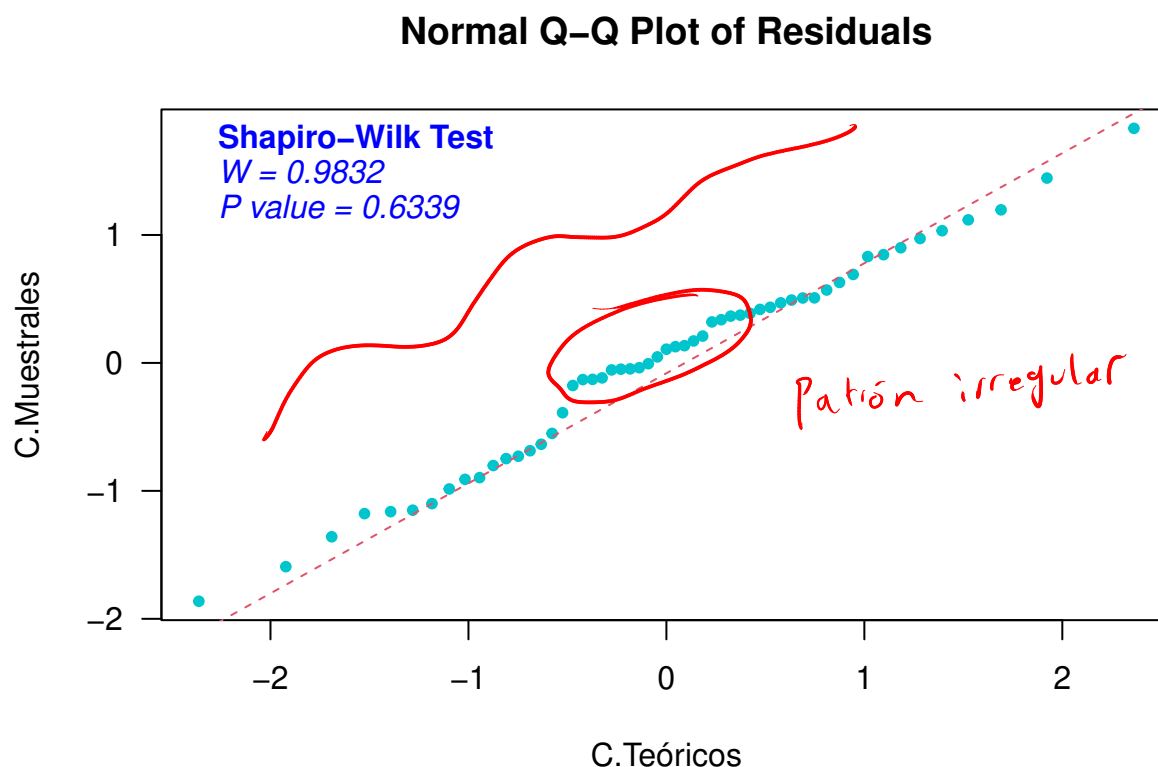


Figura 1: Gráfico cuantil-cuantil y normalidad de residuales

Todo modelo de regresión debe cumplir con sus supuestos, en este caso analizaremos el supuesto de normalidad con la prueba de shapiro wilk y para asegurar; con una gráfica.

En la gráfica de normalidad Q-Q plot of residuals Miramos en la parte superior el test y la gráfica. Según el test de shapiro wilk el valor p es aproximadamente 0.6339, se trabaja con un $\alpha = 0.05$ entonces $0.6339 > 0.05$, se acepta H_0 cero y se supone que debe distribuir normal, sin embargo, cuándo miramos esa gráfica se observa horas más pesadas y patrones anormales lo que nos haría suponer que no se cumple el supuesto de normalidad. Teniendo dos análisis acerca de la normalidad optamos por elegir la gráfica ya que tiene más qué divinidad un análisis gráfico qué un test generalizado.

De hecho el de
 ustedes no tiene colas
 pesadas

4.1.2. Varianza constante 0,5 pt

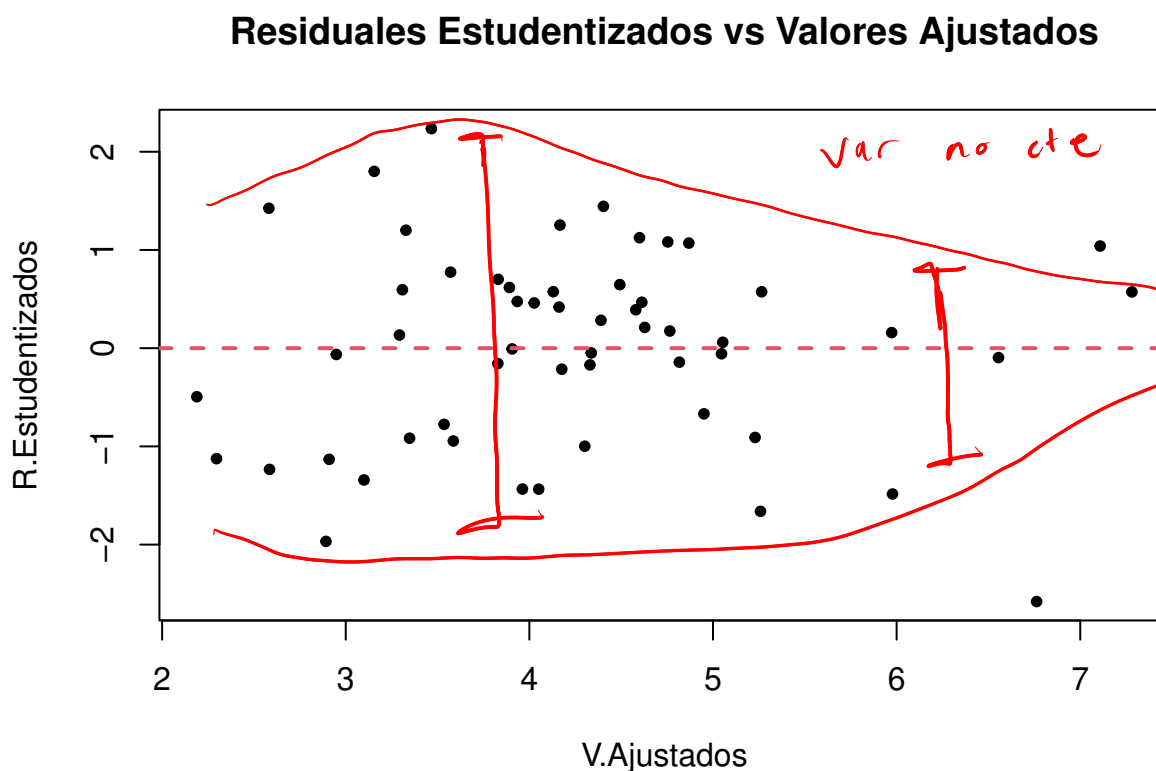


Figura 2: Gráfico residuales estudentizados vs valores ajustados

Después de analizar el supuesto de normalidad, miraremos el otro supuesto acerca de si el modelo de regresión múltiple tiene varianza constante. Analizamos la gráfica, observamos la forma de nube de puntos y no hay un patrón claro, sospecha de una posible heterocedasticidad. Identificando el patrón en la forma de la nube de puntos hay una línea que cruza los residuos; esto definitivamente indica la presencia de heterocedasticidad en el modelo y la ausencia del supuesto de varianza constante supone que la variabilidad de los errores no es la misma en todas las combinaciones de los predictores. Debido a todo lo expuesto, llegamos a la conclusión que no tiene varianza constante. X

¿por qué el no haber patrón nos da esa sospecha?

→ Nada que ver, eso es de la media 0

Concluyen efectivamente var no cte, pero los argumentos son incorrectos

4.2. Verificación de las observaciones

4.2.1. Datos atípicos

3pt

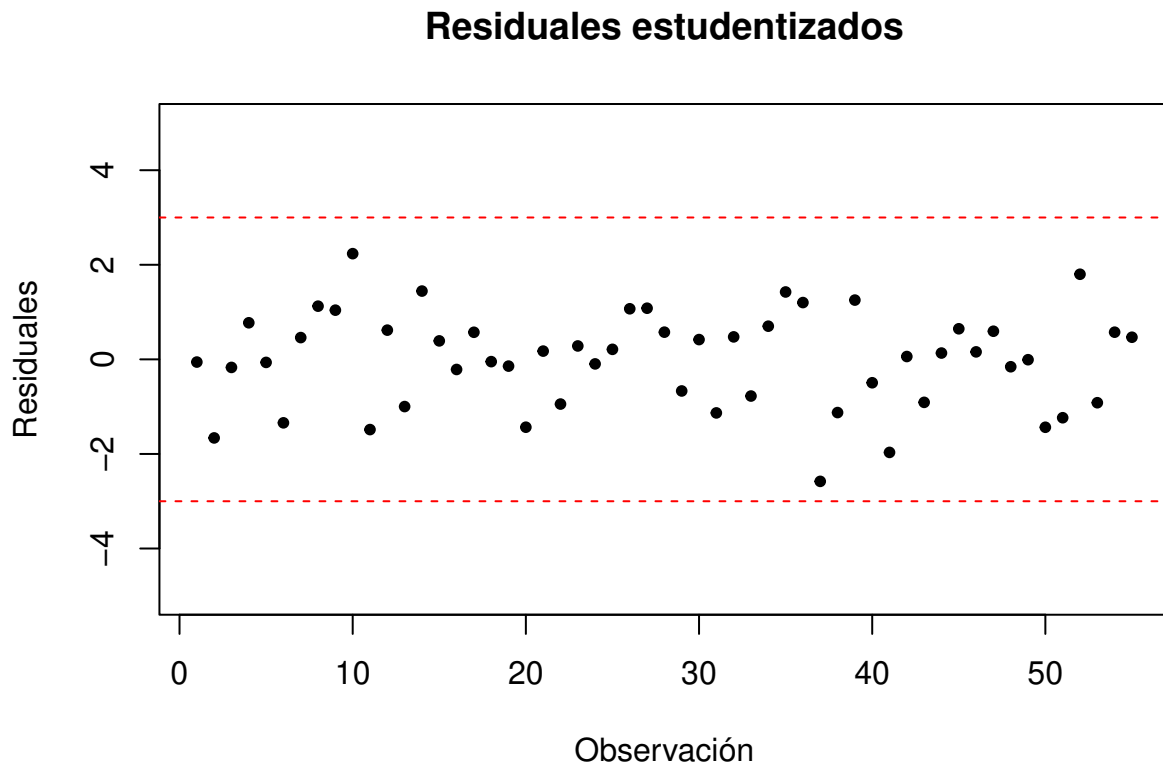


Figura 3: Identificación de datos atípicos

Verificaremos las observaciones con una gráfica que presenta datos atípicos de los residuales estudentizados por las observaciones en el eje x y los residuales en el eje y. El criterio nos dice que $|\text{restud}| > 3$, en la gráfica los puntos de nubes se mantienen entre $(-3, 3)$, a causa de eso no hay datos atípicos en el conjunto de datos. ✓

4.2.2. Puntos de balanceo

lot

Gráfica de hii para las observaciones

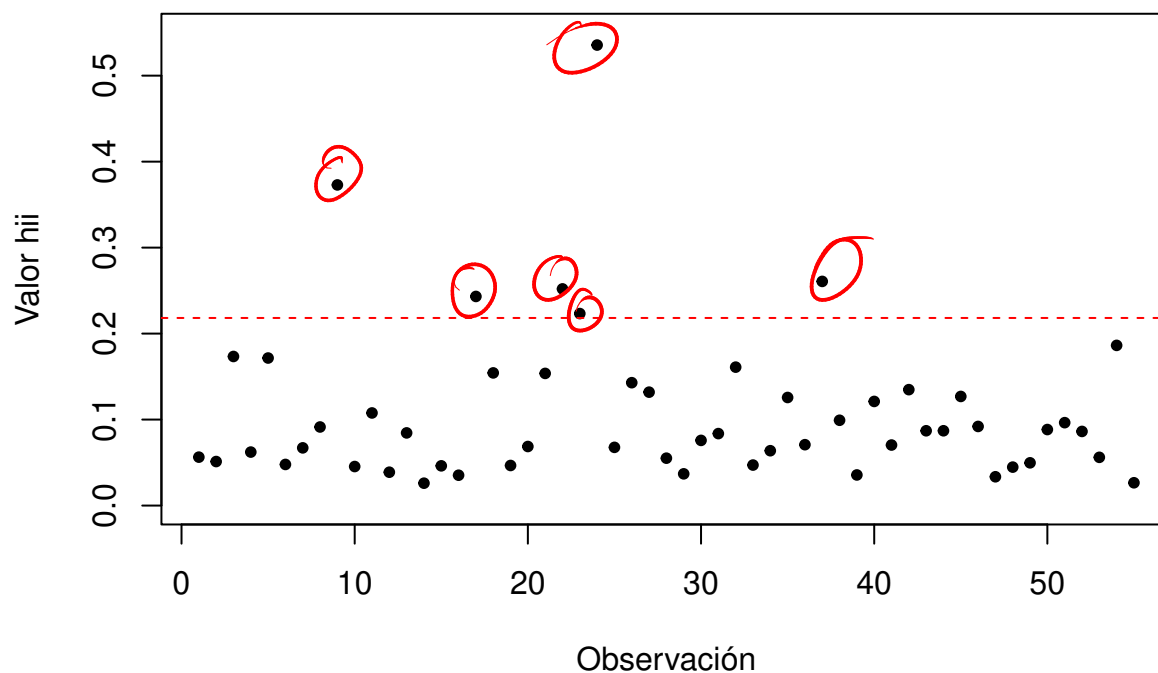
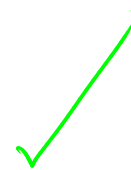


Figura 4: Identificación de puntos de balanceo

Cuadro 5: Tabla de observaciones de Balanceo

	Errores Estudentizados	D.Cook	Valor hii	DFITS
9	1.0402	0.1072	0.3729	0.8027
17	0.5728	0.0176	0.2432	0.3225
22	-0.9450	0.0502	0.2520	-0.5480
23	0.2837	0.0039	0.2233	0.1507
24	-0.0958	0.0018	0.5356	-0.1018
37	-2.5800	0.3911	0.2607	-1.6311



criterio $\{ \frac{2 \cdot 0.87}{h_{ii} > \frac{2 \cdot 0.87}{n}} \}$

punto de balanceo, no atípico

Para la figura 4 examinamos que si hay datos atípicos porque hay puntos de balanceo fuera de las líneas y de $(-3,3)$. estos datos según la tabla de puntos de balanceo son 9,17,22,23,24 y 37. Y es que estos valores extremos desvían significativamente el valor de la muestra, por lo tanto, su inclusión puede sesgar la interpretación de los resultados y en lo posible, sería

mejor eliminarlos aunque es importante tenerlos; siempre y cuando no sean muchos ya que eso representa un error en la medición o recolección de datos, y es necesario que hallan para que proporcionen información valiosa sobre la variabilidad de los datos.

Lean la teoría de qué es un punto de balanceo, lo que causa y cómo identificarlos.

4.2.3. Puntos influyentes

Gráfica de distancias de Cook

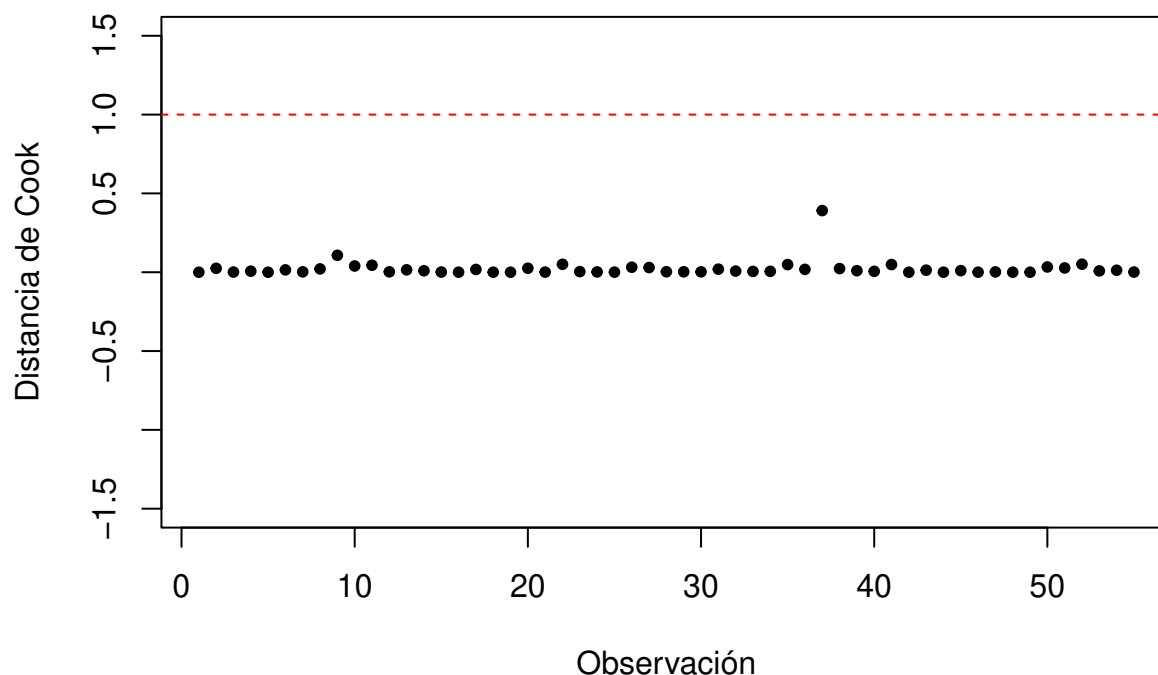


Figura 5: Criterio distancias de Cook para puntos influyentes

Con la figura 5 de distancia de Cook, con observaciones en el eje x y distancia Cook(y). El criterio de Cook dice que observaciones i será influencia le si $D_i > 1$, en la gráfica ningún punto es mayor a 1, por lo tanto no hay puntos influenciados. ✓

2pt

Gráfica de observaciones vs Dffits

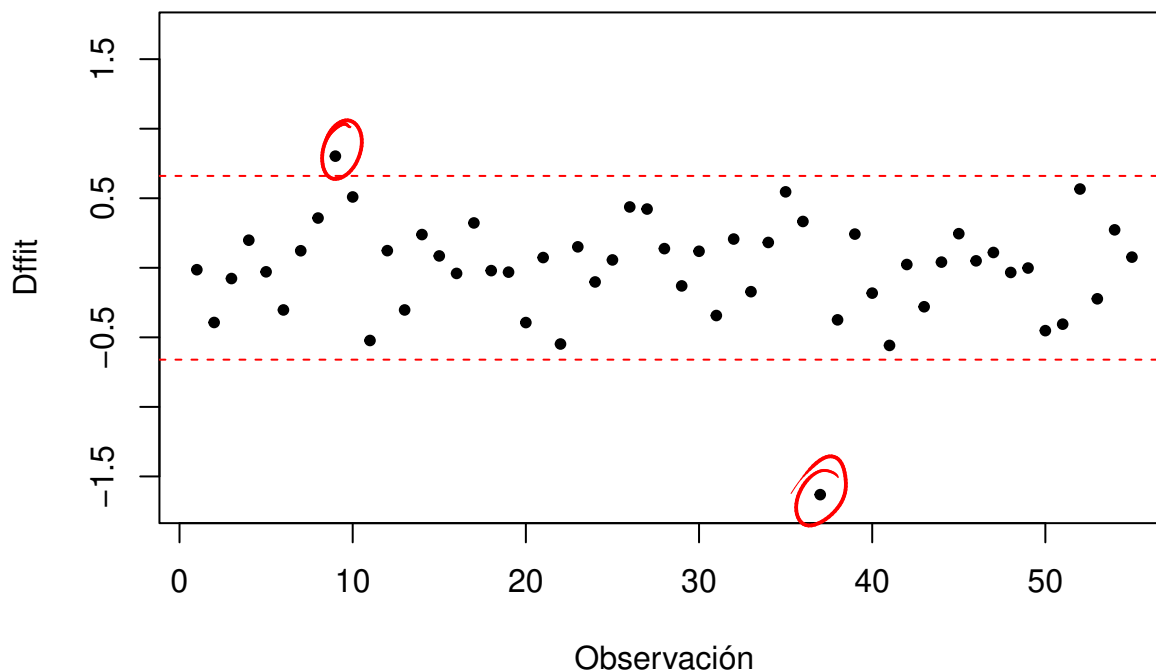


Figura 6: Criterio Dffits para puntos influyentes

0,5 pt

Cuadro 6: Punto influyente

	Errores Estudentizados	D.Cook	Valor hii	DFITS
9	1.0402	0.1072	0.3729	0.8027
37	-2.5800	0.3911	0.2607	-1.6311

En la figura 6 por medio deo criterio Dffits para puntos influyentes, se presenta una gráfica de observaciones vs Dffits, con observación en eje x y Dffits en el eje y. El criterio de Dffits nos expresa que para cualquier punto () por lo que las observaciones 9 y 37 son puntos influyentes.

¿qué es un punto influyente? ¿copy + paste! ¿plagio de nuevo? según este criterio y qué causa?

4.3. Conclusión

- De este modelo

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + \beta_4 X_{4i} + \beta_5 X_{5i} + \varepsilon_i; \quad \varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2); \quad 1 \leq i \leq 55$$

se puede deducir que β_1 β_2 β_3 β_4 β_5 que respetivamente son ((Duracion de la estadia(X_1).Numero de camas(X_2),Censo promedio diario(X_3),numero de enfermeras(X_4)) Nos explican el riesgo de y son indispensables para explicar el riesgo de infectarse en el hospital ya que son significativas. ~

- Si la causa de los valores atípicos se puede identificar y corregir, se puede proceder a realizar un análisis normal de los datos restantes. Sin embargo, si la causa de los valores atípicos no se puede identificar, puede ser necesario considerar el uso de técnicas estadísticas más avanzadas, como la regresión Robusta o el análisis de resúmenes no paramétricos. ✓

En cualquier caso, es importante tener en cuenta que los valores atípicos no siempre deben eliminarse de la muestra, ya que pueden proporcionar información valiosa sobre la variabilidad de los datos. En lugar de eliminarlos, es recomendable tratarlos en el análisis para entender su impacto en los resultados y determinar si son valores extremos legítimos o errores de medición. ✓

No responden nada sobre validez del modelo