

Regresión lineal Multiple

Se desea estudiar la relación de una variable respuesta con 2 o más variables regresoras.

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_k X_{ik} + \epsilon_i$$

$\epsilon_i \sim N(0, \sigma^2)$
 $\forall i \quad i = 1, 2, 3, \dots, n$

Matriz de Varianzas y Covarianzas.

Esta Matriz contiene la covarianza entre variables, la cual indica su grado de variación conjunta.

Y_i : Variable respuesta de observaciones i .

X_{ik} : observación i de la variable regresora k .

$\beta_0, \beta_1, \beta_2, \dots, \beta_k$: Parámetros.

$$\Sigma(X) = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1n} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \dots & \sigma_{nn} \end{bmatrix}$$

k : \rightarrow regresoras (# de x)

P : $k+1 \rightarrow$ parámetros (# de β)

Características

Interpretaciones:

β_0 : es el valor medio de y cuando todas las variables X_k son iguales a cero. (Si no contiene al 0 no es interpretable)

$\beta_k \rightarrow$ es el efecto parcial sobre la media y por cada unidad de incremento de la variable X_k cuando las demás variables permanecen constantes.
Predictora

- Es una Matriz simétrica $\sigma_{ij} = \sigma_{ji}$

- Matriz cuadrada

- De orden $n \times n$

- Los elementos de la diagonal principal son las varianzas de las Variables $\sigma_{ii} = \sigma_i^2$

- Por Fuera de la diagonal están las covarianzas entre los pares de elementos del vector aleatorio.

Forma matricial del Modelo

$$\begin{array}{c} \underline{y} \\ \left[\begin{array}{c} y_1 \\ y_2 \\ \vdots \\ y_n \end{array} \right] \\ \downarrow \\ n \times 1 \end{array} = \begin{array}{c} \underline{X} \\ \left[\begin{array}{cccc} 1 & X_{11} & \dots & X_{1k} \\ 1 & X_{21} & \dots & X_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & \dots & X_{nk} \end{array} \right] \\ \downarrow \\ n \times p \end{array} \begin{array}{c} \underline{\beta} \\ \left[\begin{array}{c} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{array} \right] \\ \downarrow \\ p \times 1 \end{array} + \begin{array}{c} \underline{\varepsilon} \\ \left[\begin{array}{c} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{array} \right] \\ \downarrow \\ n \times 1 \end{array}$$

$$\underline{y} = \underline{X} \underline{\beta} + \underline{\varepsilon}$$

matriz de
diseño

$$\underline{\varepsilon} \sim N(\underline{0}_{n \times 1}, \sigma^2 \underline{I}_{n \times n})$$

Estimación Por Mínimos cuadrados de los Parámetros del Modelo

Se buscan valores estimados de los parámetros tales que se minimice la suma de cuadrados del error:

$$S(\underline{\beta}) = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 X_{i1} - \beta_2 X_{i2} - \dots - \beta_k X_{ik})^2$$

$$\hat{\underline{\beta}} = (\underline{X}^t \underline{X})^{-1} \underline{X}^t \underline{y}$$

Significancia de los parámetros β_j

promedio β_j unidades, siempre que las demás predictoras permanezcan constantes.

Hipotesis:

$$H_0: \beta_j = 0 \vee H_a: \beta_j \neq 0$$

con $j = 0, 1, 2, \dots, K$

ESTIMACIÓN de σ^2

Bajo los supuestos $\varepsilon_i \sim N(0, \sigma^2)$

Estadístico de Prueba:

$$t_0 = \frac{\hat{\beta}_j}{S_e(\hat{\beta}_j)} \sim \text{Bajo } H_0 \quad t_{n-p}$$

el estimador insesgado para la Varianza es:

$$\hat{\sigma}^2 = MSE = \frac{SSE}{n-p}$$

Decisión:

~~Se rechaza~~ H_0 si $|t_0| > t_{\frac{\alpha}{2}, n-p}$

donde:

~~Se rechaza~~ H_0 si el valor p

$$p = P(t_{\frac{\alpha}{2}, n-p} > |t_0|)$$

es pequeño

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Interpretación:

→ Para que β_0 sea interpretable

debe de ser significativo y en el conjunto de observaciones debe

estar $(x_1, x_2, \dots, x_k) = (0, 0, \dots, 0)$

→ Para interpretar los demás β_j

Solo es necesario que sean Significativos.

Interpretan como:

- β_j : un aumento unitario en x_j se

espera que la variable respuesta (Y)

aumente o disminuya en

Análisis de la Varianza

ANOVA

Fuente de Variación	Suma de Cuadrados	Grados de libertad	Cuadrados medios	F_0
Regresión	SSR	k	$MSR = \frac{SSR}{k}$	$F_0 = \frac{MSR}{MSE}$
Error	SSE	n-p	$MSE = \frac{SSE}{n-p}$	$= \frac{SSR (n-p)}{SSE (k)}$
Total	SST	n-1		

Prueba de significancia de la regresión

Hipotesis

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0 \text{ vs } H_a: \text{Al menos } \beta_j \neq 0$$

$$j = 1, 2, 3, \dots, k$$

Estadístico de Prueba:

$$F_0 = \frac{MSR}{MSE} \quad \text{Bajo } H_0 \quad F_{k, n-p}$$

Decisión:

Se rechaza H_0 si $F_0 > F_{\alpha, k, n-p}$

Se rechaza H_0 si el valor $p = p(F_{\alpha, k, n-p} > F_0)$ es pequeño

Coefficientes de determinación:

$$R^2 = \frac{SSR}{SST} = \frac{SSR}{SSR + SSE}$$

Al ingresar la variable tiende a no decrecer, aún cuando existan dentro del grupo de variables, un subconjunto de ellas que no aportan significativamente

$$R^2_{adj} = 1 - \frac{MSE}{MST} = 1 - \frac{SSE(n-1)}{SST(n-p)}$$

Disminuye cuando al modelo de regresión ingresan variables que no logran reducir el SSE.

Penaliza al modelo al ingreso de nuevas variables.

Intervalo de confianza para los Parámetros β_j :

$$\hat{\beta}_j \pm t_{\frac{\alpha}{2}, n-p} \cdot se(\hat{\beta}_j)$$

pruebas sobre subconjuntos de Coeficientes.

se quiere probar si un subconjunto de β es significativo simultáneamente:

Hipotesis:

H_0 : El subconjunto de $\beta = 0$
vs

H_a : Al menos un β del grupo $\neq 0$

Planteamos 2 Modelos

$$\varepsilon_i \sim \text{iid } N(0, \sigma^2)$$

MC (Modelo Completo) = $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + \varepsilon_i$
MF

MR (Modelo Reducido) = Reducido Bajo H_0

Estadístico de Prueba:

$$F_0 = \frac{\text{SSE (MR)} - \text{SSE (MC)}}{\text{MSE (M.C)}} \sim F_{r, n-p}$$

Decisión:

Se rechaza H_0 si $F_0 > F_{\alpha, r, n-p}$

Se rechaza si el valor $p = P(F_{\alpha, r, n-p} > F_0)$ es pequeño.

HIPOTESIS LINEAL General.

Se quiere ver por ejemplo si el efecto de dos o más variables es igual, o si el efecto de una variable es c veces otra.

Se pueden probar varias hipótesis juntas.

Ejemplo: 5 variables regresoras:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \beta_4 x_{i4} + \beta_5 x_{i5} + \varepsilon_i$$

$\varepsilon_i \sim \text{iid } N(0, \sigma^2)$

Probar las siguientes hipótesis:

$$H_0: \beta_1 = \beta_2, \beta_3 = \beta_4, \beta_5 = 0 \quad \text{vs} \quad H_a: \beta_1 \neq \beta_2, \beta_3 \neq \beta_4 \text{ o } \beta_5 \neq 0$$

igualamos las ecuaciones a cero y tenemos:

$$H_0: \beta_1 - \beta_2 = 0, \beta_3 - \beta_4 = 0, \beta_5 = 0 \quad \text{vs} \quad H_a: \beta_1 - \beta_2 \neq 0 \text{ o } \beta_3 - \beta_4 \neq 0 \text{ o } \beta_5 \neq 0$$

Estamos probando 3 ecuaciones diferentes, entonces armamos la Matriz L de orden $m \times p$ en este caso 3×6

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

Planteamos el Modelo reducido bajo la hipótesis nula:

$$M.R: Y_i = \beta_0 + \beta_1 X_{i1} + \beta_1 X_{i2} + \beta_3 X_{i3} + \beta_3 X_{i4} + \epsilon_i$$

$$Y_i = \beta_0 + \beta_1 (X_{i1} + X_{i2}) + \beta_3 (X_{i3} + X_{i4}) + \epsilon_i$$

$$Y_i = \beta_0 + \beta_1 Z_{i1,2} + \beta_3 Z_{i3,4} + \epsilon_i$$

Estadístico de Prueba:

$$F_0 = \frac{SSE(M.R) - SSE(M.F)/m}{MSE(M.C)} \sim F_{m, n-p}$$

$m = \#$ de ecuaciones que se están probando o $\#$ de Filas L.I de matriz L

Decisión:

Se rechaza H_0 si $F_0 > F_{\alpha, m, n-p}$

Se rechaza H_0 el valor $p = P(F_{\alpha, m, n-p} > F_0)$ es pequeño

Normalidad:

si $|r_i| > 3$ i es una observación atípica.

Hipotesis:

$H_0: \varepsilon_i \sim \text{Normal}$ vs $H_a: \varepsilon_i \neq \text{Normal}$

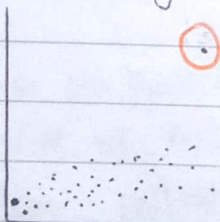
Prueba gráfica: se evalúa si los puntos en la escala normal se pueden ajustar por una línea recta.

Test de Shapiro-Wilk

Es una prueba específica de normalidad que arroja un valor p .

Observaciones atípicas

Es un valor atípico en Y que está separado del resto de datos y podría afectar los resultados del ajuste del modelo



Puntos de balanceo

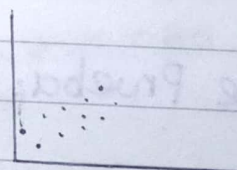
un valor en las X alejado y puede controlar algunas propiedades del modelo como el R^2

sea $H_{ii} = X_i^T (X^T X)^{-1} X_i$ es el elemento i de la Matriz hat (H) entonces:

si $h_{ii} > \frac{2p}{n}$ entonces i es un punto de balanceo.

Puntos Influyentes:

Es un valor inusual tanto en X como en Y suelen balar el modelo en su dirección.



Residuales estandarizados d_i

$$d_i = \frac{e_i}{\sqrt{MSE}}$$

si $|d_i| > 3$ indicios de atípica

Residuales estandarizados r_i

$$r_i = \frac{d_i}{\sqrt{1 - h_{ii}}}$$

Distancia de Cook:

si $COOK's d_i > 1$ entonces i es un punto influyente.

DFBETS:

si $|DFBETS_i| > \sqrt{\frac{2p}{n}}$ entonces i es un punto influyente.

Si $h_{00} < h_{\max}$ es una interpolación

Si $h_{00} > h_{\max}$ es una extrapolación

h_{\max} es el valor máximo de h_{ij}

que no sea un punto de balanceo.

Intervalo de Confianza Para la respuesta Media.

$$\hat{Y}_0 \pm t_{\frac{\alpha}{2}, n-p} \cdot Se(\hat{Y}_0)$$

Intervalo de Predicción Para Futuras observaciones

$$\hat{Y}_0 \pm t_{\frac{\alpha}{2}, n-p} \cdot \sqrt{MSE + Se(\hat{Y}_0)^2}$$

Coefficientes de regresión estandarizados.

Si las X no están en la misma escala de Medición no podemos determinar cual tiene Mayor o Menor efecto, para esto los estandarizamos así:

$$y_i^* = \frac{y_i - \bar{Y}}{\sqrt{\frac{\sum_{i=1}^n (y_i - \bar{Y})^2}{n-1}}}$$

$$x_{ij}^* = \frac{x_{ij} - \bar{x}_j}{\sqrt{\frac{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}{n-1}}}$$

$$j = 1, 2, 3 \dots k$$

$$\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

modelo sin intercepto: $\beta_1 x_{i1}^* + \beta_2 x_{i2}^* + \beta_3 x_{i3}^* + \dots + \beta_k x_{ik}^* + \varepsilon_i$