

Regresión Lineal Múltiple - Semana 07

Johnatan Cardona Jiménez

jcardonj@unal.edu.co

Profesor Asistente - Departamento de Estadística
Universidad Nacional de Colombia, Sede Medellín

Semestre 02-2024

Inferencias sobre la respuesta media y valores futuros

- Suponga que se desea estimar la respuesta media para los valores en las predictoras $X_1 = x_{01}, X_2 = x_{02}, \dots, X_k = x_{0k}$.
- Sea Y_0 la respuesta desconocida en tal conjunto de valores. **Si se define el vector fila:** $\underline{\mathbf{x}}_0 = [1 \quad x_{01} \quad x_{02} \quad \dots \quad x_{0k}]$, entonces se puede escribir $Y_0 = \underline{\mathbf{x}}_0 \underline{\boldsymbol{\beta}} + \varepsilon_0$, **por lo tanto la respuesta media en tal punto es:**

$$\mu_{Y|\underline{\mathbf{x}}_0} = E[Y|\underline{\mathbf{x}}_0] = \underline{\mathbf{x}}_0 \underline{\boldsymbol{\beta}} = \beta_0 + \beta_1 x_{01} + \beta_2 x_{02} + \dots + \beta_k x_{0k}.$$

Este valor es estimado por la correspondiente respuesta o valor ajustado, \hat{Y}_0 , que puede escribirse como:

$$\hat{Y}_0 = \underline{\mathbf{x}}_0 \hat{\underline{\boldsymbol{\beta}}} = \hat{\beta}_0 + \hat{\beta}_1 x_{01} + \hat{\beta}_2 x_{02} + \dots + \hat{\beta}_k x_{0k}.$$

\hat{Y}_0 tiene las siguientes propiedades:

- $E[\hat{Y}_0] = E[\underline{x}_0 \hat{\underline{\beta}}] = \underline{x}_0 E[\hat{\underline{\beta}}] = \underline{x}_0 \underline{\beta} = E[Y|\underline{x}_0]$, esto es, \hat{Y}_0 es un estimador insesgado de la respuesta media $E[Y|\underline{x}_0]$.
- $Var[\hat{Y}_0] = Var[\underline{x}_0 \hat{\underline{\beta}}] = \underline{x}_0 Var[\hat{\underline{\beta}}] \underline{x}_0' = \sigma^2 \underline{x}_0 (\mathbf{X}'\mathbf{X})^{-1} \underline{x}_0'$, que es estimada por:
 $\widehat{Var}[\hat{Y}_0] = MSE \underline{x}_0 (\mathbf{X}'\mathbf{X})^{-1} \underline{x}_0'$.
- Bajo el supuesto de normalidad en los errores, \hat{Y}_0 es una variable aleatoria normal, debido a que es una combinación lineal de los $\hat{\beta}_j$'s que también son normales.

En resumen:

$$\hat{Y}_0 \sim N\left(E[Y|\underline{x}_0], \sigma^2 \underline{x}_0 (\mathbf{X}'\mathbf{X})^{-1} \underline{x}_0'\right)$$

Luego, se tiene que el estadístico $T = \frac{\hat{Y}_0 - E[Y|\underline{x}_0]}{ee(\hat{Y}_0)} \sim t_{n-p}$, con $ee(\hat{Y}_0) = \sqrt{\widehat{Var}[\hat{Y}_0]}$, lo cual permite demostrar lo siguiente:

- Para la respuesta media $E[Y|\underline{x}_0]$ en un vector apropiado \underline{x}_0 .

Pruebas de hipótesis sobre la respuesta media para un nivel de significancia α

Juego de hipótesis	Estadístico de prueba	Criterio de rechazo
$H_0 : \mu_{Y \underline{x}_0} = c$ $H_1 : \mu_{Y \underline{x}_0} \neq c$ con $c \in \mathbb{R}$	$T_0 = \frac{\hat{Y}_0 - c}{ee(\hat{Y}_0)} \underset{\text{bajo } H_0}{\sim} t_{n-p}$	Se rechaza H_0 si $ T_0 > t_{1-\alpha/2, n-p}$

donde $t_{1-\alpha/2, n-p}$ es el percentil $1 - \alpha/2$ de la distribución t -Student con $n - p$ grados de libertad.

IC del $(1 - \alpha)100\%$ para la respuesta media $E[Y|\underline{x}_0]$:

Basados de nuevo en que el estadístico:

$$T = \frac{\hat{Y}_0 - E[Y|\underline{x}_0]}{ee(\hat{Y}_0)} \sim t_{n-p}$$

lo cual implica que:

$$P\left(-t_{\alpha/2, n-p} < \frac{\hat{Y}_0 - E[Y|\underline{x}_0]}{ee(\hat{Y}_0)} < t_{\alpha/2, n-p}\right) = 1 - \alpha$$

De donde se obtiene que un IC del $(1 - \alpha)100\%$ para la respuesta media:

$\mu_{Y|\underline{x}_0} = E[Y|\underline{x}_0]$ es:

$$\hat{y}_0 \pm t_{\alpha/2, n-p} ee(\hat{Y}_0).$$

- Considere ahora el problema de predecir un valor futuro Y_0 (no observado en la muestra) de la variable respuesta, en $X_1 = x_{01}, X_2 = x_{02}, \dots, X_k = x_{0k}$.
- Claramente, usando el modelo ajustado, predecimos de manera puntual tal valor por \hat{Y}_0 , pero sabemos que no es un estimador insesgado de Y_0 , por lo que siempre se genera un error de predicción dado por: $Y_0 - \hat{Y}_0$.
- Note que el error de predicción tiene media cero y dado que el valor futuro y su pronóstico son independientes, entonces la varianza del error de predicción $\hat{Y}_0 - Y_0$ está dada por:

$$\text{Var} [Y_0 - \hat{Y}_0] = \text{Var} [Y_0] + \text{Var} [\hat{Y}_0] = \sigma^2 [1 + \mathbf{x}_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0],$$

que es estimada por: $\widehat{\text{Var}} [\hat{Y}_0 - Y_0] = \text{MSE} [1 + \mathbf{x}_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0]$.

Con esto podemos hallar los siguientes resultados:

- Para un valor futuro Y_0 en un vector apropiado \underline{x}_0 :

IP del $(1 - \alpha)100\%$ para un valor futuro Y_0 :

Basados en este caso en que el estadístico:

$$T = \frac{Y_0 - \hat{Y}_0}{\text{ee}(Y_0 - \hat{Y}_0)} \sim t_{n-p},$$

con $\text{ee}(Y_0 - \hat{Y}_0) = \sqrt{\widehat{\text{Var}}[Y_0 - \hat{Y}_0]}$, lo cual implica que:

$$P \left(-t_{\alpha/2, n-p} < \frac{Y_0 - \hat{Y}_0}{\text{ee}(Y_0 - \hat{Y}_0)} < t_{\alpha/2, n-p} \right) = 1 - \alpha$$

De donde se obtiene que un IP del $(1 - \alpha)100\%$ para un valor futuro Y_0 es:

$$\hat{Y}_0 \pm t_{\alpha/2, n-p} \text{ee}(Y_0 - \hat{Y}_0)$$

Notas:

- Los intervalos de predicción estiman los posibles valores para un valor particular de la variable respuesta (no para su media) en un vector dado $\underline{\mathbf{x}}_0$.
- Asumimos que este valor particular es un valor futuro de la variable aleatoria Y , y por tanto, no fue utilizado en la regresión.
- Si Y_0 es un valor futuro y $\hat{Y}_0 = \underline{\mathbf{x}}_0 \hat{\underline{\beta}}$ es su estimador, entonces estas dos variables aleatorias son estadísticamente independientes, dado que Y_0 no fue utilizado para hallar los parámetros estimados, de ahí el estadístico y los límites del intervalo de predicción.

Precaución:

Deben evitarse las extrapolaciones por fuera del rango de experimentación en el espacio de las predictoras, para lo cual no basta con evaluar si cada valor componente del vector \underline{x}_0 se encuentra dentro del rango usado (u observado) para la correspondiente predictora, **sino que es necesario evaluar si \underline{x}_0 pertenece a la región de observación conjunta**.

Para ello basta con verificar si:

$$h_{00} = \underline{x}_0 (\mathbf{X}'\mathbf{X})^{-1} \underline{x}_0' < \max_{1 \leq i \leq n} \{h_{ii}\}$$

con h_{ij} el i -ésimo elemento de la matriz 'hat' $\mathbf{H}_{n \times n} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$.

Validación de los supuestos del modelo de RLM

Para la validación de supuestos se usan generalmente los residuales del modelo, los cuales sabemos que se definen así:

$$e_i = y_i - \hat{y}_i, \quad i = 1, \dots, n$$

Observe que, la magnitud de los residuales e_i depende de la escala de medida de la respuesta Y , lo cual no permite determinar cuando un residual es 'grande'. Para resolver este problema en lugar de usar los residuales crudos definidos arriba, se recomienda utilizar residuales escalados que transforman a los anteriores para tener media cero y varianza unitaria.

Residuales escalados

Se han definido varias versiones de residuales escalados, entre los que se destacan:

- **Residuales estandarizados:** para su definición se considera el supuesto sobre los errores, que establece que ε_i se distribuye con media cero y varianza σ^2 . Por tanto, los residuales estandarizados, denotados d_i se definen como:

$$d_i = \frac{e_i}{\sqrt{\text{MSE}}}, \quad i = 1, \dots, n$$

Si el supuesto es adecuado los valores de d_i deben oscilar entre -3 y 3. Por tanto, Un d_i grande ($|d_i| > 3$) es indicio de una observación potencialmente atípica.

- **Residuales estudentizados:** para su definición se considera el hecho de que realmente los residuales e_i en general no son independientes ni tienen varianza constante como los errores ε_i . Veamos, Sabemos que, $\underline{e} = (I - H)\underline{Y}$, donde $I - H$ es una matriz simétrica e idempotente. Luego,

$$\begin{aligned} E[\underline{e}] &= E[(I - H)\underline{Y}] = (I - H)E[\underline{Y}] = (I - H)\underline{X}\underline{\beta} \\ &= \underline{X}\underline{\beta} - H\underline{X}\underline{\beta} = \underline{X}\underline{\beta} - \underline{X}(\underline{X}'\underline{X})^{-1}\underline{X}'\underline{X}\underline{\beta} = \underline{0} \end{aligned}$$

$$Var[\underline{e}] = Var[(I - H)\underline{Y}] = (I - H)Var[\underline{Y}](I - H)' = \sigma^2(I - H)$$

De donde: $V(e_i) = \sigma^2(1 - h_{ii})$, y $cov(e_i, e_j) = -\sigma^2 h_{ij}$ para $i \neq j$.

Por tanto, mientras que los errores ε_i tienen varianza constante σ^2 y son incorrelacionados, los residuales no necesariamente tienen la misma varianza y pueden ser correlacionados.

De esta forma, los residuales estudentizados, denotados r_i , se definen como:

$$r_i = \frac{e_i}{\sqrt{\text{MSE}(1 - h_{ii})}}, \quad i = 1, \dots, n.$$

Este residual, también debe ubicarse entre -3 y 3. Se considera atípica aquella observación con un r_i grande ($|r_i| > 3$).

NOTAS:

- Si el modelo de RLM especificado es correcto los r_i tienen varianza aproximadamente constante!! igual a 1.
- En conjuntos grandes de datos la varianza de los r_i se puede estabilizar en 1 y así no habrá mucha diferencia entre éstos y los d_i .
- Si todos los supuestos del modelo se cumplen, se espera que aproximadamente el 68% de los residuales d_i ó r_i , estén entre -1 y $+1$, aproximadamente el 95% entre -2 y $+2$ y aproximadamente 99.7% entre -3 y $+3$.

La validación de los supuestos vista en regresión lineal simple se mantiene, solo que ahora se recomienda utilizar residuales escalados (d_i ó preferiblemente r_i) en lugar de utilizar los residuales crudos e_i .

Validación de los supuestos en los errores

Recuerde que en los modelos de regresión se han impuesto las siguientes cuatro condiciones sobre el término de error:

- Los errores son variables aleatorias normales.
- Los errores tienen media cero.
- Los errores tienen varianza constante.
- Los errores son mutuamente independientes.

Recuerde que en esta asignatura asumiremos el supuesto de independencia de los errores y en virtud de que usando los residuales del modelo el supuesto de media cero siempre se cumple, entonces se define lo siguiente:

- El supuesto de normalidad puede chequearse bien sea con el gráfico de probabilidad normal de los residuales o con alguna de las pruebas estadísticas de normalidad, entre las cuales se tienen las de Shapiro Wilk, Kolmogorov Smirnov, Cramér von Mises y Anderson Darling.
- Para chequear el supuesto de varianza constante, resulta útil un gráfico de residuales versus valores ajustados de la respuesta.

Medidas Remediales

Las medidas remediales descritas en el caso de RLS también son aplicables en RLM. Con el fin de superar las deficiencias del modelo se pueden realizar transformaciones sobre la variable respuesta y/o sobre las variables predictoras.

Las transformaciones sobre la respuesta pueden ayudar en el caso de que los errores no resulten normales o la varianza no sea constante. Transformaciones sobre las variables predictoras resultan útiles cuando la superficie de respuesta es curvilínea.

Si las desviaciones respecto al supuesto de normalidad son severas, y ninguna transformación resulta útil y/o interpretable, existe otra alternativa, **los llamados modelos lineales generalizados** con los cuales se pueden modelar respuestas que no se distribuyen normales; sin embargo, tales modelos están más allá del alcance de este curso.

Identificación de observaciones extremas en el modelo de RLM

Además de la validación de supuestos en los errores de un modelo de RLM, se debe chequear la presencia de observaciones extremas, tales como:

- Observaciones atípicas (outliers)
- Puntos de balanceo
- Observaciones influyentes

Observaciones atípicas

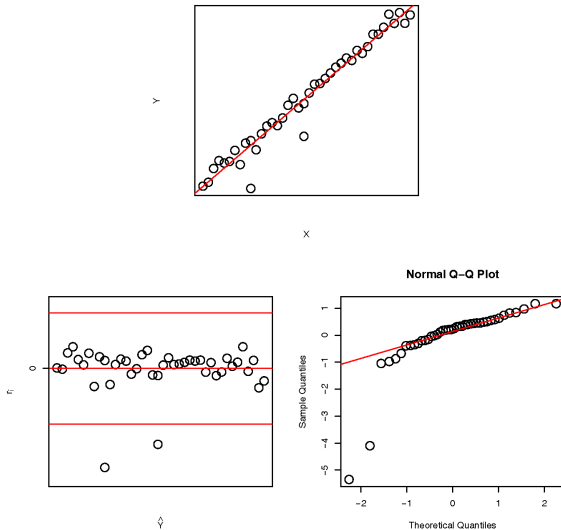
Una observación **atípica** (o **outlier**) es aquella que está separada (**en su valor de la respuesta Y**) del resto de las observaciones y por tanto puede afectar los resultados del ajuste del modelo de regresión.

Interesa identificarlas para luego, si es posible analizar si se tratan de observaciones malas (por errores de registro o medición) que pueden ser descartadas, o si realmente son datos correctos pero extraños que no deben ser eliminados del conjunto de datos.

Para detectar observaciones atípicas se usan los residuales escalados definidos anteriormente. Se considera que una observación es **atípica** **cuando su residual estudentizado r_i** , es tal que: $|r_i| > 3$.

Muchos **outliers** en los datos pueden causar niveles de confianza reales menores de lo esperado.

La siguiente Figura ilustra el caso de **dos observaciones atípicas**.



Puntos de balanceo

Un **punto de balanceo** es una observación en el espacio de las predictoras, alejada del resto de la muestra y que puede controlar ciertas propiedades del modelo ajustado.

Este tipo de observaciones posiblemente no afecte los coeficientes de regresión estimados pero sí las estadísticas de resumen como el R^2 y los errores estándar de los coeficientes estimados.

Los puntos de balanceo son detectados mediante el análisis de los elementos de la diagonal principal de la matriz H , los h_{ii} , que proporcionan una medida estandarizada de la distancia de la i -ésima observación al centro del espacio definido por las predictoras.

Se tiene lo siguiente:

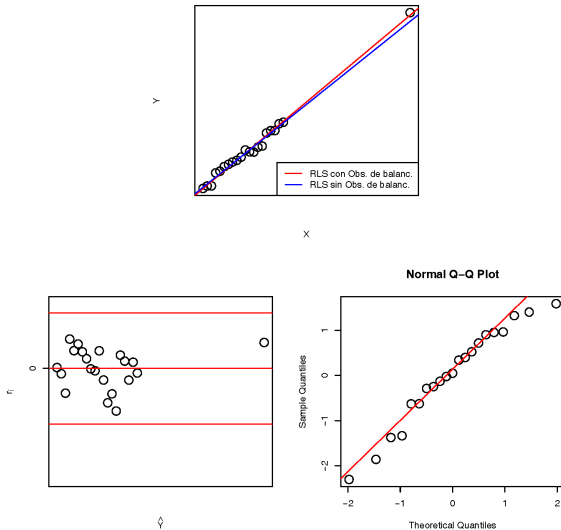
- La media de los h_{ii} es:

$$\bar{h} = \frac{\sum_{i=1}^n h_{ii}}{n} = \frac{\text{traza}(\mathbf{H})}{n} = \frac{p}{n}$$

con p el número de parámetros del modelo de RLM.

- Se asume que la observación i es un **punto de balanceo** si $h_{ii} > 2p/n$, pero si $2p/n > 1$ este criterio no funciona pues los h_{ii} siempre son menores que 1.

La próxima Figura ilustra el caso de **una observación de balanceo**.



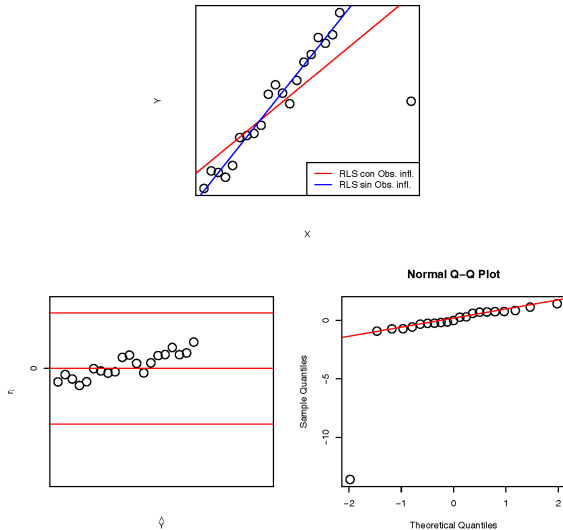
Observaciones influenciales

Una observación es **influyente** si tiene un impacto notable sobre los coeficientes de regresión ajustados, esto es, una observación influyente se dice que **hace al modelo en su dirección**, es decir, una observación es influyente si su exclusión del modelo causa cambios importantes en la ecuación de regresión ajustada.

Estas observaciones se caracterizan por tener un valor moderadamente inusual tanto en el espacio de las predictoras como en la respuesta.

Después de identificar las observaciones que están alejadas con respecto a los valores de Y (atípicas) y/o con respecto a sus valores en X (puntos de balanceo) evaluamos si éstas son influyentes.

La Figura siguiente ilustra el caso de **una observación influyente**.



Para la evaluación se cuenta con las siguientes medidas:

- Distancia de Cook.
- Diagnóstico DFFITS.
- Diagnóstico DFBETAS.

A continuación se presentan los diagnósticos para detectar observaciones influenciales.

- ❶ **Distancia de Cook:** es una medida de la distancia cuadrática entre, el estimador de $\underline{\beta}$ por mínimos cuadrados basado en las n observaciones, y el estimador de $\underline{\beta}$ obtenido eliminando la i -ésima observación, así:

$$D_i = \frac{(\hat{\underline{\beta}}_{(i)} - \hat{\underline{\beta}})' \mathbf{X}' \mathbf{X} (\hat{\underline{\beta}}_{(i)} - \hat{\underline{\beta}})}{p \text{MSE}} = \frac{r_i^2}{p} \left(\frac{h_{ii}}{1 - h_{ii}} \right), \quad i = 1, \dots, n$$

donde, $\hat{\underline{\beta}}_{(i)}$ es el vector de parámetros estimados obtenido cuando no se considera en el ajuste del modelo a la observación i .

Note que si D_i es alto, entonces la observación i tiene influencia sobre el vector de parámetros estimados $\hat{\underline{\beta}}$. Como $f_{0.5, p, n-p} \approx 1$ se dice que la observación i será influyente si $D_i > 1$

- ② **Diagnóstico DFFITS:** es el número de desviaciones estándar que el valor ajustado \hat{y}_i se mueve si la observación i es omitida:

$$\text{DFFITS}_i = \frac{\hat{Y}_i - \hat{Y}_{(i)}}{\sqrt{\text{MSE}_{(i)} h_{ii}}} = \frac{e_i}{\sqrt{\text{MSE}_{(i)} (1 - h_{ii})}} \left(\frac{h_{ii}}{1 - h_{ii}} \right)^{1/2}$$

donde, $\hat{Y}_{(i)}$ es el i -ésimo valor ajustado obtenido cuando no se considera en el ajuste del modelo a la observación i y $\text{MSE}_{(i)}$ es el cuadrado medio del error obtenido cuando no se considera en el ajuste del modelo a la observación i .

Una observación será **influencial** si $|\text{DFFITS}_i| > 2\sqrt{\frac{p}{n}}$.

- ③ **Diagnóstico DFBETAS:** indica cuánto cambia el j -ésimo coeficiente de regresión estimado $\hat{\beta}_j$ en unidades de desviación estándar, **si se omite la i -ésima observación:**

$$\text{DFBETAS}_{j(i)} = \frac{\hat{\beta}_j - \hat{\beta}_{j(i)}}{\sqrt{\text{MSE}_{(i)} c_{jj}}}$$

donde c_{jj} es el j -ésimo elemento en la diagonal principal de la matriz: $(\mathbf{X}'\mathbf{X})^{-1}$ y $\text{MSE}_{(i)}$ es el MSE de la regresión sin la observación i .

Una observación será influyente si $|\text{DFBETAS}_{j(i)}| > 2/\sqrt{n}$.

NOTA: Las observaciones que sean identificadas como puntos influenciales y/o de balanceo deben ser investigadas. Si las observaciones identificadas según los criterios anteriores son errores de registro en la base de datos, estas deben ser removidas del análisis. Caso contrario, se debe emplear métodos más robustos para mantener dichos registros en el conjunto de datos: **métodos estadísticos robustos que no se vean afectados por observaciones atípicas, de balanceo o influenciales.**