

Regresión lineal múltiple.

Nota. Los modelos de regresión funcionan tanto para datos numéricos como para datos aleatorios

Modelo de regresión simple

$$y = \underset{\substack{\uparrow \\ \text{variable dependiente}}}{\beta_0} + \underset{\substack{\uparrow \\ \text{variable independiente}}}{\beta_1} X$$

intercepto coef. de regresión
slope (pendiente)

Siempre es posible realizar una regresión lineal simple.

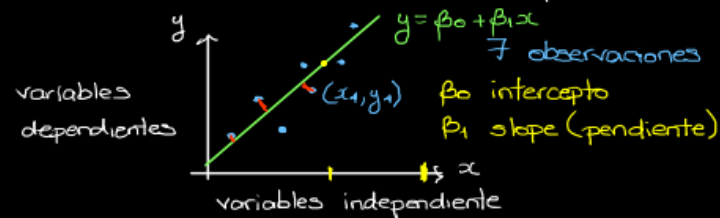
cuantitativos

El caso simple:

1.- ¿Qué tan fuerte es la relación entre dos variables?

2.- Conocer el valor de la variable dependiente de

un cierto valor de la variable independiente.



Hipótesis:

(i) Homogeneidad de la varianza: el tamaño del error en la predicción no cambia (significativamente) conforme los valores de la variable indep cambia

(ii) Independencia entre las observaciones

(iii) Normalidad Los datos siguen una distribución normal

variable dependiente variable independiente

Error total del modelo

$$e = y - (\beta_0 + \beta_1 x)$$

Hipótesis sobre el error: $E[e] = 0$
 $Var[e] = 1$

Meta: A partir de un conjunto dado de datos $\{(x_i, y_i)\}_{i=1}^n$

encontrar

$$\hat{\beta}_0, \hat{\beta}_1 = \arg \min_{\beta_0, \beta_1} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

$\hat{\beta}_0$ y $\hat{\beta}_1$ son estimadores para β_0 y β_1 respectivamente

- $\hat{\beta}_0$ y $\hat{\beta}_1$ dependen de $\{(x_i, y_i)\}$
- "Aproximan" a β_0 y a β_1

(X, Y) vector gaussiano:

$$aX + bY \sim N(\mu, \sigma^2)$$

para cualquier a, b que no sean 0 simultáneamente.

Por un resultado en prob., existe una v.a. que se llama esperanza condicional

$$E[Y|X] = \beta_0 + \beta_1 X \quad \beta_0 = E[Y - \beta_1 X]$$

Muestral

$$\beta_1 = \frac{Cov(X, Y)}{Var(X)}$$

$$\beta_0 = \bar{Y} - \beta_1 \bar{X}$$

$$\beta_1 = \frac{Cov(\hat{X}, \hat{Y})}{\hat{Var}(X)}$$

Múltiple:

$k < n$

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} \quad i=1, \dots, n \quad (1)$$

El sistema de ecuaciones lineales (1), lo podemos reescribir como

$$\underbrace{y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}}_{n \times 1} = \underbrace{X = \begin{pmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{pmatrix}}_{n \times (k+1)} \underbrace{\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}}_{(k+1) \times 1} + \underbrace{e = \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix}}_{n \times 1}$$

Obs. Los coeficientes de la matriz X pueden ser variables aleatorias pero si son variables numéricas el vector e de errores se aproxima (en promedio) a 0

Problema de optimización:

Este vector es el que necesitamos

$$(x_1, \dots, x_n) = x$$

$$S(\beta) = \sum_{i=1}^n e_i^2 = e^T e = (y - X\beta)^T (y - X\beta)$$

$$(A+B)^T = A^T + B^T$$

$$(AB)^T = B^T A^T$$

$$S(\beta) = y^T y - \beta^T X^T y - y^T X \beta + \beta^T X^T X \beta$$

$$S(\beta) = y^T y - 2\beta^T X^T y + \beta^T X^T X \beta$$

$$\frac{\partial S(\beta)}{\partial \beta} = -2X^T y + 2X^T X \beta = 0$$

$$\frac{d\beta^T \beta}{d\beta} = 2\beta$$

$$\Rightarrow \boxed{X^T X \beta = X^T y} \quad \text{Ecuaciones normales}$$

$$\text{Si } X^T X \text{ es invertible} \Rightarrow \boxed{\hat{\beta} = (X^T X)^{-1} X^T y}$$

Coefficiente de determinación

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

Medido de qué tan bien se ajusta el modelo de regresión

$$SS_{\text{res}} = \sum e_i^2$$

$R^2 \sim 1$ mejor ajuste

$$SS_{\text{tot}} = \sum (y_i - \bar{y})^2$$

numerador que aparece en la varianza

