

Tercer taller

miércoles, 18 de diciembre de 2024 2:00 p. m.

1. Determine el valor de verdad de las siguientes afirmaciones.

- (a) Bajo los supuestos del modelo de regresión lineal múltiple, se cumple que la respuesta media estimada $\hat{Y}_0 \sim N(E[Y|X_0], \sigma^2 X_0 (X'X)^{-1} X_0')$ donde además, \hat{Y}_0 es un estimador insesgado para Y_0 .
 $Y_0 - \hat{Y}_0$
 Falso
 MSE
 Varianza estimada
- (b) El error de predicción $\hat{Y}_0 - Y_0$ tiene una varianza asociada dada por $\sigma^2 X_0 (X'X)^{-1} X_0'$ al igual que \hat{Y}_0 .
 $Var[Y_0 - \hat{Y}_0] = Var[Y_0] + Var[\hat{Y}_0] = \sigma^2 + \sigma^2 X_0 (X'X)^{-1} X_0' = \sigma^2 [1 + X_0 (X'X)^{-1} X_0']$
 Falso
- (c) Si se sabe que $X_0 = [1, x_{01}, x_{02}, \dots, x_{0k}]'$ es un punto en que no se comete extrapolación, entonces es correcto afirmar que $X_0 (X'X)^{-1} X_0' < 1$.
 $X_0 (X'X)^{-1} X_0' < 1$
 Verd.
 $H = X(X'X)^{-1}X'$
- (d) Una observación atípica está separada del resto de las observaciones en su valor de respuesta Y aunque ~~no~~ afecta los resultados del ajuste. Su evaluación se realiza a través del residual estandarizado $|d_i| > 3$.
 Falso
 $0 \leq h_{ii} \leq 1$
 \hat{e}_i
 $d_i = \frac{e_i}{\sqrt{MSE}}$
 $|d_i| > 3 \vee |r_i| > 3$
 $\sim N(0,1)$
 99.75%
 -3σ
 3σ
- (e) Se cumple que una observación i es de balanceo si está definida en el espacio de la respuesta Y y se cumple que $h_{ii} > 2p/n$, controlando propiedades del modelo.
 Falso
 Están definidos en el espacio de las covariables
 No las estimaciones puntuales
 Afectar el error estándar estimado
 β_0 ; Afectar R^2
- (f) Una observación i es influyente si hala el modelo en su dirección, siendo inusual en el espacio de predictor y la respuesta.
 Verdadera

2. Seleccione las expresiones adecuadas que se muestra a continuación, interprételas y corrija las expresiones incorrectas.

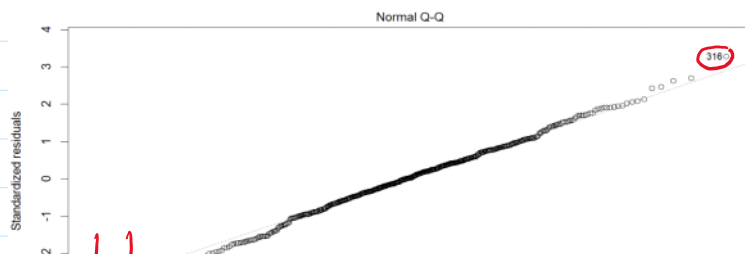
- a. $d_i = \frac{e_i}{\sqrt{MSE}}$
 Influencia
- b. $r_i = \frac{d_i}{\sqrt{1-h_{ii}}}$
 $r_i = \frac{e_i}{\sqrt{MSE(1-h_{ii})}}$
- c. $DFBETAS_{j(i)} = \frac{\hat{\beta}_j - \hat{\beta}_{j(i)}}{\sqrt{MSE_{(i)} c_{jj}}}$
 Diagonal principal $(X'X)^{-1}$
- d. $DFFITS_i = \frac{\hat{y}_i - \hat{y}_{i(i)}}{\sqrt{MSE_{(i)}}}$
 hii

• Recordar para la detección de puntos atípicos: $|d_i| > 3$; $|r_i| > 3$: **Más conservador: $|d_i| > 2$; $|r_i| > 2$**

1. Verifique los supuestos del modelo de regresión, esto es, $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$ a partir de los procedimientos apropiados para ello.

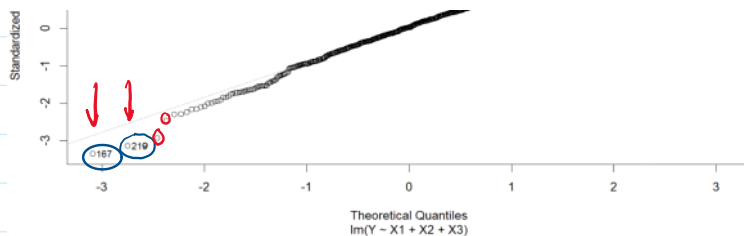
• Verificar: Normalidad; Varianza constante, media cero; (La independencia está dada)

• $\hat{y}_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + \epsilon_i$; $\epsilon_i \sim N(0, \sigma^2)$; Normalidad:



• Criterio Analítico: Shapiro-Wilks:

$H_0: \epsilon_i \sim \text{Normal}$ $H_1: \epsilon_i \not\sim \text{Normal}$	Shapiro-Wilk normality test data: modelo\$residuals W = 0.9957, p-value = 0.1877
--	---



$H_0: \epsilon_i \sim \text{Normal}$

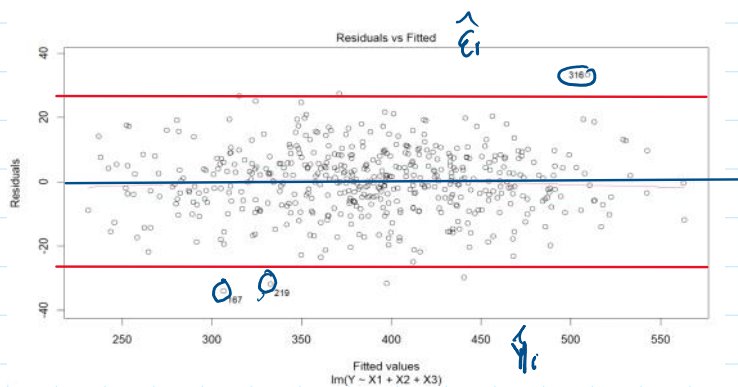
data: modelo\$residuals
W = 0.9957, p-value = 0.1877

$P_{val} < \alpha$ $\alpha = 0.05$

• No se rechaza H_0 : Se puede concluir a un nivel de sign 5% que los errores tienen una $\sim \text{Normal}$.

• Homogeneidad variancia

• Verificar el supuesto de media cono



2. Identifique puntos atípicos y puntos de balanceo a través de un criterio gráfico y analítico. ¿Podrían ser estos puntos a su vez influyentes?

$|d_i| > 3$; $|r_i| > 3$

Balanceo: $|h_{ii}| > 2p/n$: $2p/n < 1$

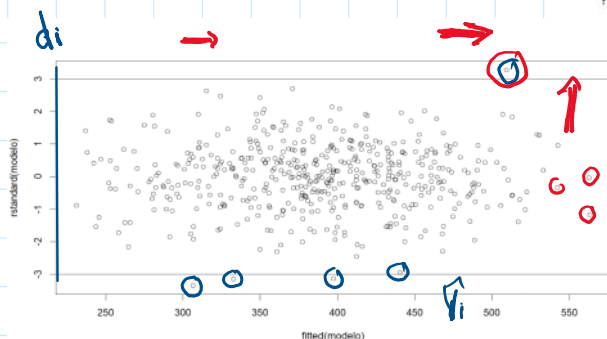
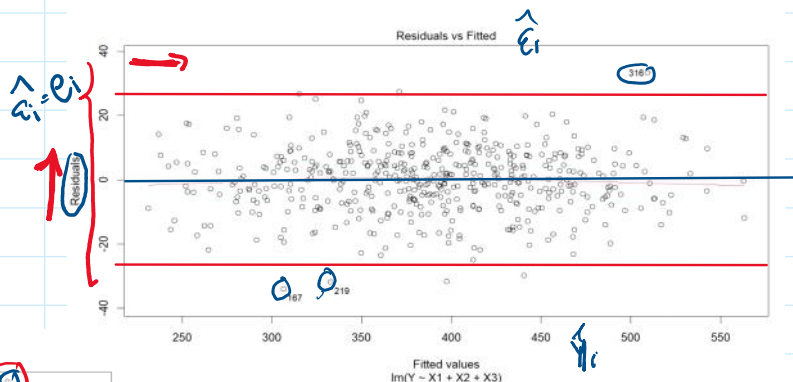
$p = \# \text{ predictores (4): } \beta_0, \beta_1, \beta_2, \beta_3$
 $n = \# \text{ muestra (500)}$

> balanceo
75 102 194 274 311

Pueden afectar R^2 y la estimación del error estándar asociado β_j

> atipicos_estandarizados
16 167 219 316
16 167 219 316
> atipicos_estudentizados
16 167 219 316
16 167 219 316

Puntos atípicos



• Una observación influyente, se cataloga como tal si es inusual en X y Y

3. Identifique puntos influyentes. Compare los criterios empleados para ello.

A partir de 3 criterios: Cook's, DFBetas, DFFITS
 $\hookrightarrow \hat{\beta}$ $\hookrightarrow \hat{\beta}_j$ $\hookrightarrow \text{Ajuste } Y_i$

- Cook's: $D_i > 1$
- DFBetas: $|DFBetas| > 2/\sqrt{n}$
- DFFITS: $|DFFITS| > 2/\sqrt{p}$

Según los criterios anteriores, para que una observación sea influyente, no necesariamente debe ser catalogada como tal en todos los criterios anteriores.

- DFBetas: $|DFBetas| > 2/\sqrt{n}$
- DFFITS: $|DFFITS| > 2\sqrt{p/n}$

sea influyente, no necesariamente debe ser catalogada como tal en todos los criterios anteriores.

No hay influencia

Según DFBetas

```
> which(cooks > 1) # Verificar cooks
named integer(0)
```

```
> which(abs(DFBetas) > 2/sqrt(n)) # Verificar DFBETAS
[1] 22 35 82 107 109 114 117 122 124 131 134
[12] 139 144 154 167 178 183 187 194 214 219 247
[23] 255 274 276 289 316 317 320 362 366 405 418
[34] 425 441 516 522 542 545 575 582 607 609 613
[45] 614 620 622 654 681 687 694 709 714 719 722
[56] 736 747 750 778 789 794 798 808 839 845 856
[67] 862 866 880 953 961 1001 1016 1030 1045 1075 1107
[78] 1113 1119 1134 1135 1173 1178 1181 1219 1236 1237 1247
[89] 1276 1280 1289 1294 1308 1316 1339 1397 1405 1435 1441
[100] 1454 1461 1516 1522 1529 1530 1534 1535 1545 1559 1575
[111] 1582 1609 1614 1622 1624 1631 1634 1644 1654 1667 1673
[122] 1681 1683 1694 1722 1735 1736 1739 1747 1754 1755 1776
[133] 1794 1816 1817 1820 1838 1849 1856 1862 1921 1925 1941
[144] 1953
```

```
> which(abs(DFFITS) > (2 * sqrt(p/n))) # Verificar DFFITS
1 16 22 29 45 75 107 109 113 122 154 167 173 214 219
1 16 22 29 45 75 107 109 113 122 154 167 173 214 219
222 236 247 276 289 294 316 317 356 362 453
222 236 247 276 289 294 316 317 356 362 453
```

Según DFFITS → \hat{y}_i

	dfb.1	dfb.X1	dfb.X2	dfb.X3	dffit	cov.r	cook.d	hat
1	-0.03	-0.08	0.12	-0.02	-0.18	0.97_*	0.01	0.01
8	0.04	-0.06	0.03	-0.06	-0.15	0.96_*	0.01	0.00
16	0.02	-0.10	0.17	-0.14	-0.28_*	0.94_*	0.02	0.01
29	0.06	0.05	-0.07	-0.13	-0.18	0.97_*	0.01	0.01
61	0.02	-0.03	-0.05	0.00	-0.11	0.98_*	0.00	0.00
113	0.01	0.13	-0.13	-0.07	-0.24	0.95_*	0.01	0.01
119	-0.08	0.01	0.12	-0.03	-0.16	0.97_*	0.01	0.01
158	-0.01	0.06	0.01	-0.08	-0.14	0.97_*	0.00	0.00
167	-0.21	0.08	0.00	0.23	-0.29_*	0.93_*	0.02	0.01
173	0.08	-0.06	-0.16	0.11	0.24	0.96_*	0.01	0.01
219	-0.23	0.24	0.18	-0.04	-0.33_*	0.94_*	0.03	0.01
254	0.07	0.01	0.03	-0.11	0.16	0.97_*	0.01	0.00
289	0.21	-0.19	-0.10	-0.05	0.25	0.96_*	0.02	0.01
311	0.01	0.00	0.00	0.00	-0.01	1.03_*	0.00	0.02
316	-0.23	0.09	0.23	0.15	0.32_*	0.93_*	0.03	0.01
362	0.16	-0.15	0.05	-0.15	0.24	0.97_*	0.01	0.01

4. Realice inferencia para $\mathbf{x}_{01} = [1, 45.03, 80.88, 60.33]$ y $\mathbf{x}_{02} = [1, 77.08, 100, 13.76]$ con su respectivo intervalo de predicción. Verifique primero si no se trata de un punto de extrapolación.

$$\mathbf{x}_0'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_0 < \max\{\hat{h}_{ii}\} \quad \checkmark$$

Hiperplano (Ajustado)

$$\mathbf{x}_0'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_0 < \max\{\hat{h}_{ii}\} \quad \mathbf{x}_0 = [1, \dots, \dots]$$

• $\{f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)\}$

```
> summary(X)
(Intercept)      X1      X2      X3
Min.   :1   Min.   :21.06   Min.   :27.57   Min.   :23.66
1st Qu.:1   1st Qu.:34.10   1st Qu.:43.76   1st Qu.:38.14
Median :1   Median :47.65   Median :63.08   Median :51.79
Mean   :1   Mean   :47.86   Mean   :62.94   Mean   :51.71
3rd Qu.:1   3rd Qu.:60.62   3rd Qu.:81.80   3rd Qu.:65.36
Max.   :1   Max.   :76.19   Max.   :99.88   Max.   :81.00
```

```
> ifelse(t(x01)%solve(t(X)%X)%x01 < max(Hat_values), "Pertence a la region de diseno", "No pertenece")
[1,]
[1,] "Pertence a la region de diseno" X01 ✓✓
> ifelse(t(x02)%solve(t(X)%X)%x02 < max(Hat_values), "Pertence a la region de diseno", "No pertenece")
[1,]
[1,] "No pertenece" X02 ✓ (No se puede realizar predicciones)
```

Intervalo predicción:

$$\hat{y}_0 \pm t_{\alpha/2, n-p} \cdot s.e.(\hat{y}_0 - y_0)$$

$$\hat{y}_0 \pm t_{\alpha/2, n-p} \sqrt{\text{Var}(\hat{y}_0 - y_0)}$$

$$\hat{y}_0 \pm t_{\alpha/2, n-p} \sqrt{MSE[1] + \mathbf{x}_0'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_0}$$

```
fit      lwr      upr
1 445.887 425.8413 465.9327
```

\hat{y}_0 \downarrow l \downarrow u

El valor que podría tomar un Punto de predicción.