

TEMA 3

MODELO DE REGRESIÓN LINEAL

| X | Y |
|----------|----------|
| x_1 | y_1 |
| x_2 | y_2 |
| \vdots | \vdots |
| x_n | y_n |

| X_1 | X_2 | \dots | X_k | Y |
|----------|----------|---------|----------|----------|
| x_{11} | x_{21} | \dots | x_{k1} | y_1 |
| x_{12} | x_{22} | \dots | x_{k2} | y_2 |
| \vdots | \vdots | | \vdots | \vdots |
| x_{1n} | x_{2n} | \dots | x_{kn} | y_n |

\uparrow var. explicativa (regresión)
 \uparrow respuesta var. dep (papeles no simétricos)

$\underbrace{\hspace{10em}}$ explicativas
 \uparrow respuesta

Recordatorio (Recta de regresión) : $x_1 y_1, \dots, x_n y_n$

Calcular β_0 y β_1 para las que se minimizan los errores verticales entre la nube de puntos y la recta $y = \beta_0 + \beta_1 x$

$$ECM(\beta_0, \beta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

$$\rightarrow \begin{cases} \hat{\beta}_1 = \frac{cov_{xy}}{V_x} \\ \hat{\beta}_0 = \bar{y} - \frac{cov_{xy}}{V_x} \cdot \bar{x} \end{cases}$$

$$\text{Recta regresión: } y - \bar{y} = \frac{cov_{xy}}{V_x} (x - \bar{x}) \iff y = \hat{\beta}_0 + \hat{\beta}_1 x$$

$$\text{Errores/residuos: } e_i = y_i - \underbrace{(\hat{\beta}_0 + \hat{\beta}_1 x_i)}_{\hat{y}_i} \quad i = 1, \dots, n$$

$$ECM(\hat{\beta}_0, \hat{\beta}_1) = \frac{1}{n} \sum_{i=1}^n e_i^2$$

1. Modelo de regresión lineal simple

$$\begin{array}{c} X \\ x_1 \\ | \\ x_n \end{array} \quad \begin{array}{c} Y \\ y_1 \\ | \\ y_n \end{array}$$

$$Y|X=x := \beta_0 + \beta_1 x + \varepsilon$$

$$\mathbb{E}(\varepsilon) = 0$$

$$\mathbb{E}(\varepsilon) = \sigma^2$$

$$\text{Habitual: } \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

Recuerdo: (X, Y) normal bidim. $\begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{yx} & \sigma_y^2 \end{pmatrix}$

$$Y|X=a \sim \mathcal{N}(\tilde{\mu}, \tilde{\sigma}^2)$$

$$\begin{cases} \tilde{\mu} = \mu_y + \sigma_{xy} \frac{1}{\sigma_x^2} (a - \mu_x) & \leftarrow \text{lineal en } a \\ \tilde{\sigma}^2 = \sigma_y^2 - \frac{\sigma_{xy}^2}{\sigma_x^2} & \leftarrow \text{no depende de } a \end{cases}$$

Modelo:

$$\begin{array}{c} x_1 \ y_1 \\ | \quad | \\ x_n \ y_n \\ \hline \text{fijos} \end{array}$$

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \quad i=1, \dots, n$$

↓ variables aleatorias

Condiciones:

$$\rightarrow \mathbb{E}(\varepsilon_i) = 0 \quad \forall i$$

$$\rightarrow V(\varepsilon_i) = \sigma^2 \quad \forall i$$

$$\rightarrow \text{cov}(\varepsilon_i, \varepsilon_j) = 0 \quad i \neq j$$

$$\text{Muy habitual: } \varepsilon_i \sim \mathcal{N}(0, \sigma^2) \text{ indep.}$$

Versión matricial

$$\begin{pmatrix} y_1 \\ | \\ y_n \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & x_1 \\ | & | \\ 1 & x_n \end{pmatrix}}_X \underbrace{\begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}}_{\vec{\beta}} + \underbrace{\begin{pmatrix} \varepsilon_1 \\ | \\ \varepsilon_n \end{pmatrix}}_{\vec{\varepsilon}}$$

$$Y = X \vec{\beta} + \vec{\varepsilon}$$

$$\mathbb{E}(\vec{\varepsilon}) = \vec{0} \quad \text{cov}(\vec{\varepsilon}) = \sigma^2 I_n$$

$$\Rightarrow \mathbb{E}(Y) = X \cdot \vec{\beta}$$

$$\text{cov}(Y) = \sigma^2 I_n$$

$$\text{Habitual: } \vec{\varepsilon} \sim \mathcal{N}(\vec{0}, \sigma^2 I_n) \Rightarrow Y \sim \mathcal{N}(X \cdot \vec{\beta}, \sigma^2 I_n)$$

$$\text{Parámetros: } \beta_0, \beta_1, \sigma^2$$

1.1. - ESTIMACIÓN para $\vec{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}$ (Mínimos cuadrados)

Buscamos $\vec{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}$ para el que $f(\vec{\beta}) = (\vec{y} - \vec{X}\vec{\beta})^T (\vec{y} - \vec{X}\vec{\beta}) =$

$$= \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

$$f(\vec{\beta}) = \vec{y}^T \vec{y} - \vec{y}^T \vec{X} \vec{\beta} - \vec{\beta}^T \vec{X}^T \vec{y} + \vec{\beta}^T \vec{X}^T \vec{X} \vec{\beta} =$$

$$= \vec{y}^T \vec{y} - 2\vec{\beta}^T \vec{X}^T \vec{y} + \vec{\beta}^T \vec{X}^T \vec{X} \vec{\beta}$$

$$\frac{\partial f}{\partial \vec{\beta}} = -2\vec{X}^T \vec{y} + 2\vec{X}^T \vec{X} \vec{\beta} = \vec{0} \Leftrightarrow \vec{X}^T \vec{X} \vec{\beta} = \vec{X}^T \vec{y} \Leftrightarrow$$

$$\Leftrightarrow \boxed{\hat{\vec{\beta}} = (\vec{X}^T \vec{X})^{-1} \vec{X}^T \cdot \vec{y}}$$

Obs: $\vec{X}^T \vec{X} = \begin{pmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \end{pmatrix} \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} = \begin{pmatrix} n & n\bar{x} \\ n\bar{x} & n\bar{x}^2 \end{pmatrix} = n \begin{pmatrix} 1 & \bar{x} \\ \bar{x} & \bar{x}^2 \end{pmatrix}$

$$\det(\vec{X}^T \vec{X}) \Leftrightarrow V_x \neq 0$$

$$(\vec{X}^T \vec{X})^{-1} = \frac{1}{n V_x} \begin{pmatrix} \bar{x}^2 & -\bar{x} \\ -\bar{x} & 1 \end{pmatrix}$$

$$\begin{cases} \hat{\beta}_1 = \frac{1}{n V_x} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \frac{\text{cov}_{xy}}{V_x} \\ \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \end{cases}$$

1.2. - VALORES AJUSTADOS, RESIDUOS y SUMAS DE CUADRADOS

Ya tenemos $\hat{\vec{\beta}} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix}$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \quad i=1, \dots, n$$

$$\hat{\vec{y}} = \begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{pmatrix} \longrightarrow \hat{\vec{y}} = \vec{X} \hat{\vec{\beta}} = \underbrace{\vec{X} (\vec{X}^T \vec{X})^{-1} \vec{X}^T}_{\text{"H matriz hat"}} \vec{y} = H \vec{y}$$

H es $\begin{cases} \text{simétrica} & H^T = H \\ \text{idempotente} & H^2 = H \end{cases}$

Residuos: $\begin{pmatrix} e_1 \\ \vdots \\ e_n \end{pmatrix} \rightarrow e = y - \hat{y} = \underbrace{(I_n - H)}_{\text{simétrica, idempotente}} y$ $e_i = y_i - \hat{y}_i \quad i=1, \dots, n$

Obs: 1) $x^T \cdot e = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$, porque $x^T (I_n - H) y = \underbrace{(x^T - x^T x (x^T x)^{-1} x^T)}_{=0} y = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$

También sabemos $\begin{pmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \end{pmatrix} \cdot e = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \Rightarrow$

$$\Rightarrow \begin{cases} 1^T \cdot e = 0 \Rightarrow \sum_{i=1}^n e_i = 0 \\ x^T e = 0 \Rightarrow \sum_{i=1}^n x_i e_i = 0 \end{cases}$$

2) $\hat{y}^T \cdot e = 0$, porque $[H \cdot (I_n - H)] y = H - H^2 = H - H = 0$

¿Cuán grandes son los residuos?

$$e^T \cdot e = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2 = y^T (I_n - H) y$$

Observar que $\sum e_i^2$ dividido por \underbrace{n}_{n-2} , será estimación de σ^2

Observación sobre sumas de cuadrados:

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n ((y_i - \hat{y}_i) + (\hat{y}_i - \bar{y}))^2 + \underbrace{2 \sum_{i=1}^n e_i (\hat{y}_i - \bar{y})}_{=0}$$

(salvo n , es V_y variación total de las y_i)

$$\underbrace{\sum_{i=1}^n (y_i - \bar{y})^2}_{\substack{\text{suma de} \\ \text{cuadrados} \\ \text{total} \\ \text{(TSS - SCT)}}} = \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}_{\substack{\text{suma de cuadrados} \\ \text{explicada por} \\ \text{modelo} \\ \text{(MSS - SCM)}}} + \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}_{\substack{\text{suma de} \\ \text{cuadrados} \\ \text{residual} \\ \text{(RSS - SCR)}}}$$

Se define entonces el coeficiente de determinación:

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

Obs: para el caso bidimensional
 $R^2 = \rho_{xy}^2$

1.3 - PROPIEDADES DE LOS ESTIMADORES

no

$$\begin{matrix} x_1 & y_1 \\ | & | \\ & \\ x_n & y_n \end{matrix}$$

sino

$$\begin{matrix} x_1 & Y_1 \\ | & | \\ & \\ x_n & Y_n \end{matrix} = Y$$

$$Y = X \cdot \beta + \epsilon = \begin{pmatrix} \epsilon_1 \\ | \\ \epsilon_n \end{pmatrix}$$

$\begin{cases} E(\epsilon) = 0 \\ \text{cov}(\epsilon) = \sigma^2 I_n \end{cases}$
 quizás en algún momento $\epsilon \sim N(0, \sigma^2 I)$

Tenemos estimadores:

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

$$\begin{cases} \hat{\beta}_1 = \frac{1}{nV_x} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\ \hat{\beta}_0 = \bar{y} - \bar{x} \hat{\beta}_1 \end{cases} \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

$$\begin{cases} E(Y) = X \cdot \beta \\ \text{cov}(Y) = \sigma^2 I_n \end{cases}$$

Miramos algunas propiedades, empezando por la media:

$$\begin{aligned} E(\hat{\beta}_1) &= \frac{1}{nV_x} \sum_{i=1}^n (x_i - \bar{x}) \underbrace{E(y_i)}_{\beta_0 + \beta_1 x_i} = \frac{1}{nV_x} \left[\underbrace{\beta_0 \sum_{i=1}^n (x_i - \bar{x})}_{=0} + \underbrace{\beta_1 \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})}_{=nV_x} \right] \\ &= \beta_1 \quad (\hat{\beta}_1 \text{ es un estimador insesgado de } \beta_1) \end{aligned}$$

$$E(\hat{\beta}) = (X^T X)^{-1} X^T \underbrace{E(Y)}_{X \cdot \beta} = \beta \quad (\hat{\beta} \text{ estimador insesgado de } \beta)$$

Hemos hecho los cálculos en una línea usando la notación matricial.

$$\begin{aligned} \text{cov}(\hat{\beta}) &= (X^T X)^{-1} X^T \underbrace{\text{cov}(Y)}_{\sigma^2 I_n} [(X^T X)^{-1} X^T]^T = \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} = \\ &= \sigma^2 (X^T X)^{-1} \quad \leftarrow \text{caso bidim.} \\ &= \sigma^2 \frac{1}{nV_x} \begin{pmatrix} \bar{x}^2 & -\bar{x} \\ -\bar{x} & 1 \end{pmatrix} \end{aligned}$$

RECUERDO:

$$\begin{array}{c} x_1 \ y_1 \\ | \quad | \\ x_n \ y_n \end{array}$$

$$Y = X\beta + \varepsilon$$

$$E(\varepsilon) = \vec{0}$$

$$\text{cov}(\varepsilon) = \sigma^2 I_n$$

$$Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}$$

$$\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}$$

$$X = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}$$

$$\varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix}$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y \longrightarrow$$

$$\left\{ \begin{array}{l} E(\hat{\beta}) = \beta \\ \text{cov}(\hat{\beta}) = \sigma^2 (X^T X)^{-1} \end{array} \right.$$

$$E(\hat{\beta}_0) = \beta_0, \quad E(\hat{\beta}_1) = \beta_1 \quad \checkmark \text{ insesgados}$$

$$V(\hat{\beta}_0) = \frac{\sigma^2}{nV_x} \bar{x}^2 = \frac{\sigma^2}{nV_x} \left[\frac{1}{n} + \frac{\bar{x}^2}{nV_x} \right]$$

$$V(\hat{\beta}_1) = \frac{\sigma^2}{nV_x}$$

$$\text{cov}(\hat{\beta}_0, \hat{\beta}_1) = -\frac{\bar{x}}{nV_x} \sigma^2$$

¿Estimador para σ^2 ?

$$e^T e = \sum_{i=1}^n e_i = Y^T (I_n - H) Y$$

$$H = X(X^T X)^{-1} X^T$$

$$\begin{aligned} \text{rang}(H) &= \text{tr}(H) = \text{tr}(X(X^T X)^{-1} X^T) = \text{tr}(X^T X (X^T X)^{-1}) = \\ &= \text{tr}(I_2) = 2 \end{aligned}$$

$$\text{rang}(I_n - H) = \text{tr}(I_n - H) = n - 2$$

¿Media de $Y^T (I_n - H) Y$?

general: A simétrica, $n \times n$

$\begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} = X$ vector aleat.

$$\left\{ \begin{array}{l} E(X) = \mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix} \\ \text{cov}(X) = \Sigma \end{array} \right.$$

forma cuadrática $X^T A X \longrightarrow$ variable aleatoria

$$\boxed{E(\mathbf{X}^T \mathbf{A} \mathbf{X}) = \text{tr}(\mathbf{A} \Sigma) + \mu^T \mathbf{A} \mu}$$

ALGO QUE TUVIAMOS QUE HACER EN EL TEMA 1

Comentario: \mathbf{A} simétrica, $\mathbf{X} \begin{matrix} \nearrow M \\ \searrow \Sigma \end{matrix}$

$$\boxed{V(\mathbf{X}^T \mathbf{A} \mathbf{X}) = 2 \text{tr}(\mathbf{A} \Sigma \mathbf{A} \Sigma) + 4 \mu^T \mathbf{A} \Sigma \mathbf{A} \mu}$$

En nuestro caso $\mathbf{Y} \begin{matrix} \nearrow \vec{0} \\ \searrow \sigma^2 \mathbf{I}_n \end{matrix}$

$$V(\mathbf{Y}^T (\mathbf{I}_n - \mathbf{H}) \mathbf{Y}) = 2\sigma^4 \underbrace{\text{tr}(\mathbf{I}_n - \mathbf{H})}_{=n-2}$$

$$V(S_R^2) = \frac{2\sigma^4}{n-2}$$

$$\mathbf{Z} = \frac{(n-2)S_R^2}{\sigma^2} \begin{matrix} \nearrow E(\mathbf{Z}) = n-2 \\ \searrow V(\mathbf{Z}) = 2(n-2) \end{matrix}$$

¿Distribución de $\hat{\beta}_0, \hat{\beta}_1, S_R^2$? A partir de ahora añadimos esta hipótesis:

$$\varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix} \sim N(\vec{0}, \sigma^2 \mathbf{I}_n)$$

$$\mathbf{Y} \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n)$$

$$\text{Resultado: } \hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \sim \mathcal{N}_2(\beta, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1})$$

$$\hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{nV_x})$$

$$\hat{\beta}_0 \sim N(\beta_0, \frac{\sigma^2}{nV_x} \bar{x}^2)$$

Corolario 2, tema 1

$$\frac{(n-1)S_R^2}{\sigma^2} = \frac{1}{\sigma^2} \mathbf{Y}^T (\mathbf{I}_n - \mathbf{H}) \mathbf{Y} \sim \chi_{n-2}^2$$

TEOREMA: En regresión lineal simple + normalidad:

- $\hat{\beta} \sim N(\beta, \sigma^2(x^T x)^{-1})$

- $\frac{(n-2)S_R^2}{\sigma^2} \sim \chi_{n-2}^2$

- $\hat{\beta}$ y S_R^2 son independientes

Obs: $\frac{\hat{\beta}_1 - \beta_1}{S_R \sqrt{\frac{1}{nV_x}}} = \frac{\frac{\hat{\beta}_1 - \beta_1}{\sigma/\sqrt{nV_x}} \cdot \sqrt{n-2}}{\sqrt{\frac{S_R^2(n-2)}{\sigma^2}}} \sim t_{n-2}$

$\frac{\hat{\beta}_1 - \beta_1}{\sigma/\sqrt{nV_x}} \sim N(0,1)$
 $\sqrt{\frac{S_R^2(n-2)}{\sigma^2}} \sim \sqrt{\chi_{n-2}^2}$

COROLARIO: $\frac{\hat{\beta}_1 - \beta_1}{S_R \sqrt{\frac{1}{nV_x}}} \sim t_{n-2}, \quad \frac{\hat{\beta}_0 - \beta_0}{S_R \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{nV_x}}} \sim t_{n-2}$

Podemos ahora conseguir intervalos de confianza para $\beta_0, \beta_1 \Rightarrow$
 \Rightarrow percentil de t-Student.

Observación (Un último detalle, sobre sumas de cuadrados)

$$\begin{aligned} TSS &= \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (Y_i^2 + \bar{Y}^2 - 2Y_i\bar{Y}) = \sum_{i=1}^n Y_i^2 + n\bar{Y}^2 - 2n\bar{Y}^2 = \\ &= \sum_{i=1}^n Y_i^2 - n\bar{Y}^2 = \left(1 - \frac{1}{n}\right) \sum_{i=1}^n Y_i^2 - \frac{1}{n} \sum_{i \neq j} Y_i Y_j \end{aligned}$$

Adicionalmente: $J_n = \text{ones}(n) = \begin{pmatrix} 1 & \dots & 1 \\ 1 & \dots & 1 \\ \vdots & & \vdots \\ 1 & \dots & 1 \end{pmatrix} \rightarrow \frac{1}{n} J_n \begin{cases} \text{simétrica} \\ \text{idempotente} \\ \text{rango } 1 \end{cases}$

$\Rightarrow TSS = Y^T \left(I_n - \frac{1}{n} J_n \right) Y$

simétrica, idempotente, rango $n-1$

$\Rightarrow RSS = Y^T (I_n - H) Y \rightarrow \text{rango } n-2$

$\Rightarrow MSS = Y^T \left(H - \frac{1}{n} J_n \right) Y \rightarrow \text{rango } 1$

Como $(I_n - H)(H - \frac{1}{n}J_n) = H - \frac{1}{n}J_n - H + \frac{1}{n} \underbrace{HJ_n}_{=J_n} = 0$

$\downarrow H = X(X^T X)^{-1} X^T$

$HX = X(X^T X)^{-1} X^T X = X$

$\begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}$

$H \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \rightarrow HJ_n = J_n$

\Rightarrow Por el teorema 3, RSS y MSS son independientes.

Con todos los preliminares expuestos, empezamos la Estadística

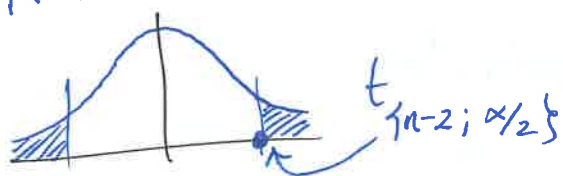
$\underbrace{\beta_0, \beta_1, \sigma^2}_{\text{desconocidos}} \xrightarrow{\text{generan}} \begin{matrix} x_1 y_1 \\ \vdots \\ x_n y_n \end{matrix} \rightarrow \hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} \quad S_R^2 = \frac{1}{n-2} \sum e_i^2$

$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \quad \hat{y} = \begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{pmatrix} \rightarrow e = \begin{pmatrix} e_1 \\ \vdots \\ e_n \end{pmatrix}$

①. Intervalos de confianza

Usamos el corolario que relaciona $\hat{\beta}_1, \hat{\beta}_0$ con t-Student.

$\mathbb{P}\left(\left|\frac{\hat{\beta}_1 - \beta_1}{S_R \sqrt{1/nV_x}}\right| > t_{\{n-2; \alpha/2\}}\right) = \alpha$



$1 - \alpha = \mathbb{P}\left(\hat{\beta}_1 - t_{\{n-2; \alpha/2\}} S_R \sqrt{1/nV_x} \leq \beta_1 \leq \hat{\beta}_1 + t_{\{n-2; \alpha/2\}} S_R \sqrt{1/nV_x}\right)$

$\Rightarrow \boxed{IC_{1-\alpha}(\beta_1) = \hat{\beta}_1 \pm t_{\{n-2; \alpha/2\}} S_R \sqrt{1/nV_x}}$

Por otro lado, desarrollando análogamente:

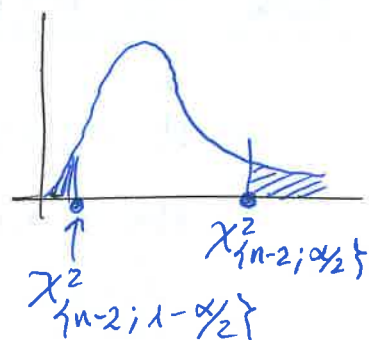
$\Rightarrow \boxed{IC_{1-\alpha}(\beta_2) = \hat{\beta}_2 \pm t_{\{n-2; \alpha/2\}} S_R \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{nV_x}}}$

Ahora usamos que $\frac{(n-2) S_R^2}{\sigma^2} \sim \chi^2_{n-2}$

$$\mathbb{P}\left(\chi^2_{n-2; 1-\alpha/2} \leq \frac{(n-2) S_R^2}{\sigma^2} \leq \chi^2_{n-2; \alpha/2}\right) = 1-\alpha$$

$$\mathbb{P}\left(\frac{(n-2) S_R^2}{\chi^2_{n-2; \alpha/2}} \leq \sigma^2 \leq \frac{(n-2) S_R^2}{\chi^2_{n-2; 1-\alpha/2}}\right) = 1-\alpha$$

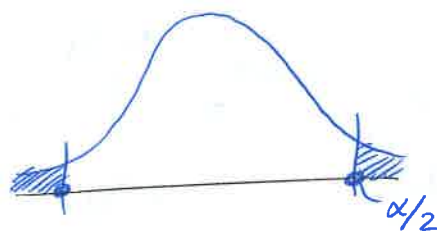
$$\Rightarrow \text{IC}_{1-\alpha}(\sigma^2) = \left(\frac{(n-2) S_R^2}{\chi^2_{n-2; \alpha/2}}, \frac{(n-2) S_R^2}{\chi^2_{n-2; 1-\alpha/2}} \right)$$



② Contraste de hipótesis

$$H_0: \beta_1 = 0$$

$$\text{Bajo } H_0, \frac{\hat{\beta}_1}{S_R \sqrt{1/n V_x}} \sim t_{n-2}$$



Región de rechazo (con nivel de significación α):

$$\mathcal{R} = \left\{ \left| \frac{\hat{\beta}_1}{S_R \sqrt{1/n V_x}} \right| > t_{n-2; \alpha/2} \right\}$$

Alternativa:

$$RSS = Y^T (I_n - H) Y \quad \leftarrow \text{independientes}$$

$$MSS = Y^T \left(H - \frac{1}{n} I_n \right) Y \quad \leftarrow$$

$$\text{Si } \beta_1 = 0: X\beta = \begin{pmatrix} 1 & x_1 \\ 1 & 1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \begin{pmatrix} \beta_0 \\ 0 \end{pmatrix} = \beta_0 \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

$$Y \sim N\left(\beta_0 \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}, \sigma^2 I_n\right) = \beta_0 \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} + X \quad \sim N(\beta_0, \sigma^2 I_n)$$

$$Y^T (I_n - H) Y = \beta_0^2 (1 \dots 1) (I_n - H) \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} + \sigma^2 I_n X^T (I_n - H) X \sim \chi^2_{n-2}$$

$$\Rightarrow \left\{ \begin{array}{l} RSS \sim \chi^2_{n-2} \\ MSS \sim \chi^2_1 \end{array} \right\} \xrightarrow{\text{indep.}} \frac{MSS/1}{RSS/(n-2)} \sim F_{1, n-2} \quad \text{F de Fisher}$$

$$\text{Región de rechazo: } \mathcal{R} = \left\{ \frac{MSS}{RSS/(n-2)} > F_{1, n-2; \alpha} \right\}$$

Queremos estimar $E(Y|X=x_0)$

$$\beta_0 + \beta_1 x_0$$

Lo hacemos con $Z = \hat{\beta}_0 + \hat{\beta}_1 x_0$

$$\begin{aligned} E(Z) &= \beta_0 + \beta_1 x_0 \\ V(Z) &= V(\hat{\beta}_0) + x_0^2 V(\hat{\beta}_1) + 2x_0 \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) \\ &= \sigma^2 \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{nV_x} \right] \end{aligned}$$

Así que $\frac{Z - (\beta_0 + \beta_1 x_0)}{S_R \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{nV_x}}} \sim t_{n-2}$

Por tanto, el intervalo de confianza para la media condicionada:

$$IC_{1-\alpha}(\beta_0 + \beta_1 x_0) = (\hat{\beta}_0 + \hat{\beta}_1 x_0) \pm t_{\{n-2; \alpha/2\}} S_R \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{nV_x}}$$

Queremos ahora estimar $Y_0 = \beta_0 + \beta_1 x_0 + \underbrace{\varepsilon_0}_{\text{nuevo}} \sim \mathcal{N}(0, \sigma^2)$ indep.

Lo hacemos con $Z = \hat{\beta}_0 + \hat{\beta}_1 x_0 \rightarrow E(Y_0 - Z) = 0$

$$\begin{aligned} V(Y_0 - Z) &= V(Y_0) + V(Z) - 2\text{Cov}(Y_0, Z) \\ &= \sigma^2 \left[1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{nV_x} \right] \end{aligned}$$

$$\Rightarrow \frac{Y_0 - Z}{S_R \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{nV_x}}} \sim t_{n-2}$$

Por tanto, el intervalo de confianza para la predicción del valor:

$$IC_{1-\alpha}(Y_0) = (\hat{\beta}_0 + \hat{\beta}_1 x_0) \pm t_{\{n-2; \alpha/2\}} S_R \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{nV_x}}$$

REGRESIÓN LINEAL MÚLTIPLE ($\vec{x}_j \in \mathbb{R}^d$)

Datos $\rightarrow \underbrace{X_1, \dots, X_K}_{\text{var. regresoras}}, \underbrace{Y}_{\text{var. respuesta}}$ $K=1 \Rightarrow$ RL simple

Consideramos $\boxed{n \geq K+2}$
Suponemos columnas lin. independientes

$$n \left\{ \begin{array}{c|c} X_1 & \dots & X_K & Y \\ \hline x_{11} & \dots & x_{1K} & y_1 \\ \vdots & & \vdots & \vdots \\ x_{n1} & \dots & x_{nK} & y_n \end{array} \right.$$

$\underbrace{\hspace{10em}}_K \quad \underbrace{\hspace{1em}}_1$

Modelo genérico: $Y | X_1=x_1 \dots X_K=x_K \equiv \beta_0 + \beta_1 x_1 + \dots + \beta_K x_K + \varepsilon$

$$\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1K} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{nK} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_K \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix} \quad \mathbb{Y} = \mathbb{X} \cdot \beta + \varepsilon$$

(notación matricial)

Hipótesis sobre ε : $\begin{cases} \mathbb{E}(\varepsilon) = \vec{0} \\ \text{cov}(\varepsilon) = \sigma^2 I_n \end{cases}$ Casi siempre $\varepsilon \sim N(\vec{0}, \sigma^2 I_n)$

① ESTIMADORES PARA PARÁMETROS

Con el mismo argumento (mínimos cuadrados) que en el caso simple, se obtienen los estimadores:

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \vdots \\ \hat{\beta}_K \end{pmatrix} \text{ de } \beta = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_K \end{pmatrix} \text{ dados por } \boxed{\hat{\beta} = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbb{Y}}$$

obs: \mathbb{X} es $n \times (K+1)$: $\begin{matrix} \boxed{\mathbb{X}^T} & \boxed{\phantom{\mathbb{X}}} \\ (K+1) \times n & n \times (K+1) \end{matrix} = \begin{matrix} \boxed{\phantom{\mathbb{X}}} \\ (K+1) \times (K+1) \end{matrix}$ simétrica y def. pos.

Esto es consecuencia de un resultado general:
Sea B matriz $n \times p$, rango p . $A = B^T B$ es $p \times p$ y simétrica
Vamos a ver que A es def. pos. Consideramos $\vec{y}^T A \vec{y} \stackrel{?}{\geq} 0$
 $\vec{y}^T A \vec{y} = \vec{y}^T B^T B \vec{y} = (B \vec{y})^T (B \vec{y}) = \|B \vec{y}\|^2 \geq 0$ y es cero solo si $B \vec{y} = \vec{0}$ (como B rango máx. solo cuando $\vec{y} = \vec{0}$).

Obs: Si suponemos normalidad, los estimadores por máx. vero son los mismos.

TEOREMA (Gauss-Markov): $Y = X \cdot \beta + \varepsilon$, X rango máximo

$$\begin{cases} E(\varepsilon) = 0 \\ \text{cov}(\varepsilon) = \sigma^2 I_n \end{cases}$$

Cada estimador $\hat{\beta}_j$ (de mínimos cuadrados) $j=0 \dots k$ tiene mínima varianza de entre todos los estimadores lineales e insesgados de β_j .

demonstración: El estimador de mínimos cuadrados es

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

Comparamos con cualquier otro:

$$\tilde{\beta} = A \cdot Y$$

Para que sea insesgado $E(\tilde{\beta}) = A \cdot E(Y) =$

$$= A X \cdot \beta = \beta \quad \text{para todo } \beta.$$

$$\Rightarrow A X = I_{k+1}$$

las varianzas de cada $\hat{\beta}_j$ están en la diagonal

$$\text{Calculamos } \text{cov}(\tilde{\beta}) = \text{cov}(A Y) = A \text{cov}(Y) A^T = \sigma^2 A A^T$$

¿Qué A minimiza los elementos de esa diagonal?

$$A A^T = (A - (X^T X)^{-1} X^T + (X^T X)^{-1} X^T) \cdot (A - (X^T X)^{-1} X^T + (X^T X)^{-1} X^T)^T =$$

$$= (A - (X^T X)^{-1} X^T) (A - (X^T X)^{-1} X^T)^T +$$

$$+ (X^T X)^{-1} X^T X (X^T X)^{-1} + \underbrace{(A - (X^T X)^{-1} X^T) X (X^T X)^{-1}}_{=0} +$$

$$+ \underbrace{\boxed{\hspace{2cm}}}_{=0} =$$

$$= [A - (X^T X)^{-1} X^T] [A - (X^T X)^{-1} X^T]^T + (X^T X)^{-1}$$

matriz semidef. pos. \rightarrow los elementos de la diagonal son $\geq 0 \Rightarrow$

$$\Rightarrow \text{se hacen } 0 \text{ si } A = (X^T X)^{-1} X^T$$

2. Pronóstico Y Residuos

$$\hat{Y} = X \cdot \hat{\beta} = \underbrace{X(X^T X)^{-1} X^T}_{\substack{\text{"} \\ H \text{ matriz hat } n \times n}} \cdot Y = HY$$

$HX = X$
simétrica
idempotente
rango es $k+1$

Los residuos son: $\begin{pmatrix} e_1 \\ \vdots \\ e_n \end{pmatrix} = \vec{e} = Y - \hat{Y} = (I_n - H)Y =$

$= (I_n - H)(X\beta + \epsilon) =$

$= X\beta + \epsilon - HX\beta - H\epsilon =$

$= \underline{(I_n - H)\epsilon}$

Propiedades

Perpendicularidad:

$$X^T \cdot e = 0_{k+1}$$

$$\hat{Y}^T \cdot e = 0$$

$$E(e) = \vec{0}$$

$$\text{cov}(e) = \text{cov}((I_n - H)\epsilon) = (I_n - H) \overbrace{\text{cov}(\epsilon)}^{\sigma^2 I_n} (I_n - H)^T =$$
$$= \sigma^2 (I_n - H)$$

\rightarrow no es diagonal!

$$SCR = e^T e = \sum_{i=1}^n e_i^2 = Y^T (I_n - H) Y$$

El estimador para σ^2 sería $S_R^2 = \frac{1}{n-k-1} SCR =$

$$= \frac{1}{n-k-1} \sum_{i=1}^n e_i^2 = \frac{1}{n-k-1} Y^T (I_n - H) Y$$

3. SUMAS DE CUADRADOS

$$SCT = TSS = \sum_{i=1}^n (Y_i - \bar{Y})^2 = Y^T (I_n - \frac{1}{n} J_n) Y \quad (\text{RANGO } n-1)$$

matriz de unos

$$SCM = MSS = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 = Y^T (H - \frac{1}{n} J_n) Y \quad (\text{RANGO } K)$$

$$SCR = RSS = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = Y^T (I_n - H) Y \quad (\text{RANGO } n-K-1)$$

$$\boxed{TSS = MSS + RSS}$$

las 3 matrices
son sim. e idempot

$$R^2 = \frac{MSS}{TSS} = 1 - \frac{RSS}{TSS} \quad \text{entre 0 y 1}$$

Nota: si se añaden variables regresoras (sin influencia sobre Y), aumenta R^2 .

Por esta razón hay quien usa: $\bar{R}^2 = 1 - \frac{RSS/n-K-1}{TSS/n-1}$

Necesitamos distribución de los estimadores.

MODELO NORMAL: Suponemos que $E \sim N_n(0, \sigma^2 I_n)$.

Por tanto, $Y \sim N(X \cdot \beta, \sigma^2 I_n)$.

Como $\hat{\beta} = (X^T X)^{-1} X^T Y$, $\hat{\beta} \sim N_{K+1}(\beta, \sigma^2 (X^T X)^{-1})$

Como $SCR = Y^T \underbrace{(I_n - H)}_{\text{rango } n-K-1} Y$, $H = X(X^T X)^{-1} X^T$

y además

$$\begin{aligned} (X\beta)^T (I_n - H) (X\beta) &= \beta^T X^T (I_n - H) X \beta = \\ &= \beta^T X^T X \beta - \beta^T X^T X (X^T X)^{-1} X^T X \beta = 0 \end{aligned}$$

Concluimos que

$$\frac{SCR}{\sigma^2} = \frac{(n-K-1) S_R^2}{\sigma^2} \sim \chi^2_{n-K-1}$$

Además, como $(X^T X)^{-1} \times (I_n - H) = 0$

resulta que $\hat{\beta}$ y SCR (o bien S_R^2) son independientes.

Llamemos $(X^T X)^{-1} = \begin{pmatrix} q_{00} & \dots & \dots \\ \vdots & q_{11} & \dots \\ \vdots & \vdots & \ddots & q_{KK} \end{pmatrix}$ q_{ii} es el elemento en la posición $i+1$ de esta matriz.

\uparrow
(K+1) x (K+1)

Cada $\hat{\beta}_j \sim N(\beta_j, \sigma^2 q_{jj}) \quad j=0, \dots, K$

Y por tanto, $\frac{\hat{\beta}_j - \beta_j}{S_R \sqrt{q_{jj}}} \sim t_{n-K-1} \quad j=0, \dots, K.$

Inmediatamente, dada una muestra $\begin{array}{c|c} 1 & \dots & 1 \\ \vdots & & \vdots \\ 1 & \dots & 1 \end{array} \begin{array}{c} y_1 \\ \vdots \\ y_n \end{array}$

A) $IC_{1-\alpha}(\beta_j) = \hat{\beta}_j \pm t_{n-K-1; \alpha/2} S_R \sqrt{q_{jj}} \quad j=0, \dots, K$

B) $H_0: \beta_j = 0 \quad j=1, \dots, K$

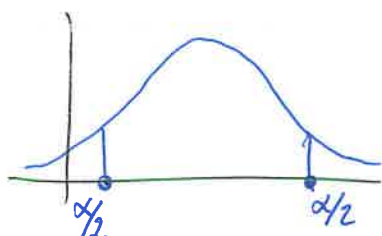
Región de rechazo para nivel de significación α

$R_j(\alpha) = \left\{ \left| \frac{\hat{\beta}_j}{S_R \sqrt{q_{jj}}} \right| > t_{n-K-1; \alpha/2} \right\}$

Nota 1: Bajo normalidad, $\hat{\beta}$ es estimador max. VERO de β
Gauss-Markov dice que $\hat{\beta}$ es de mínima varianza de entre estimadores lineales e insesgados de β .
(bajo normalidad \nearrow)

c) Intervalo de confianza para σ^2 :

$IC_{1-\alpha}(\sigma^2) = \left(\frac{(n-K-1) S_R^2}{\chi^2_{n-K-1; \alpha/2}}, \frac{(n-K-1) S_R^2}{\chi^2_{n-K-1; 1-\alpha/2}} \right)$



CONTRASTE GLOBAL DE LA REGRESIÓN

$$H_0: \beta_1 = \dots = \beta_k = 0$$

Bajo H_0 , $Y \sim N(X\beta, \sigma^2 I_n) \sim N\left(\beta_0 \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \sigma^2 I_n\right)$

Se analizan: por un lado $\frac{RSS}{\sigma^2} = \frac{1}{\sigma^2} Y^T (I_n - H) Y \sim \chi^2_{n-k-1}$ tema 1
tema 2

por otro lado $\frac{MSS}{\sigma^2} = \frac{1}{\sigma^2} Y^T \left(H - \frac{1}{n} J_n\right) Y \sim \chi^2_k$

↓ rango k
En general, no es χ^2_k
pero bajo H_0 , sí:

$$\beta_0^2 (1 \dots 1) \left(H - \frac{1}{n} J_n\right) \begin{pmatrix} 1 \\ 1 \end{pmatrix} = 0$$

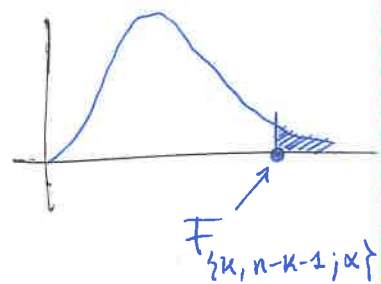
recuerdo $\rightarrow H \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$

$\rightarrow \frac{1}{n} J_n \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$

corolario 2, tema 1

Conclusión, bajo H_0 :

$$\frac{MSS/k}{RSS/(n-k-1)} \sim F_{k, n-k-1}$$



Región de rechazo para H_0 con nivel de sign. α :

$$R = \left\{ \frac{MSS/k}{RSS/(n-k-1)} > F_{k, n-k-1; \alpha} \right\}$$

Habitual: tabla ANOVA

| | Grados | Sumas cuadrados | | |
|--------------------|--------|-----------------|---------------|-----------------------------|
| MODELO (regresión) | k | MSS | MSS/k | $\frac{MSS/k}{RSS/(n-k-1)}$ |
| RESIDUOS | n-k-1 | RSS | $RSS/(n-k-1)$ | p-valor |
| TOTAL | n-1 | TSS | | |

PREDICCIÓN

| | |
|------------------------------------|----------|
| $x_1 \dots x_k$ | y |
| $x_{11} \dots x_{k1}$ | y_1 |
| $\vdots \quad \ddots \quad \vdots$ | \vdots |
| $x_{n1} \dots x_{kn}$ | y_n |

Nuevo vector $x_0^T = (x_{01}, \dots, x_{0k})$

$y_0?$

$$\Rightarrow \tilde{x}_0^T = (1, x_{01}, \dots, x_{0k}) \quad , \quad \hat{y}_0?$$

a) Queremos estimar $E(Y_0 | x_1 = x_{01}, \dots, x_k = x_{0k}) = \tilde{x}_0^T \beta =$
 $= \beta_0 + \beta_1 x_{01} + \dots + \beta_k x_{0k}$

$$\hat{z} = \tilde{x}_0^T \cdot \hat{\beta} = \hat{\beta}_0 + \hat{\beta}_1 x_{01} + \dots + \hat{\beta}_k x_{0k} = \tilde{x}_0^T (X^T X)^{-1} X^T Y$$

$$E(\hat{z}) = \tilde{x}_0^T (X^T X)^{-1} X^T X \beta = \tilde{x}_0^T \beta$$

$$V(\hat{z}) = \tilde{x}_0^T (X^T X)^{-1} X^T \sigma^2 I_n X (X^T X)^{-1} \tilde{x}_0 = \sigma^2 \tilde{x}_0^T (X^T X)^{-1} \tilde{x}_0$$

Bajo normalidad $\frac{\tilde{x}_0^T \hat{\beta} - \tilde{x}_0^T \beta}{S_R \sqrt{\tilde{x}_0^T (X^T X)^{-1} \tilde{x}_0}} \sim t_{n-k-1}$

Así que

$$\boxed{IC_{1-\alpha}(\tilde{x}_0^T \beta) = \tilde{x}_0^T \hat{\beta} \pm t_{\{n-k-1; \alpha/2\}} S_R \sqrt{\tilde{x}_0^T (X^T X)^{-1} \tilde{x}_0}}$$

b) Predecir $Y_0 = \tilde{x}_0^T \beta + \epsilon_0 \sim N(0, \sigma^2)$
 indep. de las $\epsilon_1, \dots, \epsilon_n$

con $\hat{z} = \tilde{x}_0^T \hat{\beta}$

$$E(Y_0 - \hat{z}) = 0$$

$$V(Y_0 - \hat{z}) = V(Y_0) + V(\hat{z}) = \sigma^2 (1 + \tilde{x}_0^T (X^T X)^{-1} \tilde{x}_0)$$

Así que,

$$\frac{Y_0 - \tilde{x}_0^T \hat{\beta}}{S_R \sqrt{1 + \tilde{x}_0^T (X^T X)^{-1} \tilde{x}_0}} \sim t_{n-k-1}$$

$$\boxed{IC_{1-\alpha}(Y_0) = \tilde{x}_0^T \hat{\beta} \pm t_{\{n-k-1; \alpha/2\}} S_R \sqrt{1 + \tilde{x}_0^T (X^T X)^{-1} \tilde{x}_0}}$$

① INTERPRETACIÓN DE LOS CONTRASTES

| Contraste global (F) | Contrastes individuales (t) | |
|----------------------------------|--------------------------------|----------------------------------|
| Modelo explicativo | Todas las X_i explicativas | → ✓ |
| | Algunas de las X_i explicat. | → te quedas con las explicativas |
| | Ninguna explicativa | → ups! colinealidad |
| Modelo ^{no} explicativo | Todas | → ups! colineat. |
| | Alguna | → ups! " |
| | Ninguna | → basura |

② INTERVALOS DE CONFIANZA SIMULTÁNEOS

$$Y = X \cdot \beta + \epsilon$$

β parámetros

$$\hat{\beta} = (X^T X)^{-1} X^T Y \quad (\text{estimador})$$

$X^T X$ es simétrica, def. pos.

Tomamos raíz cuadrada:

$$X^T X = P \Lambda P^T = \underbrace{P \Lambda^{1/2} P^T}_{\substack{\uparrow \\ \text{ortogonal}}} \underbrace{P \Lambda^{1/2} P^T}_{(X^T X)^{1/2}}$$

Consideramos el vector aleatorio

$$V = (X^T X)^{1/2} (\hat{\beta} - \beta) \Rightarrow$$

$$\Rightarrow \begin{cases} E(V) = 0 \\ \text{cov}(V) = (X^T X)^{1/2} \text{cov}(\hat{\beta}) (X^T X)^{1/2} = \sigma^2 I \end{cases} \quad \text{y } V \text{ es normal.}$$

$$\text{cov}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$$

$$\Rightarrow \frac{V^T V}{\sigma^2} = \frac{1}{\sigma^2} (\hat{\beta} - \beta)^T (X^T X) (\hat{\beta} - \beta) \sim \chi^2_{k+1}$$

por otro lado

$$\frac{n-k-1}{\sigma^2} S_R^2 \sim \chi^2_{n-k-1}$$

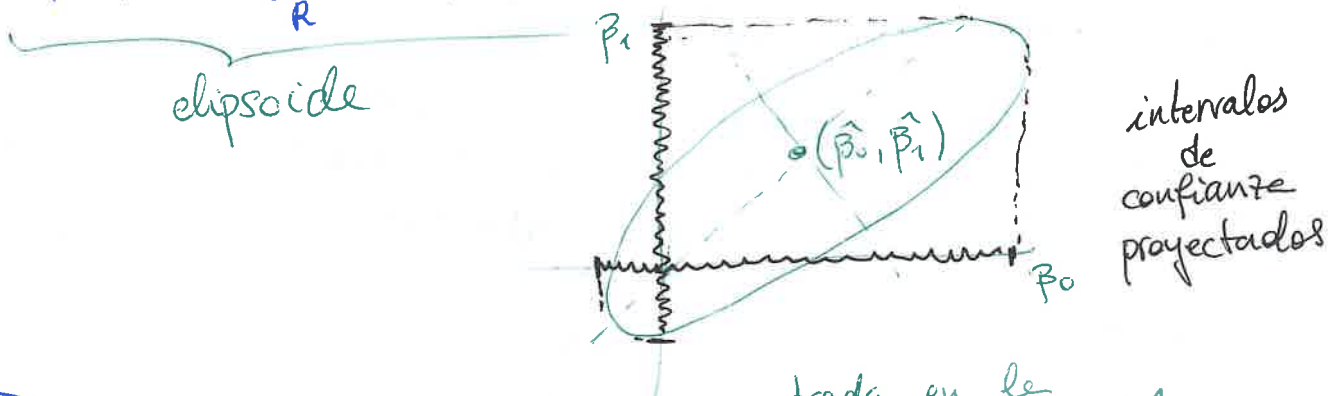
Así que:

$$\frac{V^T V / (k+1)}{S_R^2} \sim F_{k+1, n-k-1}$$

independientes

Entonces, con probabilidad $1-\alpha$:

$$(\hat{\beta} - \beta)^T \frac{(X^T X)}{S_R^2} (\hat{\beta} - \beta) \leq (k+1) F_{1, k+1, n-k-1; \alpha}$$



Proyectando, se obtiene

→ entrada en la diagonal de $(X^T X)^{-1}$

$$IC_{1-\alpha}(\beta_j) = \hat{\beta}_j \pm S_R \sqrt{f_{jj}} \sqrt{(k+1) F_{1, k+1, n-k-1; \alpha}}$$

③ VALIDACIÓN DEL MODELO

Lista de hipótesis que hemos hecho:

- Linealidad de los parámetros
- $n \geq k+2$
- no colinealidad de las variables regresoras
- ε_i media 0
- ε_i varianza σ^2
- ε_i incorreladas
- ε_i normales (→ indep.)

Tenemos los residuos ε_i $i=1, \dots, n$

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

$$\hat{Y} = X \hat{\beta} = \underbrace{X (X^T X)^{-1} X^T}_{= H} Y$$

$$e = (I_n - H)Y = (I_n - H)E$$

$$E(e) = 0$$

$$\text{cov}(e) = \sigma^2 (I_n - H)$$

$$H = (h_{ij})$$

$$\sum_{j=1}^n h_{ij} = 1 \quad \sum_{i=1}^n h_{ii} = k+1$$

$$\frac{1}{n} \leq h_{ii} \leq 1 \quad \frac{1}{2} \leq h_{ij} \leq \frac{1}{2} \quad i \neq j$$

Si todos h_{ij} son "pequeños", los residuos son "casi independientes"

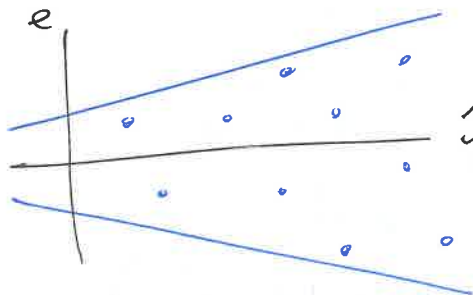
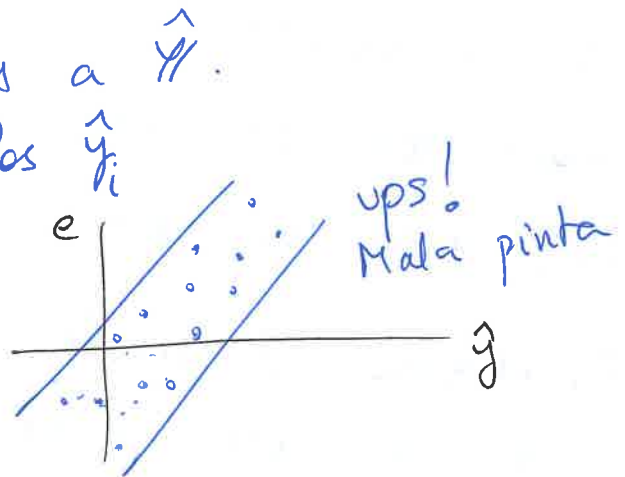
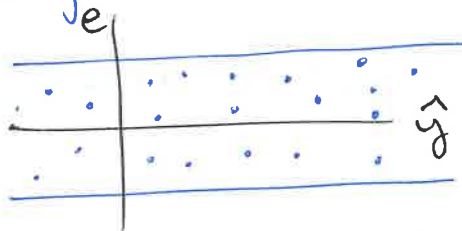
$$V(e_i) = \sigma^2(1-h_{ii})$$

Residuos estandarizados ("a la t")

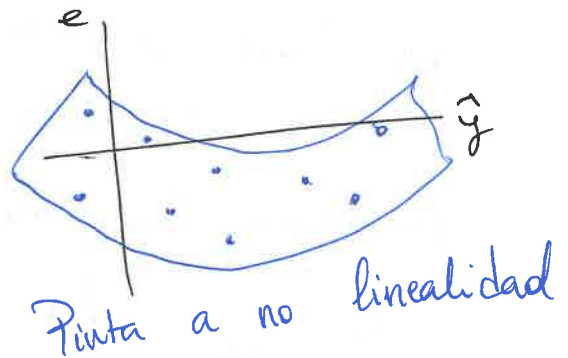
$$e_i^* = \frac{e_i}{S_R \sqrt{1-h_{ii}}} \sim \begin{matrix} \text{parecido} \\ \text{a } N(0,1) \end{matrix}$$

Residuos e son perpendiculares a \hat{Y} .

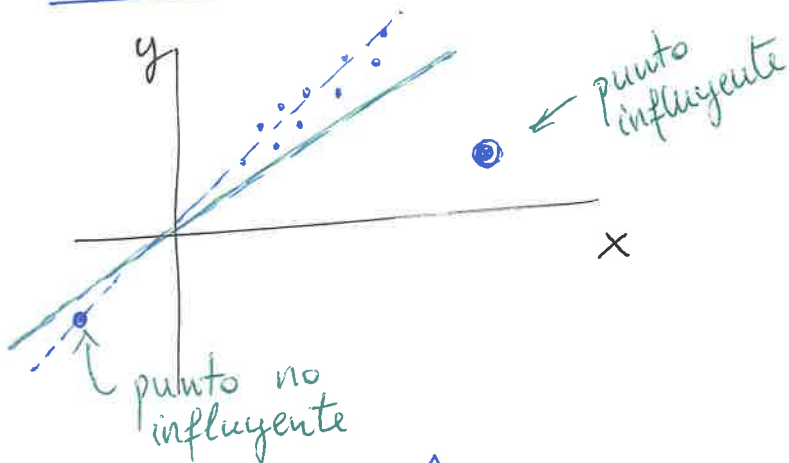
Dibujamos los e_i frente a los \hat{y}_i
perfecto!



puede que σ^2 no sea constante



④ PUNTOS "INFLUYENTES"



Una observación es influyente si tiene mucho impacto en $\hat{\beta}$ y $\hat{y} = x\hat{\beta}$

$$\hat{Y} = HY$$

$$\hat{y}_i = h_{ii} y_i + \sum_{j \neq i} h_{ij} y_j$$

Si h_{ii} es "grande" (próximo a 1), entonces \hat{y}_i viene casi determinado por y_i .

A h_{ii} se le llama el "leverage" (apalancamiento) de la observación i .

Todos los h_{ii} suman $k+1 \Rightarrow$ "en media" son $\frac{k+1}{n}$.
Hay quien dice que observación influyente es cuando $h_{ii} > 2 \frac{k+1}{n}$.

Formalizamos el concepto de "influencia" midiendo el efecto que tiene quitar esa observación:

Con todos los datos $\rightarrow \hat{\beta}, \hat{y}$

Quitando el dato " i " $\rightarrow \hat{\beta}_{(i)}, \hat{y}_{(i)}$

Comparamos (distancia de Cook)

$$D_i = \frac{1}{(k+1) S_R^2} (\hat{\beta} - \hat{\beta}_{(i)})^T (X^T X) (\hat{\beta} - \hat{\beta}_{(i)}) =$$
$$= \frac{1}{k+1} \left(\frac{e_i}{S_R \sqrt{1-h_{ii}}} \right)^2 \frac{h_{ii}}{1-h_{ii}}$$

TEMA 4 CLASIFICACIÓN

Contexto: tenemos unos individuos que pertenecen a ciertas poblaciones $\rightarrow \Pi_0, \Pi_1, \Pi_2, \dots$

De cada individuo tenemos unas medidas de ciertas magnitudes: X_1, X_2, \dots, X_k

Datos:

| | X_1 | X_2 | \dots | X_k | Población |
|----------|-------|-------|----------|-------|-----------|
| 1 | • | • | \dots | • | • |
| 2 | • | • | \dots | • | • |
| \vdots | | | \vdots | | \vdots |
| n | • | • | \dots | • | • |

Desarrollamos un procedimiento que clasifique individuos en las poblaciones

- el sistema se calibra con los datos disponibles
- se usa para clasificar datos nuevos

Sistemas

- reglas lineales
- regresión logística
- árboles de clasificación
- redes neuronales

Ejemplos

Marketing

Poblaciones
compradores
no compradores

Variables
- ingresos
- educación
- tipo familia
- compras anteriores

Medicine

desarrolla tumor
no

edad
hábitos
datos analíticos

Seguros

fraude
no fraude

ingresos
edad
#tarjetas de crédito
#partes anteriores

Finanzas

Quiebra
No quiebra

precio acciones
nivel endeudamiento
resultados anteriores

Formalizamos lo anterior:

$X = (X_1, \dots, X_k) \rightarrow$ vector de observaciones

Dos poblaciones $\rightarrow \Pi_0$ y Π_1

$\left\{ \begin{array}{l} X \text{ se distribuye en } \Pi_0 \text{ con } f_0(x_1 \dots x_k) \\ X \text{ se distribuye en } \Pi_1 \text{ con } f_1(x_1 \dots x_k) \end{array} \right.$

Esperamos que f_0 y f_1 sean "distintos"
(si iguales, da igual todo)

Normalmente, de f_0 y f_1 se sabe (o se postula)
cierta información:

\rightarrow parcial $\left\{ \begin{array}{l} \mu_1 \\ \mu_0 \end{array} \right\} \left\{ \begin{array}{l} \Sigma_0 \\ \Sigma_1 \end{array} \right.$

\rightarrow (semi) completa
 f_0 y f_1 normales con $\left\{ \begin{array}{l} \mu_0, \Sigma_0 \\ \mu_1, \Sigma_1 \end{array} \right.$

\rightarrow completa
 f_0, f_1 normales con $\left\{ \begin{array}{l} \mu_0, \Sigma_0 \\ \mu_1, \Sigma_1 \end{array} \right\}$ todo conocido

Objetivo: desarrollar un sistema/procedimiento que clasifique
"mal" con baja probabilidad
(quizás incluir un "coste" de mala clasificación)

① La idea de Fisher (1938)

Transformar los datos k -dimensionales en 1-dim,
proyectando sobre una dirección adecuada

$X = (X_1, \dots, X_k) \begin{cases} \rightarrow f_0 \text{ en } \Pi_0 \\ \rightarrow f_1 \text{ en } \Pi_1 \end{cases}$

Suponemos conocidos $\left\{ \begin{array}{l} \mu_0, \Sigma_0 \\ \mu_1, \Sigma_1 \end{array} \right.$ $\Sigma_0 = \Sigma_1 = \Sigma$

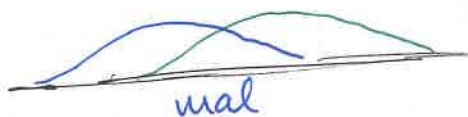
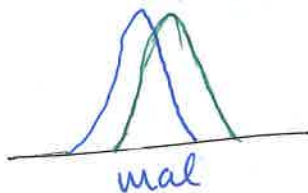
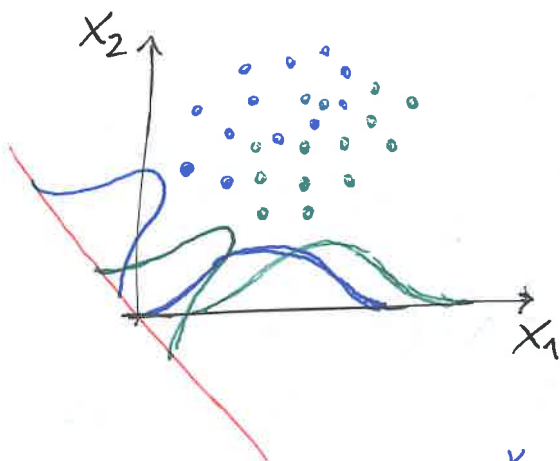
Proyectamos datos sobre dirección $a \in \mathbb{R}^k$

$$Z = a^T X \rightarrow \text{media} \quad \mathbb{E}(Z) = \begin{cases} a^T \mu_0 & \text{en } \Pi_0 \\ a^T \mu_1 & \text{en } \Pi_1 \end{cases}$$

varianza $V(Z) = a^T \Sigma a$

Queremos encontrar $a \in \mathbb{R}^k$:

- distancia entre las medias proyectadas sea grande
- varianza pequeña.



Buscamos $a \in \mathbb{R}^k$ tal que:

$$f(a) = \frac{|a^T(\mu_0 - \mu_1)|^2}{a^T \Sigma a} \text{ sea "m\u00e1xima" (maximizar } f(a))$$

Obs: si $\lambda \neq 0$, $f(\lambda a) = f(a) \rightarrow$ buscamos direcciones, y no tanto vectores

Recordatorio: (Desigualdad de Cauchy-Schwarz)

$x, y \in \mathbb{R}^k$ no nulos

$$(x^T y)^2 \leq (x^T x)(y^T y) \text{ con igualdad } \Leftrightarrow x = cy \text{ para cierto } c \in \mathbb{R}$$

Generalización: $x, y \in \mathbb{R}^k$, A matriz $k \times k$ sim. def. pos.,
entonces $(x^T y)^2 \leq (x^T A x)(y^T A^{-1} y)$ con igualdad si
y solo si $\begin{cases} x = c A^{-1} y \\ y = c A x \end{cases}$ para cierto $c \in \mathbb{R}$

demo

Escribimos $A = P\Lambda P^T$, $A^{-1} = P\Lambda^{-1}P^T$

Raíces cuadradas:

$$A = P\Lambda^{1/2}\Lambda^{1/2}P^T = (P\Lambda^{1/2})(P\Lambda^{1/2})^T$$

pero también

$$A = \underbrace{P\Lambda^{1/2}P^T}_{=C} P\Lambda^{1/2}P^T = C.C \quad (C.C^{-1} = Id)$$

$$A^{-1} = P\Lambda^{1/2}P^T P\Lambda^{-1/2}P^T = C^{-1}C^{-1}$$

Con esto:

$$(x^T y)^2 = (x^T I_k y)^2 = (x^T C C^{-1} y)^2 = [(Cx)^T (C^{-1}y)]^2 \leq$$

$$\leq (Cx)^T Cx \cdot (C^{-1}y)^T (C^{-1}y) = \underbrace{(x^T C C x)}_{=A} \cdot \underbrace{(y^T C^{-1} C^{-1} y)}_{=A^{-1}}$$

TEOREMA: El máximo de $f(a) = \frac{|a^T(\mu_0 - \mu_1)|^2}{a^T \Sigma a}$ se alcanza en múltiplos de $\boxed{a_m = \Sigma^{-1}(\mu_0 - \mu_1)}$

demonstración c-s general, con Σ

$$\frac{|a^T(\mu_0 - \mu_1)|^2}{a^T \Sigma a} \leq \frac{(\cancel{a^T \Sigma a}) (\mu_0 - \mu_1)^T \Sigma^{-1} (\mu_0 - \mu_1)}{\cancel{a^T \Sigma a}}$$

Si tomamos $a_m = \Sigma^{-1}(\mu_0 - \mu_1)$:

$f(a_m) = (\mu_0 - \mu_1)^T \Sigma^{-1} (\mu_0 - \mu_1)$ y ya está. \square

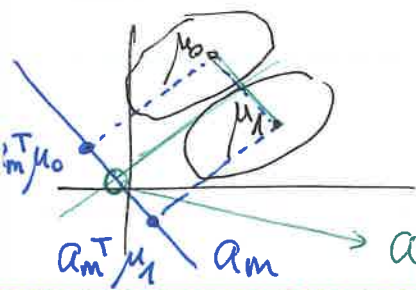
Esto nos da una REGLA DE CLASIFICACIÓN:

Calculamos $a_m = \Sigma^{-1}(\mu_0 - \mu_1)$

Dada observación $x = (x_1 - x_k)^T$,

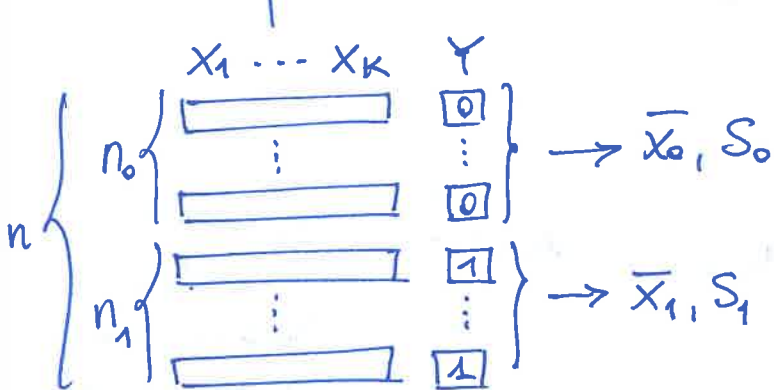
si $a_m^T x \geq a^T \left(\frac{\mu_0 + \mu_1}{2} \right) \Rightarrow x \text{ en } \Pi_1$

si $a_m^T x < \quad \quad \quad \Rightarrow x \text{ en } \Pi_0$



Esto último requiere conocer $\rightarrow \mu_0, \mu_1, \Sigma_0 = \Sigma_1 = \Sigma$

En la práctica:



Estimamos $\begin{array}{l} \hat{\mu}_0 = \bar{x}_0 \\ \hat{\mu}_1 = \bar{x}_1 \end{array}$

$$\hat{\Sigma} = \frac{(n_0-1)S_0 + (n_1-1)S_1}{n_0 + n_1 - 2}$$

S_{pooled}

La regla de clasificación sea:

① Calculamos $\hat{a}_m = S_{\text{pooled}}^{-1} (\bar{x}_0 - \bar{x}_1)$

② Dado $x = (x_1 \dots x_k)^T$,

si $\hat{a}_m^T x \geq \hat{a}_m^T \left(\frac{\bar{x}_0 + \bar{x}_1}{2} \right) \rightarrow x \text{ en } \Pi_1$

si $\hat{a}_m^T x < \hat{a}_m^T \left(\frac{\bar{x}_0 + \bar{x}_1}{2} \right) \rightarrow x \text{ en } \Pi_0$

② ¿Qué tal funciona?

Tasa de error aparente: TEA = porcentaje de puntos mal clasificados
(quizás convenga distinguir mal clasificados entre los de Π_0 y los de Π_1)

Trampa: usar los mismos datos para $\begin{array}{l} \rightarrow \text{crear la regla} \\ \rightarrow \text{validar su funcionamiento} \end{array}$

Alternativas: dividir los datos $\begin{array}{l} \rightarrow \text{training set} \\ \rightarrow \text{test set} \end{array}$

\rightarrow validación cruzada

REGRESIÓN LOGÍSTICA

$X_1, \dots, X_k \leftarrow$ variables regresoras

$Y \leftarrow$ respuesta (0,1 en nuestro caso)

Estructura de datos

| | X_1 | ... | X_k | Y |
|----------|----------|-----|----------|----------|
| dato 1 | x_{11} | ... | x_{1k} | y_1 |
| \vdots | \vdots | | \vdots | \vdots |
| dato n | x_{n1} | ... | x_{nk} | y_n |

Modelo $Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$

y_i serán Bernoulli's indeps. cada uno con su param.

$y_i \sim \text{Ber}(p_i)$

$p_i = P(Y_i = 1 | x_i)$

$Y \sim \text{Ber}(p(x)) \quad x = (x_1, \dots, x_k)$

¿ $p(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$?

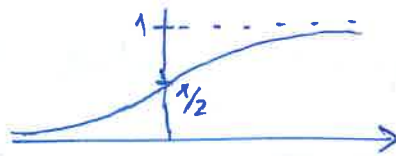
valores en $[0, 1]$ valores en \mathbb{R}

no funciona

Buscamos función $\mathbb{R} \rightarrow [0, 1]$

$h(x) = \frac{1}{1 + e^{-x}}$

función logística (logit)

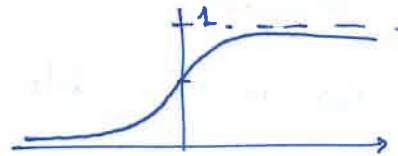


$h(0) = 1/2$

$h(-x) = 1 - h(x)$

$h'(x) = h(x)(1 - h(x))$

Alternativa, $\Phi(x)$ (probit)



Modelo Logit

$$p(x) = P(Y=1|x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k)}} = h(\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k)$$

① Interpretación de los parámetros $\beta_0, \beta_1, \dots, \beta_k$

Razón de probabilidades (odds): $O(x) = \frac{p(x)}{1-p(x)} = \frac{\frac{1}{1+e^{-(\dots)}}}{1 - \frac{1}{1+e^{-(\dots)}}} =$

$$= e^{\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k} \Rightarrow \ln\left(\frac{p(x)}{1-p(x)}\right) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$

Variaciones: $x = (x_1, \dots, x_k) \rightarrow x + \Delta_j = (x_1, \dots, x_j + 1, \dots, x_k)$

¿Cuánto cambian los parámetros?

$$\frac{O(x + \Delta_j)}{O(x)} = \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k} \cdot e^{\beta_j}}{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k}} = e^{\beta_j}$$

② Estimación de parámetros

Datos:

$$\begin{array}{ccc} x_1 & [x_{11} \dots x_{1k}] & [y_1] \\ \vdots & \vdots & \vdots \\ x_n & [x_{n1} \dots x_{nk}] & [y_n] \end{array}$$

Queremos estimar $\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}$

Usamos máx. verosimilitud

$$VERO(\beta_0, \dots, \beta_k) = \prod_{i=1}^n p(x_i)^{y_i} (1-p(x_i))^{1-y_i}$$

Variamos β_0, \dots, β_k para maximizar $VERO(\cdot)$
explícitamente es una expresión muy complicada \rightarrow solver (o un mejor maximizador)

(*) Un poco de notación:

$$\beta = (\beta_0, \dots, \beta_k)^T$$

$$x_i = (x_{i1}, \dots, x_{ik})^T \quad i = 1, \dots, n$$

$$\tilde{x}_i = (1, x_{i1}, \dots, x_{ik})^T \quad i = 1, \dots, n \quad (\text{añadimos columna de 1's})$$

$$\beta^T \cdot \tilde{x}_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik}$$

$$p(x_i) = \frac{1}{1 + e^{-\beta^T \tilde{x}_i}} = h(\beta^T \tilde{x}_i)$$

$$\frac{p(x_i)}{1 - p(x_i)} = e^{\beta^T \tilde{x}_i}$$

$$\begin{cases} h'(x) = h(x)(1-h(x)) \\ \frac{\partial}{\partial \beta_j} (\beta^T \tilde{x}_i) = x_{ij} \end{cases}$$

$$\frac{\partial}{\partial \beta_j} p(x_i) = h(\beta^T \tilde{x}_i) (1 - h(\beta^T \tilde{x}_i)) x_{ij} \quad \begin{matrix} i = 1, \dots, n \\ j = 1, \dots, k \end{matrix}$$

$$\text{VERO}(\beta) = \prod_{i=1}^n p(x_i)^{y_i} (1 - p(x_i))^{1-y_i}$$

$$\log \text{VERO}(\beta) = \sum_{i=1}^n (y_i \ln(p(x_i)) + (1-y_i) \ln(1-p(x_i))) =$$

$$= \sum_{i=1}^n y_i \ln\left(\frac{p(x_i)}{1-p(x_i)}\right) + \sum_{i=1}^n \ln(1-p(x_i)) =$$

$$= \sum_{i=1}^n y_i \cdot \beta^T \tilde{x}_i + \sum_{i=1}^n \ln(1-p(x_i))$$

$$\frac{\partial}{\partial \beta_j} \log \text{VERO}(\beta) = \sum_{i=1}^n y_i x_{ij} - \sum_{i=1}^n \frac{p(x_i)(1-p(x_i))}{1-p(x_i)} x_{ij} =$$

$$= \sum_{i=1}^n x_{ij} (y_i - p(x_i))$$

$$\nabla \log \text{VERO}(\beta) = \begin{pmatrix} \frac{\partial}{\partial \beta_0} LV \\ \vdots \\ \frac{\partial}{\partial \beta_k} LV \end{pmatrix} = \begin{pmatrix} 1 & \dots & 1 \\ x_{11} & x_{21} & \dots & x_{n1} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1k} & x_{2k} & \dots & x_{nk} \end{pmatrix} \begin{pmatrix} y_1 - p(x_1) \\ \vdots \\ y_n - p(x_n) \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$$

$k+1$ ecuaciones no lineales

$k+1$ incógnitas $\beta_0, \beta_1, \dots, \beta_k$

↓ numérico

$\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$ solución

3. Regla de clasificación

observación $\rightarrow a = (a_1, \dots, a_k)$

$\tilde{a} = (1, a_1, \dots, a_k)$

calculamos $\hat{p}(a) = h(\hat{\beta}^T \tilde{a}) = \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 a_1 + \dots)}}$

si $\hat{p}(a) > \frac{1}{2}$
↓
clase 1

si $\hat{p}(a) < \frac{1}{2}$
↓
clase 2

$$\hat{p}(a) > \frac{1}{2} \iff \frac{1}{1 + e^{-\hat{\beta}^T \tilde{a}}} > \frac{1}{2} \iff \boxed{\hat{\beta}^T \tilde{a} > 0}$$

obs: ¿Distribución de los estimadores?

Si n grande, como son estimadores de máx-vero,

$$\hat{\beta} \simeq N(\beta, \text{cov}(\hat{\beta})) \rightarrow (X^T \hat{W} X)^{-1} \rightarrow \begin{pmatrix} \hat{p}(x_1)(1-\hat{p}(x_1)) & 0 \\ 0 & \hat{p}(x_n)(1-\hat{p}(x_n)) \end{pmatrix}$$

Esto daña un intervalo de confianza (aprox.)

$$IC_{1-\alpha}(\beta_j) = \hat{\beta}_j \pm z_{\alpha/2} \cdot \text{SE}(\hat{\beta}_j) \rightarrow \begin{matrix} \text{standard error} \\ \text{raíz cuadrada del elemento} \\ j+1 \text{ de la diag. de } (X^T W X)^{-1} \end{matrix}$$

O bien, contrastar $H_0: \beta_j = 0$

Usamos que $\frac{\hat{\beta}_j}{SE(\hat{\beta}_j)} \approx N(0,1)$

Wald \rightarrow

REGLAS "ÓPTIMAS" DE CLASIFICACIÓN

Planteamiento: dos poblaciones $\begin{matrix} \pi_0 \\ \pi_1 \end{matrix}$

Queremos clasificar objetos (en π_0 o en π_1) en función de unos cuantos atributos/observaciones/medidas $X = (x_1, \dots, x_k)^T$

Este vector X se distribuye como $\begin{matrix} \rightarrow f_0(x) \text{ en } \pi_0 \\ \rightarrow f_1(x) \text{ en } \pi_1 \end{matrix}$ $\vec{x} = (x_1, \dots, x_k)^T$
(razonablemente "diferentes")

Además, tenemos probabilidades "a priori" $p_0, p_1: p_0 + p_1 = 1$

Una regla de clasificación $g: g: \Omega \subset \mathbb{R}^k \rightarrow \{0,1\}$
 $\vec{x} \mapsto g(\vec{x}) = \begin{cases} 0 \\ 1 \end{cases}$

Una tal regla divide Ω en: $\Omega = R_0^{(g)} \cup R_1^{(g)}$

Consideramos: $P(X \in R_0^{(g)} | \pi_1) = \int_{R_0^{(g)}} f_1(\vec{x}) d\vec{x}$

$P(X \in R_1^{(g)} | \pi_0) = \int_{R_1^{(g)}} f_0(\vec{x}) d\vec{x}$

$P(\text{mala clasificación con regla } g) = p_1 \int_{R_0^{(g)}} f_1(\vec{x}) d\vec{x} + p_0 \int_{R_1^{(g)}} f_0(\vec{x}) d\vec{x}$

Objetivo: hallar g que minimice esta probabilidad

Observación: $1 = \int_{\Omega} f_1(\vec{x}) d\vec{x} = \int_{R_0(g)} f_1(\vec{x}) d\vec{x} + \int_{R_1(g)} f_1(\vec{x}) d\vec{x}$

$$\Rightarrow P(\text{mala clasific. regla } g) = P_1 \left(1 - \int_{R_1(g)} f_1(\vec{x}) d\vec{x} \right) + P_0 \int_{R_1(g)} f_0(\vec{x}) d\vec{x} =$$

$$= P_1 + \int_{R_1(g)} (P_0 f_0(\vec{x}) - P_1 f_1(\vec{x})) d\vec{x}$$

Esta probabilidad es mínima cuando

$$R_1 = \left\{ x \in \Omega : P_0 f_0(\vec{x}) - P_1 f_1(\vec{x}) \leq 0 \right\} =$$

$$= \left\{ x \in \Omega : \frac{f_1(\vec{x})}{f_0(\vec{x})} \geq \frac{P_0}{P_1} \right\}$$

R_0 es la contraria (con \leq aquí)

Nota: Argumento bayesiano

Observamos un $\vec{x} = (x_1 \dots x_n)$

clasificamos \vec{x} en función de cual sea la mayor probabilidad a posteriori

$$P(\pi_1 | \vec{x}) = \underbrace{P(\vec{x} | \pi_1)}_{f_1(\vec{x})} \cdot \frac{P_1}{P(\vec{x})} = \frac{P_1 f_1(\vec{x})}{P(\vec{x})}$$

$$P(\pi_0 | \vec{x}) = P(\vec{x} | \pi_0) \cdot \frac{P_0}{P(\vec{x})} = \frac{P_0 f_0(\vec{x})}{P(\vec{x})}$$

comparamos
estos dos
valores

Nota 2 : Incorporamos costes de mala clasificación

| | π_0 | π_1 |
|---------|----------|----------|
| π_0 | 0 | C_{10} |
| π_1 | C_{01} | 0 |

Interesa minimizar coste medio de mala clasificación =

$$= P_1 C_{01} \int_{R_0^{(g)}} f_1(\vec{x}) d\vec{x} + P_0 C_{10} \int_{R_1^{(g)}} f_0(\vec{x}) d\vec{x}$$

Sea g óptima:

$$R_1^{(g)} = \left\{ \vec{x} \in \Omega \mid P_0 C_{10} f_0(\vec{x}) - P_1 C_{01} f_1(\vec{x}) \leq 0 \right\} =$$

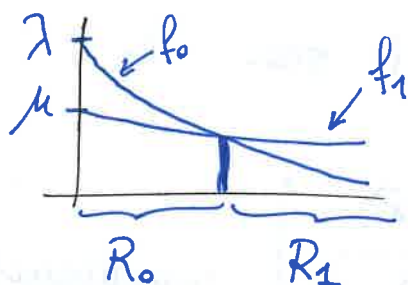
$$= \left\{ \vec{x} \in \Omega \mid \frac{f_1(\vec{x})}{f_0(\vec{x})} \geq \frac{C_{10}}{C_{01}} \cdot \frac{P_0}{P_1} \right\}$$

Ejemplo: $K=1$, $\Omega = \mathbb{R}^+$, x variable

$$\begin{cases} f_0(x) = \lambda e^{-\lambda x}, & x > 0 \\ f_1(x) = \mu e^{-\mu x}, & x > 0 \end{cases}$$

$$\lambda > \mu \quad E(x) = \frac{1}{\lambda}$$

$$P_0 = P_1 = \frac{1}{2}$$



$$f_0(x) = f_1(x)$$

$$\lambda e^{-\lambda x} = \mu e^{-\mu x} \Rightarrow x = \frac{1}{\lambda - \mu} \ln\left(\frac{\lambda}{\mu}\right)$$

$$\text{Si } P_0 = \frac{1}{10}, \quad P_1 = \frac{9}{10}$$

$$\Rightarrow \frac{f_1(x)}{f_0(x)} \geq \frac{1/10}{9/10} = \frac{1}{9} \quad \dots$$

El caso de dos poblaciones normales

$$x = (x_1, \dots, x_k) \quad f_0(x) = N_k(\mu_0, \Sigma_0) \quad , \quad P_0 \quad (\text{quizás coste } C_{01})$$

$$f_1(x) = N_k(\mu_1, \Sigma_1) \quad , \quad P_1 \quad ("C_{10"})$$

$$f_i(x) = \frac{1}{(2\pi)^{k/2} |\Sigma_i|^{1/2}} \cdot \exp\left(\frac{-1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i)\right) \quad i=0,1$$

$$R_1(x) = \left\{ x \in \mathbb{R}^k : \frac{f_1(x)}{f_0(x)} \geq \frac{P_0}{P_1} \right\}$$

$$\frac{P_0}{P_1} \leq \frac{|\Sigma_0|^{1/2}}{|\Sigma_1|^{1/2}} \exp\left\{\frac{-1}{2} \left((x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) - (x - \mu_0)^T \Sigma_0^{-1} (x - \mu_0) \right)\right\}$$

$$\ln\left(\frac{P_0}{P_1}\right) \leq \frac{-1}{2} \ln\left(\frac{|\Sigma_1|}{|\Sigma_0|}\right) - \frac{1}{2} \left[x^T (\Sigma_1^{-1} - \Sigma_0^{-1}) x - 2(\mu_1^T \Sigma_1^{-1} - \mu_0^T \Sigma_0^{-1}) x + \mu_1^T \Sigma_1^{-1} \mu_1 - \mu_0^T \Sigma_0^{-1} \mu_0 \right]$$

$$\boxed{\ln\left(\frac{P_0}{P_1}\right) \leq \underbrace{\frac{-1}{2} x^T (\Sigma_1^{-1} - \Sigma_0^{-1}) x}_{\text{forma cuadrática}} + \underbrace{(\mu_1^T \Sigma_1^{-1} - \mu_0^T \Sigma_0^{-1}) x}_{\text{lineal}} - \underbrace{\frac{1}{2} \left[\ln\left(\frac{|\Sigma_1|}{|\Sigma_0|}\right) + \mu_1^T \Sigma_1^{-1} \mu_1 - \mu_0^T \Sigma_0^{-1} \mu_0 \right]}_{\text{constante}}}$$

Caso particular $\Sigma_1 = \Sigma_0 = \Sigma$

$$\ln\left(\frac{P_0}{P_1}\right) \leq (\mu_1 - \mu_0)^T \Sigma^{-1} x + \frac{1}{2} (\mu_1 - \mu_0)^T \Sigma^{-1} (\mu_1 + \mu_0)$$

Si $P_0 = P_1 \longrightarrow$ Fisher

