

# Trabajo 1

3,7

Estudiantes

**Arturo De Jesus Rangel Julio**

Equipo #48

Docente

**Francisco Javier Rodriguez Cortes**

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 pruebas de hipótesis y modelo reducido . . . . . | 5        |
| 2.2. Estadístico de prueba y conclusión . . . . .                   | 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 influyentes . . . . .                                 | 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 para los beta . . . . .         | 3 |
| 2. | Tabla ANOVA para el modelo . . . . .             | 4 |
| 3. | Resumen de los coeficientes beta . . . . .       | 4 |
| 4. | Resumen tabla de todas las regresiones . . . . . | 6 |

# 1. Pregunta 1

14,5 pt

Teniendo en cuenta la base de datos 48, encontramos que tiene 5 variables regresoras las cuales 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

Así, al plantear el modelo de regresión lineal múltiple tenemos:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + \beta_4 X_{4i} + \beta_5 X_{5i} + \varepsilon_i; \varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2); 1 \leq i \leq 50$$

## 1.1. Modelo de regresión

1,5 pt

De esta manera al traer el modelo ajustado, obtenemos los siguientes coeficientes para los valores de beta:

Cuadro 1: Tabla de valores para los beta

| V. parámetros |         |
|---------------|---------|
| $\beta_0$     | -1.8188 |
| $\beta_1$     | 0.1857  |
| $\beta_2$     | 0.0357  |
| $\beta_3$     | 0.0511  |
| $\beta_4$     | 0.0165  |
| $\beta_5$     | 0.0012  |

entonces, el resultado del modelo de regresión ajustado es:

$$\hat{Y}_i = -1.8188 + 0.1857X_{1i} + 0.0357X_{2i} + 0.0511X_{3i} + 0.0165X_{4i} + 0.0012X_{5i} + \varepsilon_i, \varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2); 1 \leq i \leq 50$$

No va en ec. ajustada

## 1.2. Significancia de la regresión

4 pt

Así si queremos analizar la significancia de la regresión, se debe plantear un juego de hipótesis que ayuden a interpretar su validez:

$$\begin{cases} H_0 : \beta_1 = \beta_2 = \beta_3 = \beta_4 = \beta_5 = 0 \\ H_1 : \text{Algún } \beta_j \text{ distinto de } 0 \text{ para } j=1, 2, \dots, 5 \end{cases}$$

SE utiliza un estadístico de prueba que es:

$$F_0 = \frac{\cancel{MST}}{MSE} \stackrel{H_0}{\sim} f_{5,44} \quad F_0 = \frac{MSR}{MSE} \neq \frac{MST}{MSE} \quad (1)$$

Presentando la tabla Anova:

Cuadro 2: Tabla ANOVA para el modelo

|           | Suma Cuadratica | G.libertad. | media cuadratica | $F_0$   | P-valor     |
|-----------|-----------------|-------------|------------------|---------|-------------|
| Regresión | 47.573          | 5           | 9.51460          | 10.6996 | 9.03818e-07 |
| Error     | 39.127          | 44          | 0.88925          |         |             |

Viendo los resultados obtenidos en la tabla Anova, se observa un valor P es un valor muy pequeño que lo hace muy cercano a 0, por lo que para cualquier nivel de significancia ( $\alpha$ ) se rechaza la hipótesis nula en la que  $\beta_j = 0$  con  $1 \leq j \leq 5$ , por tanto se aceptan la hipótesis alternativa que nos dice que algún  $\beta_j \neq 0$ , por esta razón la regresión general es significativa.

### 1.3. Significancia de los parámetros

En el siguiente cuadro podemos encontrar detalladamente la información de los parámetros estudiados, de esta manera podremos determinar con mayor precisión cuáles de ellos son significativos y cuáles no.

Prueba de hipótesis para los parámetros:

$$\begin{cases} H_0 : \beta_j = 0 \\ H_1 : \beta_j \neq 0 \end{cases}$$

Cuadro 3: Resumen de los coeficientes beta

|           | $\hat{\beta}_j$ | $SE(\hat{\beta}_j)$ | $T_{0j}$ | P-valor |
|-----------|-----------------|---------------------|----------|---------|
| $\beta_0$ | -1.8188         | 1.8419              | -0.9875  | 0.3288  |
| $\beta_1$ | 0.1857          | 0.0812              | 2.2857   | 0.0271  |
| $\beta_2$ | 0.0357          | 0.0335              | 1.0665   | 0.2920  |
| $\beta_3$ | 0.0511          | 0.0150              | 3.4040   | 0.0014  |
| $\beta_4$ | 0.0165          | 0.0074              | 2.2220   | 0.0315  |
| $\beta_5$ | 0.0012          | 0.0007              | 1.6785   | 0.1003  |

De la tabla anterior al analizar los P-valores resultantes, permiten llegar a la conclusión de que con un nivel de significancia  $\alpha = 0.05$ , para los parámetros  $\beta_2$  y  $\beta_5$  se acepta la hipótesis nula ( $H_0$ ) y para los Parámetros  $\beta_1$ ,  $\beta_3$  y  $\beta_4$  se acepta la hipótesis alternativa ( $H_1$ ), por lo que son significativos. Ya que sabemos que para que los P-valores cumplan la prueba de significancia deben ser menores a  $\alpha$ .

#### 1.4. Interpretación de los parámetros 0 pt

$\hat{\beta}_1$ : Por cada unidad de aumento en la Probabilidad promedio estimada de adquirir infección en el hospital ( $Y$ ), la Duración promedio de la estadía de todos los pacientes en el hospital ( $X_1$ ) aumenta significativamente en 0.1857 unidades, cuando los valores en las demás predictoras se mantiene fijo. X

$\hat{\beta}_3$ : Por cada unidad de aumento en la Probabilidad promedio estimada de adquirir infección en el hospital ( $Y$ ), el Número promedio de camas en el hospital durante el periodo del estudio ( $X_3$ ) aumenta significativamente en 0.0511 unidades, cuando los valores en las demás predictoras se mantiene fijo. X

$\hat{\beta}_4$ : Por cada unidad de aumento en la Probabilidad promedio estimada de adquirir infección en el hospital ( $Y$ ), el Número promedio de pacientes en el hospital por día durante el periodo del estudio ( $X_4$ ) aumenta significativamente en 0.0165 unidades, cuando los valores en las demás predictoras se mantiene fijo. X

unidad de aumento en  $X_j$ ,  
y aumenta  $\beta_j$  unidades  
3 pt

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

El modelo tiene un coeficiente de determinación múltiple  $R^2 = 0.5487$ , de lo cual identificamos que el porcentaje de la variabilidad de el Riesgo de infección explicado por el modelo propuesto en el presente es 54.77%.

¿cómo se calcula?

## 2. Pregunta 2 3 pt

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

Encontramos que las 3 covariable con el P-valor más alto en el modelo fueron  $X_2$ ,  $X_4$ ,  $X_5$ , así vemos que a partir de la tabla de todas las regresiones posibles se pretende hacer la siguiente prueba de hipótesis:

$$\begin{cases} H_0 : \beta_2 = \beta_4 = \beta_5 = 0 \\ H_1 : \text{Algún } \beta_j \text{ distinto de } 0 \text{ para } j = 2, 4, 5 \end{cases}$$

Cuadro 4: Resumen tabla de todas las regresiones

|                 | $SSE$  | Covariables en el modelo |    |    |    |    |
|-----------------|--------|--------------------------|----|----|----|----|
| Modelo completo | 39.127 | X1                       | X2 | X3 | X4 | X5 |
| Modelo reducido | 45.921 | X1                       | X3 |    |    |    |

asi encontramos que un modelo reducido para la prueba de significancia del subconjunto es:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_3 X_{3i} + \varepsilon_i; \varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2); 1 \leq i \leq 50$$

## 2.2. Estadístico de prueba y conclusión

Se construye el estadístico de prueba de la siguiente forma:

$$\begin{aligned}
 F_0 &= \frac{(SSE(\beta_0, \beta_1, \beta_3) - SSE(\beta_0, \dots, \beta_5))/3}{MSE(\beta_0, \dots, \beta_5)} \stackrel{H_0}{\sim} F_{3,44} \\
 &= \frac{45.921 - 39.127}{0.88925} \\
 &= 1.921
 \end{aligned}$$

$$F_0 < F_{0.95, 3, 44}$$

Ahora, al realizar una comparacion entre  $F_0$  con  $f_{0.95, 3, 44} = 2.8165$ , se puede encontrar que  $F_0 < f_{0.95, 3, 44}$ . Debido a esto, el subconjunto es significativo. por tal razon no es posible descartar las variables del modelo ya que se rechaza la hipotesis nula ( $H_0$ ) y almenos unas de sus variables es distinto de 0, por lo que se acepta la hipotesis alternativa ( $H_1$ ).

## 3. Pregunta 3

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

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

$$\begin{cases} H_0 : \beta_3 = 5\beta_4; \beta_2 = 2\beta_5 \\ H_1 : \text{Alguna de las igualdades no se cumple} \end{cases}$$

reescribiendo matricialmente:

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

donde la  $\mathbf{L}$  dada por

$$L = \begin{bmatrix} 0 & 0 & 0 & 1 & -5 & 0 \\ 0 & 0 & 1 & 0 & 0 & -2 \end{bmatrix}$$

El modelo reducido está dado por:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_4 X_{3,4i}^* + \beta_5 X_{2,5i}^* + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2); \quad 1 \leq i \leq 50$$

Donde  $X_{3,4i}^* = 5X_{3i} + X_{4i}$  y  $X_{2,5i}^* = 2X_{2i} + X_{5i}$

### 3.2. Estadístico de prueba

El estadístico de prueba  $F_0$  está dado por:

$$F_0 = \frac{(SSE(MR) - SSE(MF))/2}{MSE(MF)} \stackrel{H_0}{\sim} f_{2,44} \quad (3)$$

&=

$$F_0 = \frac{(SSE(MR) - 39.127)/2}{0.88925} \stackrel{H_0}{\sim} f_{2,44} \quad (4)$$

## 4. Pregunta 4

### 4.1. Supuestos del modelo

#### 4.1.1. Normalidad de los residuales

Para la verificación de este supuesto, se hará el planteamiento de la siguiente prueba de hipótesis ~~shapiro-wilk~~, de igual manera se agregará un gráfico cuantil-cuantil:

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



### Normal Q-Q Plot of Residuals

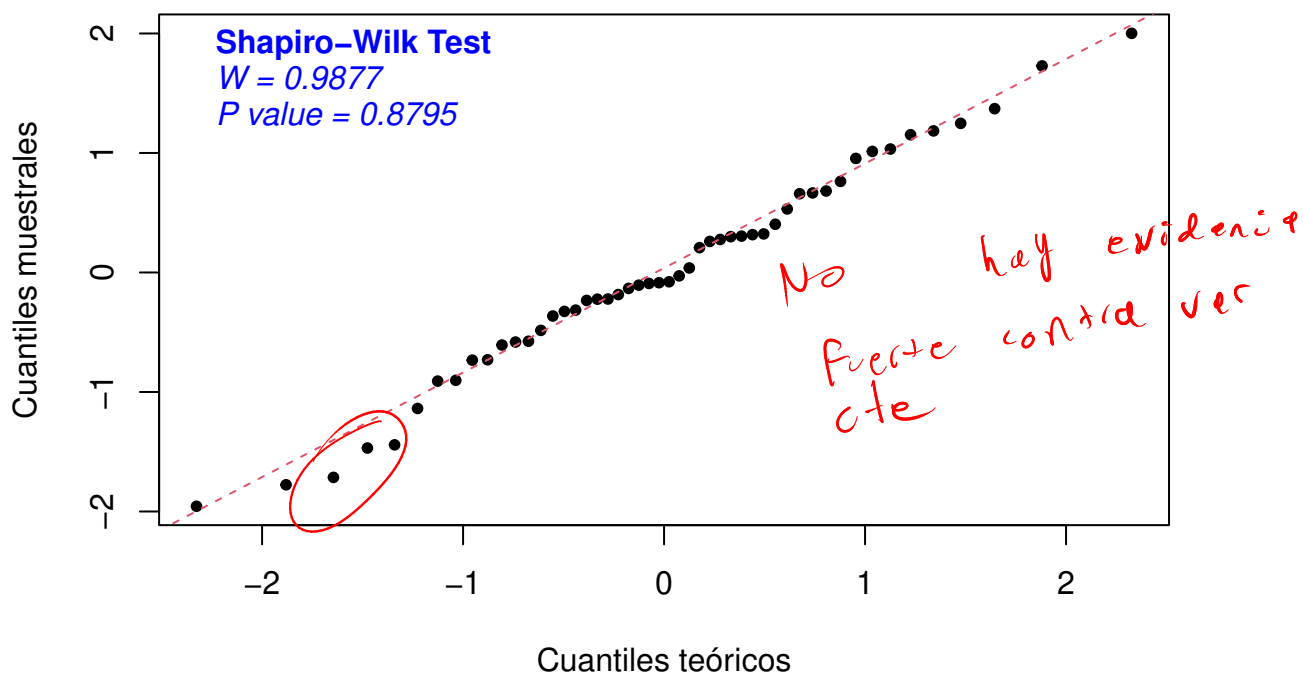


Figura 1: Gráfico cuantil-cuantil y normalidad de residuales

Al ser el P-valor aproximadamente igual a 0.8795 y al ser  $\alpha = 0.05$ , el p-valor es mucha mas alto que el  $\alpha$  por lo que no se rechaza la hipótesis nula ( $H_0$ ) y al ser tan cercano a 1 esto muestra que cumple con la prueba de normalidad, por lo que los datos se distribuyen normal con ~~media  $\mu$  y varianza  $\sigma^2$~~ . Pero en este caso hay que tener en cuenta que la gráfica aunque no se logran identificar colas muy pesadas, claramente se pueden ver patrones marcados e irregulares, por lo que se toma la decisión de rechazar el cumplimiento del supuesto ya que la gráfica tiene mas importancia a la hora de decidir. ahora vamos a proceder a ver si la varianza es constante.

Te valgo el análisis pero en esta base no se rechaza.

No se está probando media y var ctes.

## 4.1.2. Varianza constante

10+

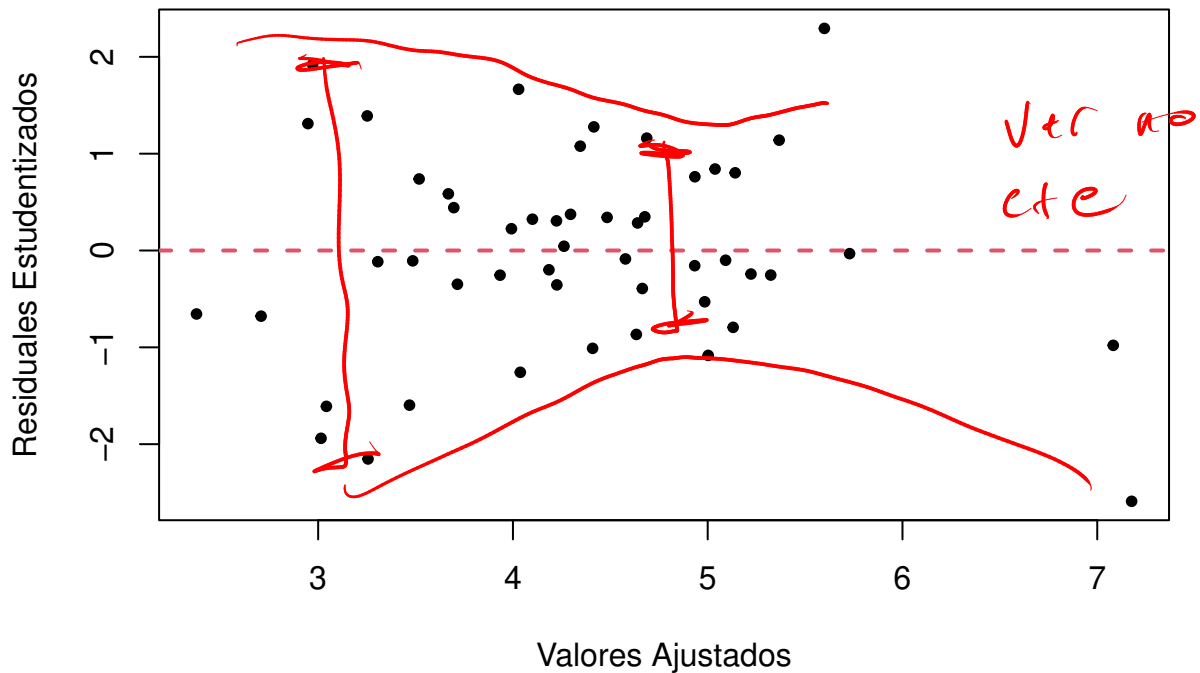
**Residuales Estudentizados vs Valores Ajustados**

Figura 2: Gráfico residuales estudentizados vs valores ajustados

Al trazar las dos líneas una por debajo y una arriba de los datos, se puede identificar que no existen comportamientos que me hagan rechazar una varianza constante. De igual forma no vemos patrones de aumento o disminución de la varianza. también podemos observar que se tiene media 0.

X

## 4.2. Verificación de las observaciones

### 4.2.1. Datos atípicos

3pt

#### Residuales estudentizados

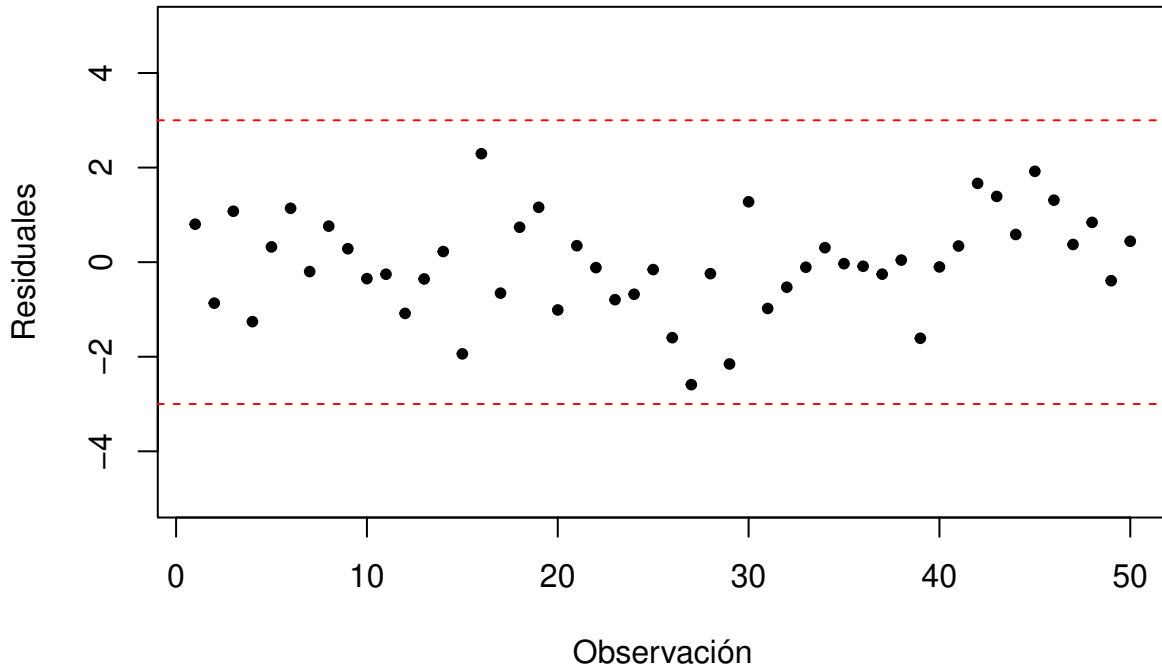


Figura 3: Identificación de datos atípicos

podemos ver en la grafica anterior que no se observan datos atipicos en le conjunto de datos pues ningún residual estudentizado sobrepasa el criterio de  $|r_{estud}| > 3$ .



## 4.2.2. Puntos de balanceo

1p+

## Gráfica de hii para las observaciones

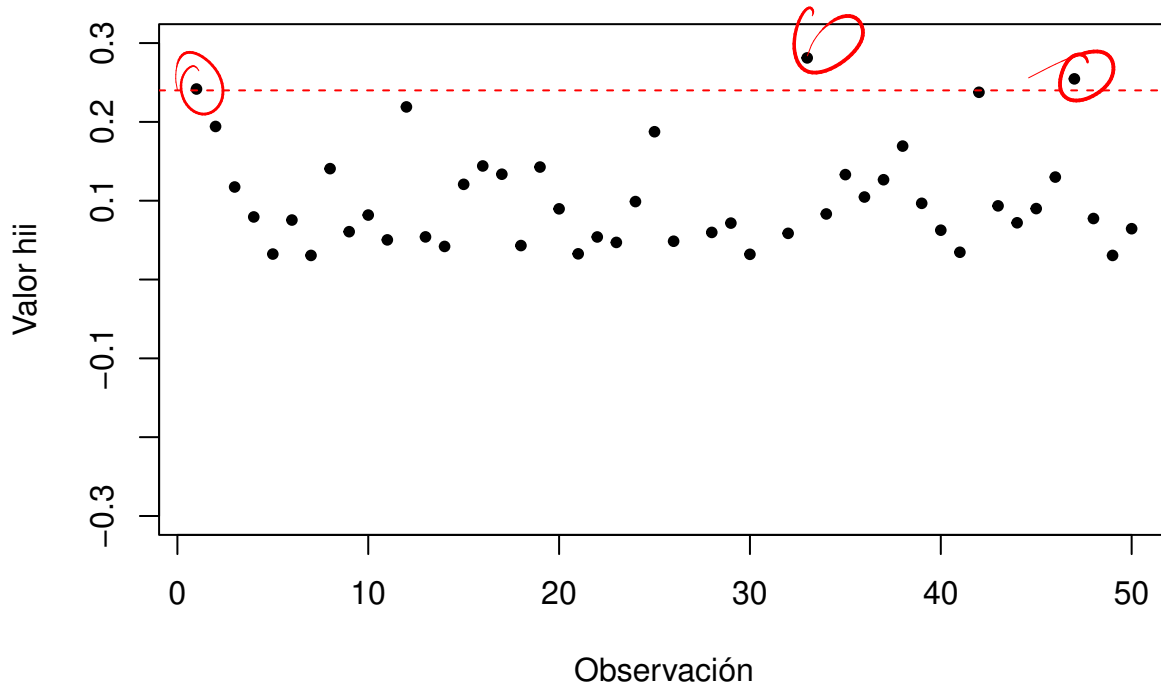


Figura 4: Identificación de puntos de balanceo

Al observar la gráfica de observaciones vs valores  $h_{ii}$ , donde la línea punteada roja representa el valor  $h_{ii} = 2\frac{p}{n} = 0.24$ , se puede apreciar que existen 4 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.

Tienes de hecho 5  
puntos de balanceo.

↓  
¿wá! No  
está.

|    | hii.value | Dffits  |
|----|-----------|---------|
| 1  | 0.2417    | 0.4512  |
| 27 | 0.4710    | -2.6243 |
| 31 | 0.6027    | -1.2056 |
| 33 | 0.2812    | -0.0662 |
| 47 | 0.2546    | 0.2162  |

## 4.2.3. Puntos influyentes

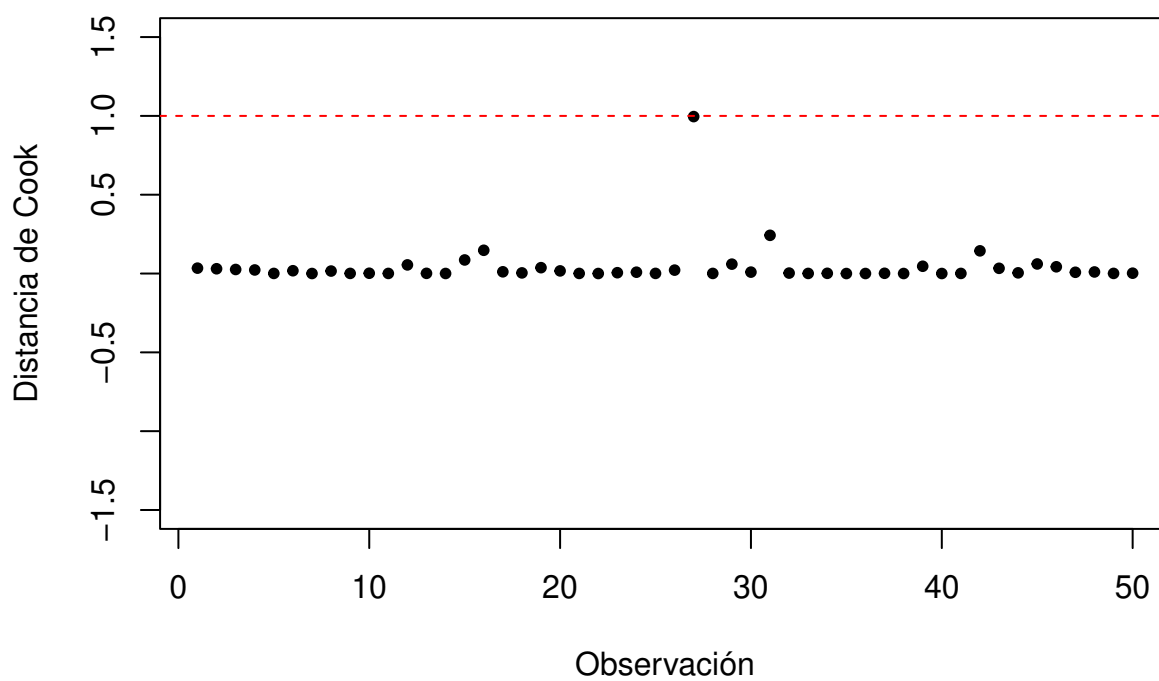
**Gráfica de distancias de Cook**

Figura 5: Criterio distancias de Cook para puntos influyentes

Se ven 4 en  
gráfica y tienen 5,  
no son congruentes.

### Gráfica de observaciones vs Dffits

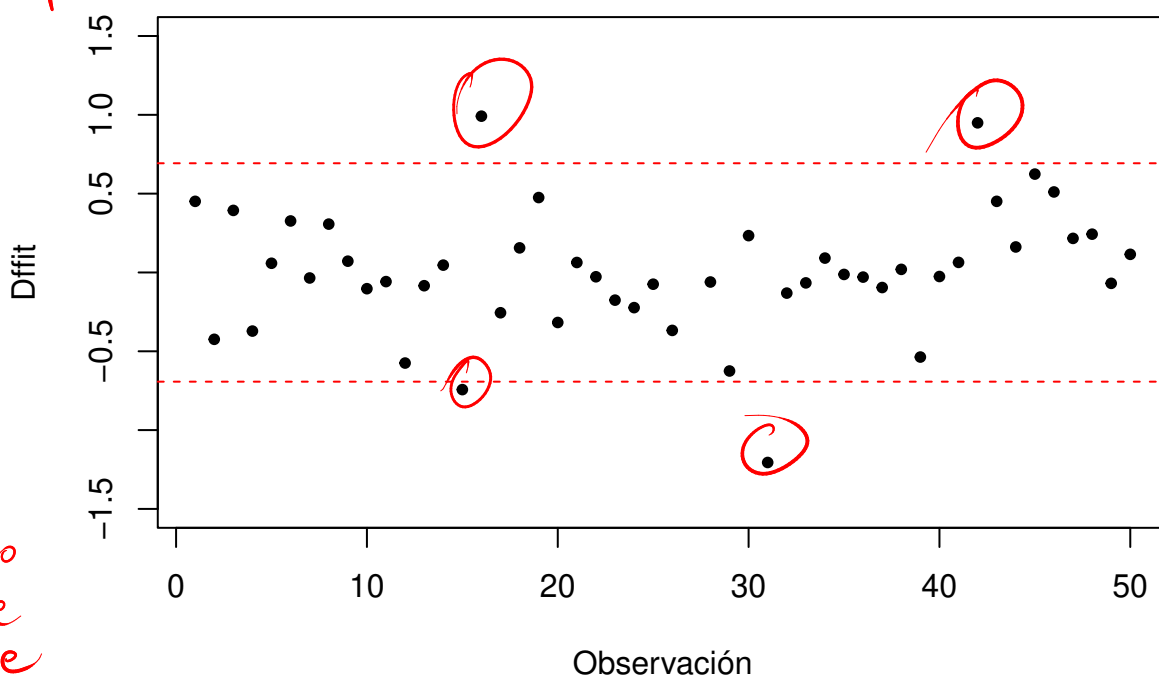


Figura 6: Criterio Dffits para puntos influyentes

| ##    | res.stud | Cooks.D | hii.value | Dffits  |
|-------|----------|---------|-----------|---------|
| ## 15 | -1.9393  | 0.0861  | 0.1208    | -0.7431 |
| ## 16 | 2.2936   | 0.1476  | 0.1441    | 0.9916  |
| ## 27 | -2.5901  | 0.9955  | 0.4710    | -2.6243 |
| ## 31 | -0.9793  | 0.2425  | 0.6027    | -1.2056 |
| ## 42 | 1.6650   | 0.1440  | 0.2376    | 0.9493  |

→ tabla, no salida de R  
¿cuáles? → ¿cuánto da?

Como se puede ver, las observaciones 15, 16, 27, 31, 42 son puntos influyentes según el criterio de Dffits, el cual 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, ninguno de los datos cumple con serlo.

¿Qué causan?

### 4.3. Conclusión

3 pt

El modelo no es válido debido a que no se cumple el supuesto de normalidad ya que a través de gráfica llegamos a dicha conclusión, esto se puede dar debido a los puntos de balance que pueden ver por fuera de la condición propuesta en gráfica de que observa los Valores de hii.

