

4,3

## Trabajo 1

Estudiantes

Elizabeth Laverde Sierra  
Maria Clara González Arismendi  
Paula Andrea Cifuentes David  
Daniel Felipe Zapata Patron  
Equipo 11

Docente

Julieth Verónica Guarín Escudero

Asignatura

Estadística II



UNIVERSIDAD  
**NACIONAL**  
DE COLOMBIA

Sede Medellín  
30 de marzo de 2023

# Índice

<b>1. Pregunta 1</b>	<b>3</b>
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 . . . . .	5
1.5. Coeficiente de determinación múltiple $R^2$ . . . . .	5
<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 . . . . .	6
<b>3. Pregunta 3</b>	<b>6</b>
3.1. Prueba de hipótesis y prueba de hipótesis matricial . . . . .	6
3.2. Estadístico de prueba . . . . .	7
<b>4. Pregunta 4</b>	<b>7</b>
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 . . . . .	11
<b>5. Conclusiones</b>	<b>13</b>

## Índice de figuras

1.	Gráfico cuantil-cuantil y normalidad de los 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 . . . . .	12

## Índice de cuadros

1.	Tabla de valores de los coeficientes estimados . . . . .	3
2.	Tabla anova de significancia de la regresión . . . . .	4
3.	Resumen de los coeficientes . . . . .	4
4.	Resumen de todas las regresiones . . . . .	6
5.	Tabla de puntos de balanceo . . . . .	11
6.	Tabla de puntos influenciales bajo el criterio de Dffits . . . . .	13

## 1. Pregunta 1 18 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  $R^2$ .

### 1.1. Modelo de regresión 3 p+

Teniendo en cuenta la base de datos asignada, la cual es **Equipo11.txt**, en la que las covariables son:

$Y$ : Riesgo de infección

$X_1$ : Duración de la estadía

$X_2$ : Rutina de cultivos

$X_3$ : Número de camas

$X_4$ : Censo promedio diario

$X_5$ : Número de enfermeras

el modelo de regresión lineal múltiple que se propone es:

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

al ajustar el modelo anterior, se obtienen los siguientes coeficientes:

Cuadro 1: Tabla de valores de los coeficientes estimados

	valor del parámetro
$\hat{\beta}_0$	-2.2385
$\hat{\beta}_1$	0.0579
$\hat{\beta}_2$	0.0484
$\hat{\beta}_3$	0.0317
$\hat{\beta}_4$	0.0284
$\hat{\beta}_5$	0.0021

por lo tanto, el modelo de regresión ajustado es:

$$\hat{Y}_i = -2.2385 + 0.0579X_{1i} + 0.0484X_{2i} + 0.0317X_{3i} + 0.0284X_{4i} + 0.0021X_{5i}$$

### 1.2. Significancia de la regresión 4 p+

Para la significancia de la regresión se hará uso de la siguiente tabla anova:

Cuadro 2: Tabla anova de significancia de la regresión

	Sumas de cuadrados	g.l	Cuadrado medio	$F_0$	Valor-P
Modelo de regresión	52.4755	5	10.495096	11.668	1.87366e-07
Error	44.0743	49	0.899476		

Para analizar la significancia de la regresión se establece la siguiente prueba de hipótesis:

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

cuyo estadístico de prueba es:

$$F_0 = \frac{\cancel{MST}}{MSE} \stackrel{H_0}{\sim} f_{5,49} \quad (1)$$

Con un valor de significancia de  $\alpha = 0,05$  y a partir de la tabla Anova se concluye que el vP  $< 0.05$ , por lo que se rechaza la hipótesis nula por lo cual es posible afirmar que el modelo de regresión lineal múltiple sí es significativo. Esto quiere decir, que el riesgo de infección depende significativamente de al menos una de las predictoras del modelo.

### 1.3. Significancia de los parámetros

En el siguiente cuadro se presenta la información de los parámetros, la cuál permitirá determinar los parámetros individuales que son significativos o no para el modelo.

Cuadro 3: Resumen de los coeficientes

	Estimación $\beta_j$	$se(\hat{\beta}_j)$	$T_{0j}$	Valor-P
$\beta_0$	-2.2385	1.4607	-1.5325	0.1318
$\beta_1$	0.0579	0.1028	0.5638	0.5755
$\beta_2$	0.0484	0.0267	1.8081	0.0767
$\beta_3$	0.0317	0.0152	2.0904	0.0418
$\beta_4$	0.0284	0.0096	2.9460	0.0049
$\beta_5$	0.0021	0.0008	2.4583	0.0175

Estas pruebas establecen la siguiente prueba de hipótesis:

$$\begin{cases} H_0 : \beta_j = 0 \\ H_a : \beta_j \neq 0 \end{cases} \quad j = 1, 2, \dots, 5$$

De acuerdo al valor p, en el cuadro 3 y teniendo una significancia de  $\alpha = 0.05$  se puede concluir que  $\beta_3, \beta_4, \beta_5$  son significativos cada uno en presencia de los demás parámetros.

Asimismo, es posible deducir que  $\beta_0, \beta_1, \beta_2$  son parámetros individualmente no significativos cada uno en presencia de los demás parámetros. ✓

#### 1.4. Interpretación de los parámetros

- $\hat{\beta}_3$  indica que por cada unidad que aumenta <sup>al</sup> número de camas, el promedio del riesgo de infección aumenta en 0.0317 cuando las demás se mantienen fijas. <sup>2 pt</sup>
- $\hat{\beta}_4$  indica que por cada unidad que aumenta de censo promedio diario, el promedio del riesgo de infección aumenta en 0.0284 cuando las demás se mantienen fijas. <sup>el porcentaje promedio " probabilidad"</sup>
- $\hat{\beta}_5$  indica que por cada unidad que aumenta de número de enfermeras, el promedio del riesgo de infección aumenta en 0.0021 cuando las demás se mantienen fijas.

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

3 pt

A partir de la tabla ANOVA calculamos el

$$R^2 = \frac{SSR}{SST} = \frac{52.4755}{52.4755 + 44.0743} = 0.5435 \quad \checkmark$$

lo que quiere decir que aproximadamente el 54.35 % de la variabilidad total en el riesgo de infección, es explicado por las variables del modelo de regresión lineal múltiple propuesto. ✓

## 2. Pregunta 2

4 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

2 pt

Los parámetros cuyos valores P fueron los más altos corresponden a  $\beta_1, \beta_2, \beta_3$ . Por tanto, se plantea la siguiente prueba de hipótesis:

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

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

$$\text{MR: } Y_i = \beta_4 X_{4i} + \beta_5 X_{5i} + \varepsilon_i; \varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2); 1 \leq i \leq 55 \quad \times \text{ y } \beta_0 ?$$

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

	$SSE$	Covariables en el modelo				
Modelo completo	44.074	X1	X2	X3	X4	X5
Modelo reducido	50.501			X4	X5	

## 2.2. Estadístico de prueba y conclusiones

Se construye el estadístico de prueba como:

$$F_0 = \frac{(SSE(\beta_3, \beta_4, \beta_5) - SSE(\beta_0, \beta_1, \dots, \beta_5))/3}{MSE(\beta_0, \beta_1, \dots, \beta_5)} \stackrel{H_0}{\sim} f_{3,49}$$

$$= \frac{(50.501 - 44.074)/3}{44.074/49} = 2.38$$

entonces  $f_{0.95,3,49} = 2.7939$ , Ahora bien, teniendo el valor crítico podemos compararlo con el estadístico de prueba:  $F_0 = 2.38 < f_{0.95,3,49}$ . Finalmente, se concluye que no se rechaza la hipótesis nula  $H_0$ , por lo que las variables no son significativas en presencia de los demás parámetros, es decir, se pueden retirar del modelo.

## 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 sumas de cuadrados extra). 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:

$$\begin{cases} H_0 : \beta_2 = \beta_5 \text{ y } \beta_1 = \beta_3 \\ H_a : \beta_2 \neq \beta_5 \text{ y } \beta_1 \neq \beta_3 \end{cases}$$

Reescribiendo matricialmente:

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

Donde  $L$  está dada por:

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

Por lo que  $H_0$  se puede plantear de la siguiente manera:

$$\mathbf{H}_0 = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 & -1 \\ 0 & 1 & 0 & -1 & 0 & 0 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \beta_5 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \checkmark$$

En donde el modelo reducido está dado por:

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

Donde

$$X_{1i}^* = X_{1i} + X_{3i} \quad \text{y} \quad X_{2i}^* = X_{2i} + X_{5i} \quad \checkmark$$

### 3.2. Estadístico de prueba

El estadístico de prueba  $F_0$  está dado por:

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

$$F_0 = \frac{(SSE(MR) - 44.074)/2}{0.8995} \stackrel{H_0}{\sim} f_{2,49} \quad \checkmark$$

2 pt

## 4. Pregunta 4 16,5

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.

### 4.1. Supuestos del modelo

#### 4.1.1. Normalidad de los residuales 4 pt

Para la validación de este supuesto, se plantea la siguiente prueba de hipótesis:

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

Muy bien que se dieron cuenta que Shapiro-wilk no es la prueba de hipótesis sino un test.



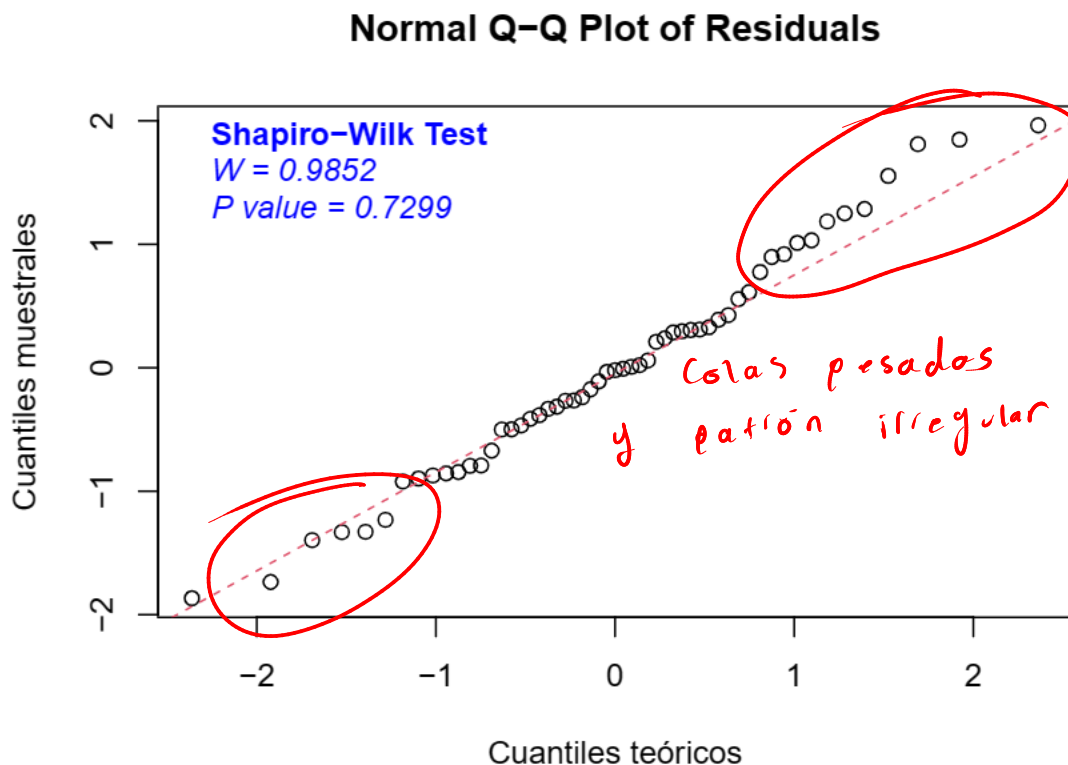


Figura 1: Gráfico cuantil-cuantil y normalidad de los residuales

Mediante el test de Shapiro-Wilk podemos concluir que el valor- $p=0.7299 > \alpha(0.05)$  por lo que la hipótesis nula no se rechaza, lo que quiere decir que los datos distribuyen de manera normal con media  $\mu$  y varianza  $\sigma^2$  sin embargo, es posible notar del análisis de la gráfica anterior, que existen datos dispersos los cuales no siguen la línea roja que representa el ajuste a la distribución de los residuales a una distribución normal. Así, no distribuye de manera normal.

*-> analicen más a fondo que hay un patrón y colas pesadas, pero concluyen bien*

## 4.1.2. Varianza constante

lpt

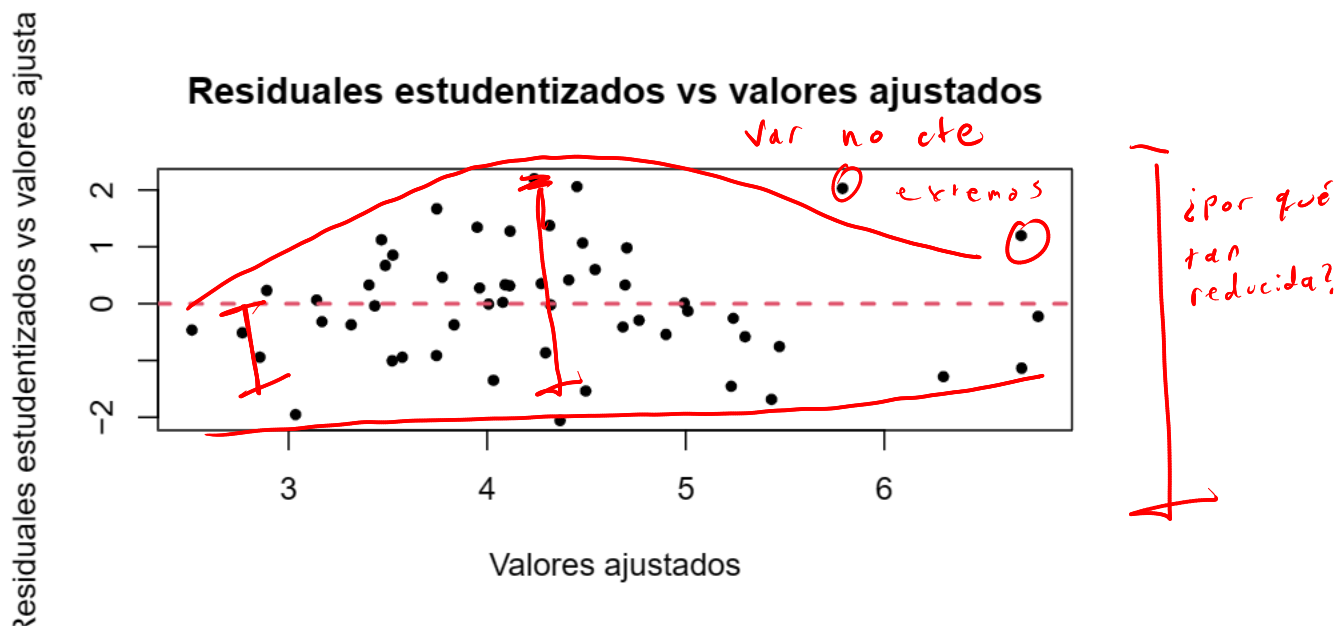


Figura 2: Gráfico residuales estudentizados vs valores ajustados

Después de analizar el gráfico, no encontramos ningún patrón que sugiera un aumento o disminución en la variabilidad de los datos dado que los puntos en la gráfica se encuentran muy dispersos, por lo que no encontramos ninguna señal que indique que la variable es constante.

→ y al contrario

→ Dicen que es constante

¿Al final a qué llegan?

## 4.2. Verificación de las observaciones

### 4.2.1. Datos atípicos 3 pt

#### Residuales estudentizados

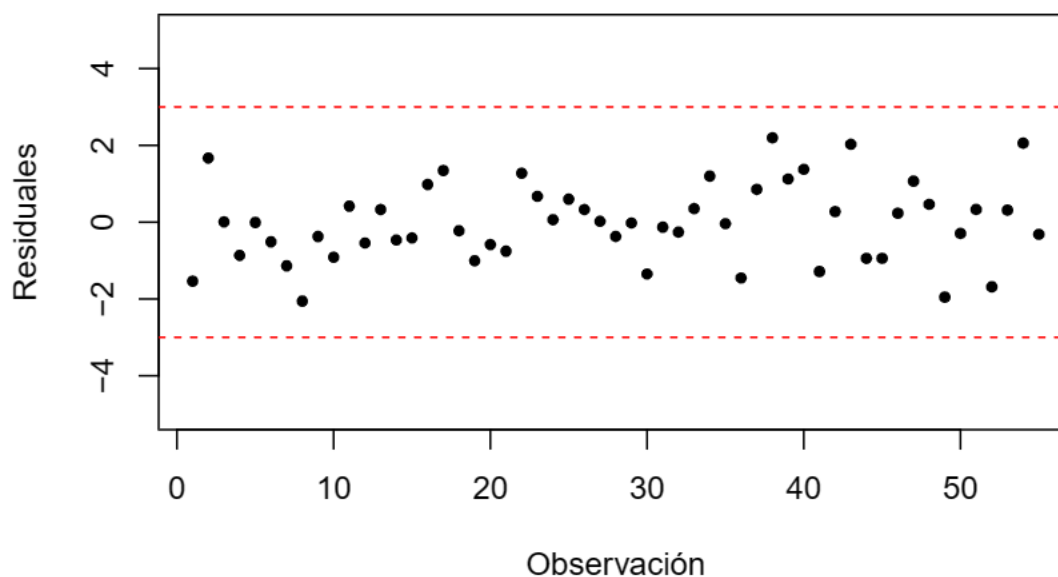


Figura 3: Identificación de datos atípicos

Según el criterio para evaluar datos atípicos, se considera que una observación es atípica cuando  $|r_{estud}| > 3$  (el valor absoluto del residual estudiantizado es mayor que 3). En este caso, al observar la gráfica se puede evidenciar que los valores del residual están entre -3 y 3, por lo que se considera que no hay datos atípicos. ✓

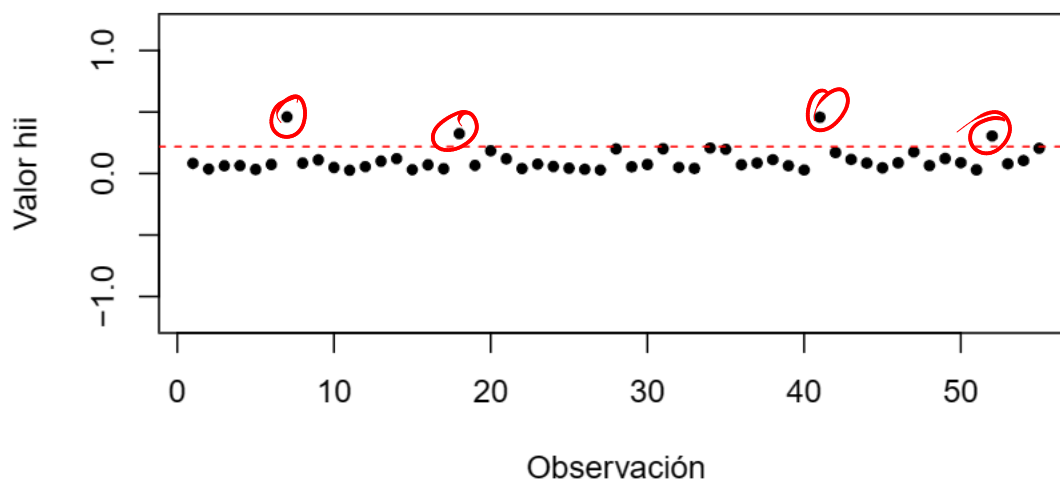
4.2.2. Puntos de balanceo 2 pt**Gráfica de hii para las observaciones**

Figura 4: Identificación de puntos de balanceo

Cuadro 5: Tabla de puntos de balanceo

	Residuales estudentizados	D.cook	Valor hii	Dffits
7	-1.1362	0.1835	0.4603	-1.0524
18	-0.2243	0.0040	0.3242	-0.1538
41	-1.2850	0.2321	0.4575	-1.1882
52	-1.6845	0.2075	0.3049	-1.1377

Al observar la gráfica anterior  $h_{ii}$  donde la línea punteada roja representa el valor

la línea representa  
la igualdad  $h_{ii} = 2p/n$   
no la desigualdad

$$h_{ii} > 2p/n = 2p/n = 2(6/55) = 0.21818 \quad \checkmark$$

nos permite notar que existen 4 puntos (7,18,41,52) de balanceo los cuales se pueden evidenciar en el cuadro 5 “tabla de puntos de balanceo”.

¿Qué causan estos puntos?

## 4.2.3. Puntos influyentes

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

### Gráfica de distancias de cook

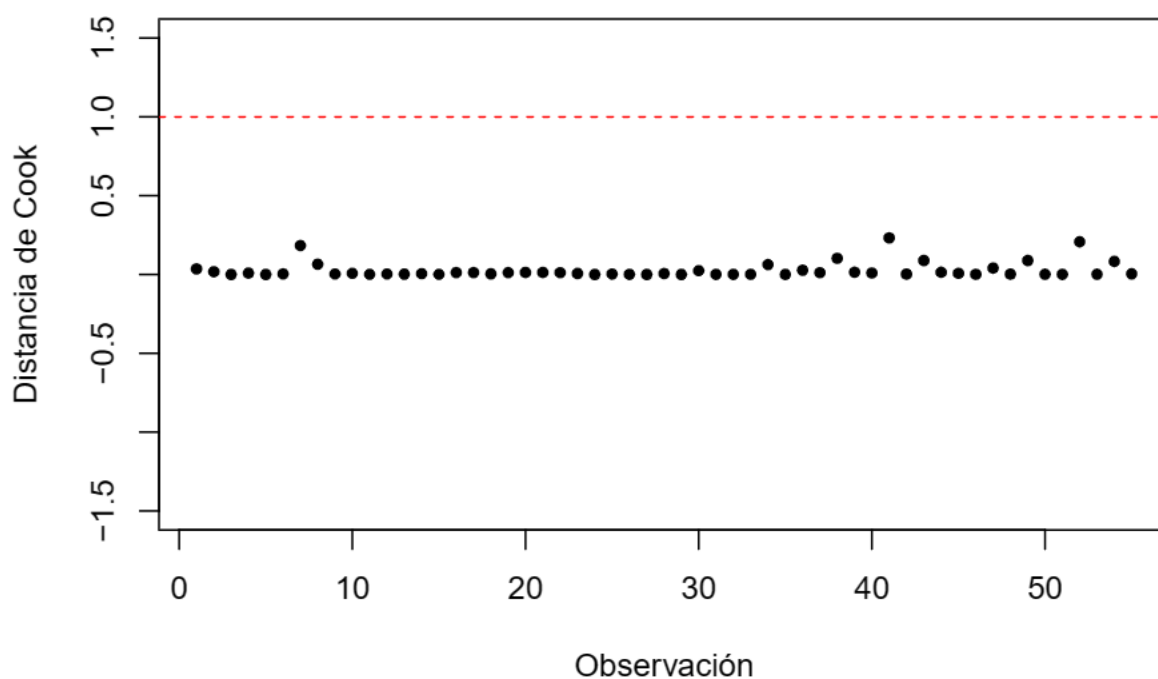


Figura 5: Criterio distancias de Cook para puntos influyentes

Según el criterio de Cook que dice que la observación  $i$  será influyente si  $Di > 1$ , podemos concluir que no hay puntos influyentes. ✓ 2 p +

### Gráfica de observaciones vs Dffits

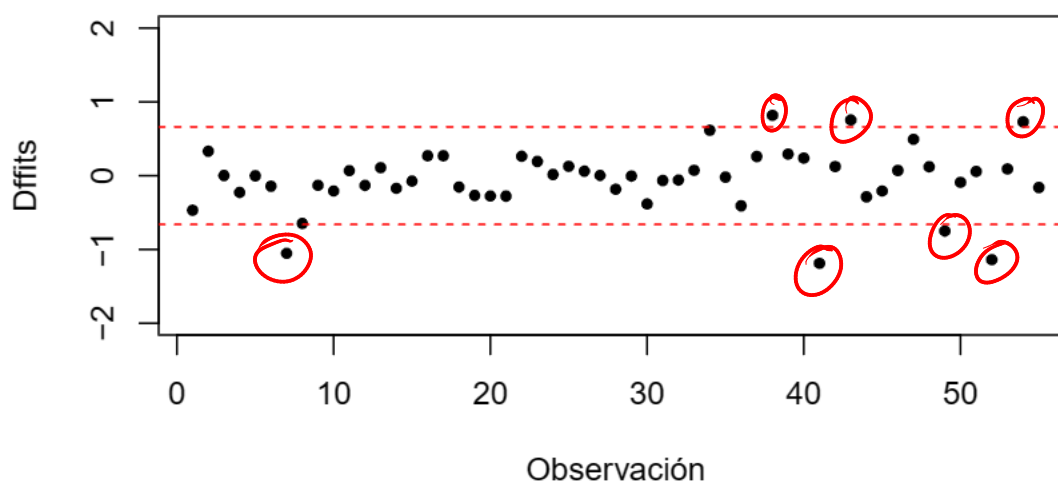


Figura 6: Criterio Dffits para puntos influyentes

De acuerdo con la columna Dffits del cuadro 6 “Tabla de puntos influyentes bajo el criterio

de Dffits” mostrada a continuación y según el criterio de este mismo el cual dice que para cualquier punto cuyo  $|Dffits| > 2\sqrt{\frac{p}{n}}$  la observación  $i$  será un punto influyente, es decir, si  $|Dffits| > 2\sqrt{\frac{p}{n}} = 0.6606$ , por lo que las observaciones 7,38,41,43,49,52 y 54 son puntos influyentes. ✓

Cuadro 6: Tabla de puntos influyentes bajo el criterio de Dffits

	Residuales estudentizados	D.cook	Valor hii	Dffits
7	-1.1362	0.1835	0.4603	-1.0524
38	2.1980	0.1026	0.1130	0.8180
41	-1.2850	0.2321	0.4575	-1.1882
43	2.0288	0.0889	0.1147	0.7554
49	-1.9526	0.0884	0.1221	-0.7504
52	-1.6845	0.2075	0.3049	-1.1377
54	2.0593	0.0829	0.1050	0.7303

¿Qué causan estos puntos?

## 5. Conclusiones

### 1. Del modelo de regresión

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

es posible concluir que las variables  $\beta_3$ ,  $\beta_4$  y  $\beta_5$  correspondientes al número de camas, censo promedio diario y número de enfermeras respectivamente, son indispensables para estimar el porcentaje del riesgo de infección en los hospitales, puesto que estos parámetros son significativos y a su vez interpretables. ✓

2. Al realizar la prueba de significancia para las tres variables predictoras con los valores P más grandes, se encontró que la duración de la estadía ( $\beta_1 =$ ), la rutina de cultivos ( $\beta_2 =$ ) y el número de camas ( $\beta_3 =$ ), no son estadísticamente significativos para predecir el promedio del riesgo de infección hospitalaria, lo que lleva a la conclusión de que estas variables se descartan del modelo por su falta de aporte significativo. ✓

3. Tras la realización de la prueba de hipótesis en forma matricial, fue posible resolver el sistema de ecuaciones que nos llevó a una nueva expresión del modelo reducido utilizando solamente las variables que son significativas. ✓

4. Después de realizar un análisis de la validez de los supuestos del modelo, incluyendo la normalidad de los residuales y la varianza constante, se encontró que no se cumplen dichos supuestos. Además, se detectaron puntos de balanceo y puntos influyentes, lo que sugiere que el modelo no es el más adecuado para describir el comportamiento de la variable respuesta (riesgo de infección). ✓