

# Análisis de Regresión Lineal Múltiple

Juan Sosa, PhD



I - 2018

## Objetivo

- Caracterizar o explicar el valor medio de una variable [dependiente] bajo condiciones específicas de otra(s) variable(s) [independientes].

## Observaciones

- Establecer **asociación**, mas **no causalidad**.
- Aunque predecir [el valor promedio] no es el centro de atención, pero es posible hacerlo.
- Modelo [de carácter probabilístico] con muchas posibilidades!

## Estructura del conjunto de datos

- Matriz de datos  $\mathbf{X}$  de tamaño  $n \times p$  : variables independientes/explicativas.
- Vector de datos  $\mathbf{y}$  de tamaño  $n \times 1$  : variable dependiente/explicada/respuesta.

## Formulación

$$y_1 = \beta_1 x_{1,1} + \beta_2 x_{1,2} + \dots + \beta_p x_{1,p} + \epsilon_1$$

$$y_2 = \beta_1 x_{2,1} + \beta_2 x_{2,2} + \dots + \beta_p x_{2,p} + \epsilon_2$$

$$\vdots$$

$$y_n = \beta_1 x_{n,1} + \beta_2 x_{n,2} + \dots + \beta_p x_{n,p} + \epsilon_n$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,p} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,p} \\ \vdots & \vdots & \vdots & \vdots \\ x_{n,1} & x_{n,2} & \cdots & x_{n,p} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

## Terminología

- $x$  : variables independientes/explicativas.
- $y$  : variable dependiente/explicada/respuesta.
- $\beta$  : coeficientes de regresión.
- $\epsilon$  : errores/perturbaciones aleatorias.

## Supuestos

$\epsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n)$  donde  $\mathbf{I}_n$  es la matriz identidad de tamaño  $n$ .

## Observaciones

- Los parámetros del modelo son  $\beta_1, \beta_2, \dots, \beta_p$  y  $\sigma^2$ .
- Las variables explicativas no tienen carácter aleatorio y pueden ser continuas o categóricas.
- El supuesto de normalidad y varianza constante no es una “camisa de fuerza”.
- El carácter aleatorio del modelo recae en la variable dependiente a través del error.
- Los errores se asumen como independientes.

Como  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$  se obtiene que  $\mathbf{y} \mid \mathbf{X} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$  y además

$$\begin{aligned}\mathbb{E}[y_i \mid \mathbf{x}_i] &= \beta_1 x_{i,1} + \beta_2 x_{i,2} + \dots + \beta_p x_{i,p} = \boldsymbol{\beta}^T \mathbf{x}_i \\ \text{Var}[y_i \mid \mathbf{x}_i] &= \sigma^2\end{aligned}$$

## Interpretación de $\beta_i$

Cambio en  $y$  correspondiente a una unidad de cambio en  $x_i$  cuando el resto de las variables explicativas se mantienen constantes.

## Consecuencias

- El modelo caracteriza el valor medio de  $y$  para valores dados de  $x$ .
- La linealidad se encuentra en los coeficientes de regresión.
- La interpretación de los coeficientes no está dada en términos de cargas/pesos.

## Estimación

Encontrar los valores de  $\beta$  y  $\sigma^2$  para reproducir  $\mathbf{y}$  tan precisamente como sea posible.

### Método de máxima verosimilitud

Encontrar los valores de  $\beta$  y  $\sigma^2$  que maximicen (optimicen) la función de verosimilitud:

$$\prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2} (y_i - \beta^T \mathbf{x}_i)^2\right\}$$

### Método de mínimos cuadrados ordinarios

Encontrar los valores de  $\beta$  que minimicen (optimicen) la función cuadrática:

$$\sum_{i=1}^n (y_i - \beta^T \mathbf{x}_i)^2$$

### Estimador de $\beta$

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$
$$\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1})$$

### Valores ajustados/predichos

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}$$

Estimador del valor medio de  $\mathbf{y}$ .

### Residuales

$$\mathbf{r} = \mathbf{y} - \hat{\mathbf{y}} = \mathbf{I}_n \mathbf{y} - \mathbf{H}\mathbf{y} = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$$

Permiten evaluar la **bondad de ajuste** y los **supuestos** del modelo.

### Sumas de cuadrados

La **variabilidad total de la variable respuesta** se descompone en la **variabilidad debida a los valores ajustados** y la **variabilidad debida a los residuales**:

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2$$
$$SCT = SCR + SCE$$

Observación: **SCE** es la cantidad que se debe minimizar para ajustar el modelo.

### Coeficiente de determinación

$$R^2 = \frac{SCR}{SCT} = 1 - \frac{SCE}{SCT}$$

Proporción de la variabilidad total que explica el modelo ajustado.

### Coeficiente de determinación ajustado

$$R_a^2 = 1 - \frac{n-1}{n-p}(1 - R^2)$$

### Estimador insesgado de $\sigma^2$

$$\hat{\sigma}^2 = \frac{SCE}{n-p} = CME$$

$CME$  : cuadrado medio del error.

### Observación

Para apreciar el valor agrado de una regresión se recomienda comparar  $s_y$  con  $\hat{\sigma}$ .



## Significancia global

Existe una relación significativa entre la variable respuesta y las variables explicativas en general?

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0 \quad \text{frente a} \quad H_1 : \beta_j \neq 0 \text{ para algún } j$$

$$\text{Estadístico de prueba: } F = \frac{SCR/(p-1)}{SCE/(n-p)} = \frac{CMR}{CME} \sim F_{p-1, n-p}$$

## Significancia particular

Existe una relación significativa entre la variable respuesta y una variable explicativa en particular?

$$H_0 : \beta_i = 0 \quad \text{frente a} \quad H_1 : \beta_i \neq 0$$

$$\text{Estadístico de prueba: } t = \frac{\hat{\beta}_i}{s_{\hat{\beta}_i}} \sim t_{n-p}$$

$$\text{Intervalo de confianza: } \hat{\beta} \pm t_{n-p} s_{\hat{\beta}_i}$$

### Observación promedio para $x_0$

$$\hat{\beta}^T x_0 \pm t_{n-p} \hat{\sigma} \sqrt{x_0^T (\mathbf{X}^T \mathbf{X})^{-1} x_0}$$

### Observación particular para $x_0$

$$\hat{\beta}^T x_0 \pm t_{n-p} \hat{\sigma} \sqrt{1 + x_0^T (\mathbf{X}^T \mathbf{X})^{-1} x_0}$$

### Observación

La incertidumbre es mayor cuando se predice una observación particular que al predecir la respuesta media.

Correlación entre las variables explicativas, i.e, información redundante en el modelo.

## Observaciones

- Sobre carga de información.
- Una correlación absoluta mayor a 0.7 sugiere problemas de este estilo.
- Hay un aumento del error estándar de los coeficientes de regresión.
- Se hace difícil determinar la significancia de las variables explicativas en particular.

## Supuestos

- $\mathbb{E}[\epsilon_i] = 0$ .
- $\text{Var}[\epsilon_i] = \sigma^2$ .
- La distribución de los errores es normal.
- Los errores son mutuamente independientes.

## Residuales estandarizados

$$r_i^* = \frac{r_i}{s_{r_i}} = \frac{r_i}{\hat{\sigma}\sqrt{1-h_{ii}}}$$

## Recomendaciones

- Histograma y gráfico cuantil-cuantil (*qqplot*) de los residuales (normalidad).
- Graficar residuales (estandarizados) frente a valores predichos (varianza constante).
- Graficar residuales (estandarizados) frente a cada variable explicada (tendencias).

## Definición

Observación que diverge la tendencia general de la data. Hacen crecer el CME y por lo tanto los márgenes de error de los intervalos de confianza.

## Residuales estudentizados

$$\tilde{r}_i = \frac{y_i - \hat{y}_i}{\hat{\sigma}_{(i)} \sqrt{1 - h_{ii}}}$$

## Detección

Una observación se considerada (potencialmente) como outlier si:

$$|r_i^*| > 2 \quad \text{o} \quad |\tilde{r}_i| > 3$$

## Definición

Observación que cambia sustancialmente el análisis cuando hace parte de la data.

## DFFITS (*difference in fits*)

$$DFFITS_i = \frac{\hat{y}_i - \hat{y}_{(i)}}{\hat{\sigma}_{(i)} \sqrt{h_{ii}}}$$

## Distancia de Cook

$$D_i = \frac{(y_i - \hat{y}_i)^2}{p \hat{\sigma}^2} \frac{h_{ii}}{(1 - h_{ii})^2}$$

## Detección

Una observación se considerada (potencialmente) como influyente si:

$$DFFITS_i > 2\sqrt{\frac{p}{n-p}} \quad \text{o} \quad D_i > 1$$

## Selección de variables

- *Forward regression.*
- *Backward regression.*
- *Stepwise regression.*

## Comparación de modelos

- Criterio de información de Akaike (AIC).
- Criterio de información Bayesiano (BIC)
- Pruebas de hipótesis para modelos anidados.

## Todo tipo de variables explicativas

- Variables dummy.
- Variables polinómicas.

- Transformar la variable respuesta (logaritmo, raíz cuadrada, Box-Cox).
- Regresión lineal general.
- Regresión rígida.
- Regresión Lasso.
- Regresión ElasticNet.
- Regresión no lineal.
- Regresión no paramétrica.
- Regresión cuantílica.
- Regresión robusta.
- Modelos lineales generalizados.