

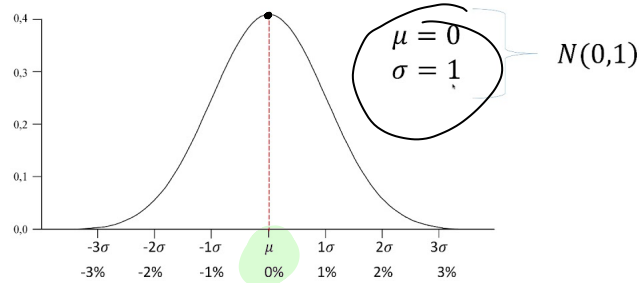
REPASO ESTADÍSTICA - INFERENCIAL

1

DISTRIBUCIONES PROBABILIDAD - DAN LA PROBABILIDAD DE OCURRENCIA DE EVENTOS PARA UN EXPERIMENTO

NORMAL -
CONTINUA

La Distribución Normal Estandarizada



✎

POISSON - DISCRETA - NÚMERO DE OBSERVACIONES EN UN PERIODO DE TIEMPO

BERNOULLI - LANZAR UNA MONEDA - EJEMPLO TIPILO

$$X \in \{0, 1\}$$

$$f(x) = P(X=w) = p^w (1-p)^{1-w}$$

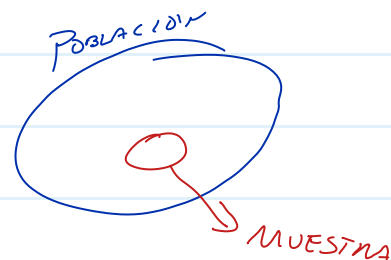
BINOMIAL - EN UN EXPERIMENTO BERNOULLI QUE REPETIMOS n VECES CUÁNTAS VECES OBTENGA UN ÉXITO

EXPONENCIAL - TIEMPO QUE PASA ANTES DE QUE OCURRA UN EVENTO EN UN PROCESO POISSON.

¿OS INTERESA SABER LA EDAD PROMEDIO DE ALUMNOS DE MAESTRÍAS MATEMÁTICAS EN MÉXICO

POBLACIÓN - TODOS ESOS ALUMNOS

MUESTRA - SUBCONJUNTO DE ESA POBLACIÓN



(a) X_1, X_2, \dots, X_n VARIABLES ALEATORIAS n : # DATOS MUESTRA

(b) INDEPENDENCIA $\parallel X_i$ NO AFECTA AL VALOR DE X_j

$$P(A|B) = P(A) \parallel$$

(c) PROVIENEN DE UNA DISTRIBUCIÓN IDENTICA

$$X_i \sim N(\mu, \sigma^2)$$

↓
DISTRIBUYEN
COMO

μ : MEDIA POBLACIONAL

σ^2 : VARIANZA POBLACIONAL

3 MANERAS DE INFERIRLOS:

→ ESTIMACIÓN PUNTUAL: CON UNA FUNCIÓN DE LA MUESTRA QUIERO ESTIMAR ESOS PARÁMETROS

(a) MEDIA MUESTRAL

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{\sum_{i=1}^n x_i}{n}$$

- ESPERANZA

$$E[\bar{x}] = \mu$$

PROPIEDAD DE
INSESGAMIENTO

- VARIANZA

$$\text{Var}[\bar{x}] = \sigma^2/n$$

DISTRIBUCIÓN

$$\bar{x} \sim N(\mu, \sigma^2/n)$$

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} \Rightarrow \underbrace{n\bar{x}}_{n\text{-COMPONENTES}} - \underbrace{\sum_{i=1}^n x_i}_{n\text{-COMPONENTES}} = 0$$

$$\Rightarrow n\bar{x} - [x_1 + x_2 + \dots + x_n] = 0$$

$$\Rightarrow (\bar{x} - x_1) + (\bar{x} - x_2) + \dots + (\bar{x} - x_n) = 0$$

Si yo conozco \bar{x} y $(n-1)$ DATOS

YA PUEDO CONOCER EL ÚLTIMO VALOR FALTANTE

$n-1$: GRADOS DE LIBERTAD

(b) VARIANZA MUESTRAL

$$s^2 = \frac{(x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n-1} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} = \frac{s_{xx}}{n-1}$$

$$\cancel{E[s^2] = \sigma^2} \quad \text{INSES GADA}$$

(c) DESVIACIÓN ESTÁNDAR MUESTRAL

$$s = \sqrt{s^2}$$

→ ESTIMACIÓN POR INTERVALOS

DEPENDIÉNDOLA DE LA MUESTRA. PARA CONSTRUIRLOS SE UTILIZAN CANTIDADES PIVOTALES

→ EN SU FÓRMULA
INCLUYEN A LOS
PARÁMETROS POBLACIONALES

PERO

SU DISTRIBUCIÓN DE
PROBABILIDAD NO
LOS INCLUYE

ESTANDARIZACIÓN : $X \sim N(\mu, \sigma^2)$

Paso 1: RESTAR LA MEDIA

$$X - \mu$$

Paso 2: DIVIDIRLO ENTRE DESV

$$\frac{X - \mu}{\sigma} \sim N(0, 1)$$

INTERVALOS DE CONFIANZA

→ ENCONTRAR μ CONOCIENDO σ^2

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

ESTANDARIZAR

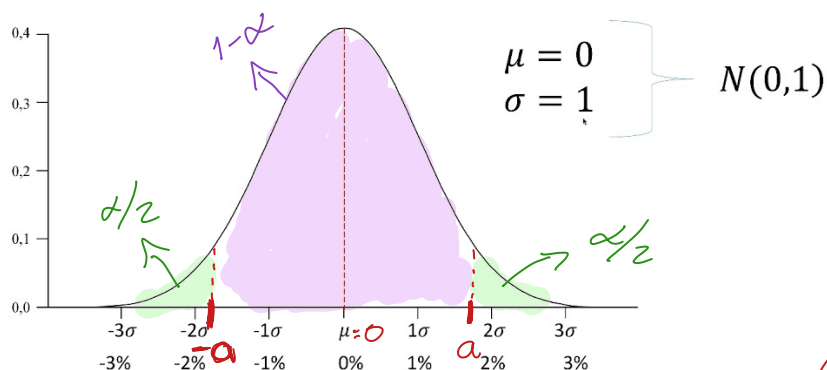
$$\frac{\bar{X} - \mu}{\sqrt{\frac{\sigma^2}{n}}} = \sqrt{n} \left(\frac{\bar{X} - \mu}{\sigma} \right) \sim N(0, 1)$$

$$P(-a \leq \sqrt{n} \left(\frac{\bar{X} - \mu}{\sigma} \right) \leq a) = 1 - \alpha$$

→ NIVEL DE CONFIANZA
Ejemplo: 95%

→ NIVEL DE SIGNIFICANCIA
Ejemplo: 5%

La Distribución Normal Estandarizada



NOTACIÓN

$$a = z_{(\alpha/2)}$$

CUANTIL DE UNA NORMAL
ESTÁNDAR EN $\alpha/2$
LA PROBABILIDAD DE SER MAYOR
ES DE $\alpha/2$

INTERVALO DE CONFIANZA

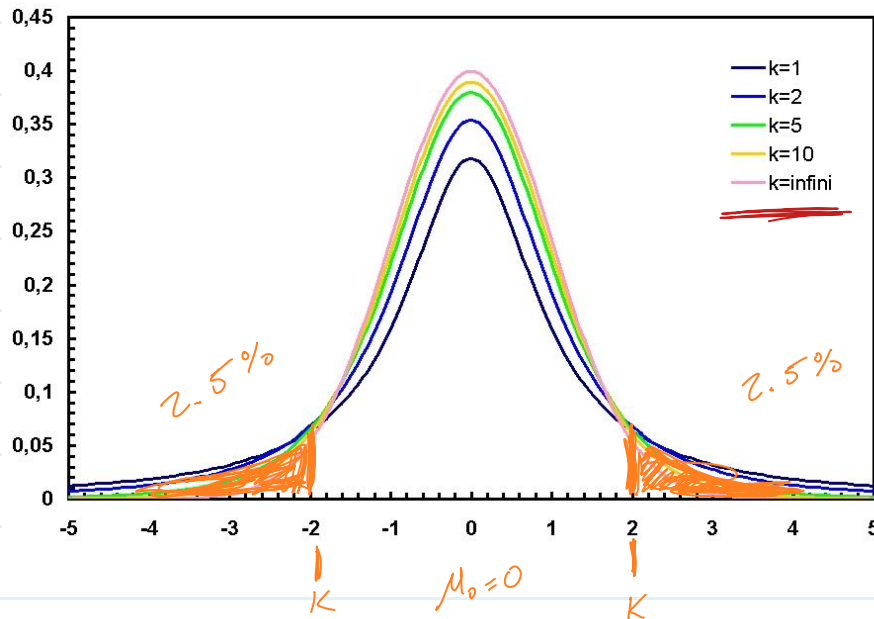
$$\mu \in \left[\bar{x} - \frac{\sigma}{\sqrt{n}} z^{(\alpha/2)}, \bar{x} + \frac{\sigma}{\sqrt{n}} z^{(\alpha/2)} \right]$$

// SI REALIZAMOS EL EXPERIMENTO 100 VECES Y $1-\alpha$ ES 95%
VAMOS A ENCONTRAR A μ 95 DE ESAS 100 VECES EN EL INTERVALO //

// SI $1-\alpha$ AUMENTA \rightarrow EL INTERVALO AUMENTA //

Caso 2: CON σ^2 DESCONOCIDA

$$\sqrt{n} \left(\frac{\bar{x} - \mu}{s} \right) \sim t_{(n-1)} \rightarrow \text{T-STUDENT CON } n-1 \text{ GRADOS DE LIBERTAD}$$



CONFORME n TIENDA

A INFINITO

$$t_{(n-1)} \sim N(0, 1)$$

$$n \rightarrow \infty$$

INTERVALO:

$$\left[\bar{x} - \frac{s}{\sqrt{n}} t_{(n-1)}^{(\alpha/2)}, \bar{x} + \frac{s}{\sqrt{n}} t_{(n-1)}^{(\alpha/2)} \right]$$

→ PRUEBAS DE HIPÓTESIS

CONJETURA RESPECTO A LOS PARÁMETROS POBLACIONALES

CONJETURA INICIAL

HIPÓTESIS NULA (H_0)

CONJETURA ALTERNATIVA

HIPÓTESIS ALTERNATIVA (H_1)

POS COLAS

COLA IZQUIERDA

COLA DERECHA

$$H_0: \mu = \mu_0$$

$$H_0: \mu > \mu_0$$

$$H_0: \mu < \mu_0$$

YA ES
ESPECÍFICO

$$H_a: \mu \neq \mu_0$$

$$H_a: \mu \leq \mu_0$$

$$H_a: \mu \geq \mu_0$$

DEBEMOS CREAR UNA REGLA DE DECISIÓN PARA SABER SI RECHAZAMOS NUESTRA H_0

	H_0 VERDADERA	H_0 FALSA
RECHAZAR H_0	Error TIPO I (α)	SIN ERROR ✓
NO RECHAZAR H_0	SIN ERROR ✓	Error TIPO II (β)

SE CONTROLA EL ERROR TIPO I (α)

$$\alpha = P(\text{RECHAZAR } H_0 \mid H_0 \text{ CIERTA})$$

¿SI RECHAZAMOS O NO? CONSTRUIR REGIÓN DE RECHAZO

UTILIZAMOS UN ESTADÍSTICO DE PRUEBA

$$T = \sqrt{n} \left(\frac{\bar{x} - \mu_0}{s} \right) \sim t_{(n-1)}$$

↓
SUPONIENDO H_0 VERDADERA

S, T ES GRANDE (POSITIVO O NEGATIVO) NOS DA EVIDENCIA PARA RECHAZAR H_0

$$\begin{matrix} 5\% \\ 2\% \\ 1\% \end{matrix} \leftarrow \alpha = P(|T| > K \mid H_0 \text{ VERDADERA})$$

RECHAZAR H_0

$$K = t_{(n-1)}^{\alpha/2}$$

REGIÓN DE RECHAZO: $H_0: \mu = \mu_0$

$S, |T| > t_{(n-1)}^{\alpha/2}$ HAY EVIDENCIA PARA RECHAZAR H_0 A UN NIVEL DE SIGNIFICANCIA α

REGRESIÓN LINEAL SIMPLE

$x \rightarrow$ PREDICTOR (VARIABLE INDEPENDIENTE) // ESTATURA // $x_1, x_2, x_3 \dots$
 $y \rightarrow$ RESPUESTA (VARIABLE DEPENDIENTE) // PESO // $y_1, y_2, y_3 \dots$

ANÁLISIS INDIVIDUAL

MEIA MUESTRA

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

$$\bar{y} = \frac{\sum_{i=1}^n y_i}{n}$$

VARIANZA MUESTRA

$$s_x^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

$$s_y^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}$$

DESVIACIÓN ESTÁNDAR

$$s_x = \sqrt{s_x^2}$$

$$s_y = \sqrt{s_y^2}$$

COVARIANZA MUESTRAL

$$\hat{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1} = \frac{S_{xy}}{n-1} = \hat{Cov}(y, x)$$

Si $\hat{Cov}(x, y) > 0 \Rightarrow$ POSITIVA PROPORCIONAL

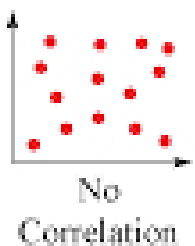
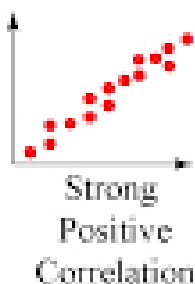
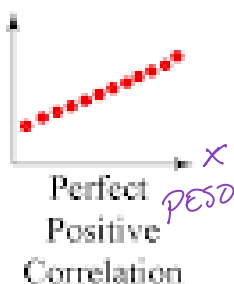
Si $\hat{Cov}(x, y) < 0 \Rightarrow$ NEGATIVA PROPORCIONAL

COEFICIENTE DE CORRELACIÓN

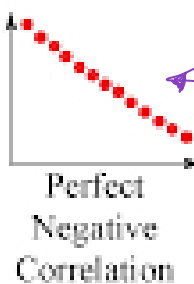
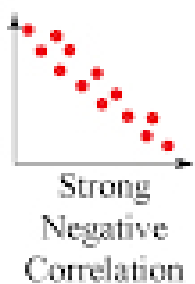
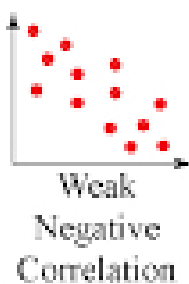
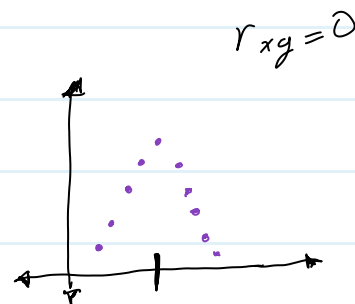
$$r_{xy} = \frac{\hat{Cov}(x, y)}{S_x S_y}$$

$$-1 \leq r_{xy} \leq 1$$

ESTATURA y
 $r_{xy} = 1$



$r_{xy} = 0$



$r_{xy} = -1$

PODEMOS PROPONER UNA RELACIÓN PARAMÉTRICA

$$y = \beta_0 + \beta_1 x + \varepsilon \rightarrow \text{ERRORES ESTOCÁSTICO} \\ \text{// ALEATORIO //$$

↓
COEFICIENTE

CONSIDERA 5 SUPUESTOS:

(1) LA RELACIÓN VERDADERA ES

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon$$

CONOCIDO

↑ ALEATORIEDAD

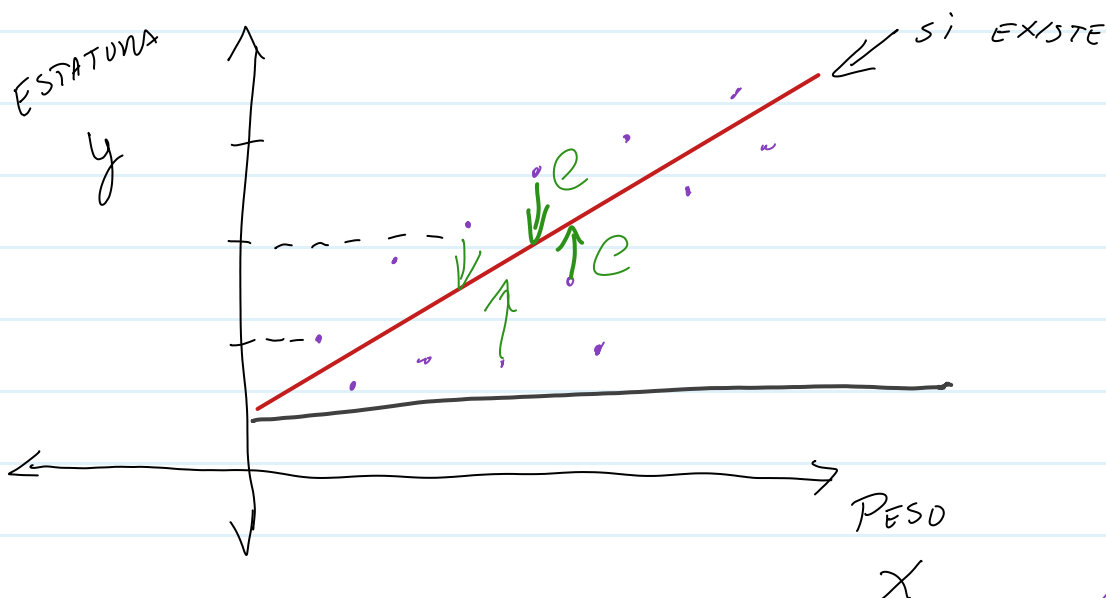
(2) x NO ES ESTOCÁSTICA // NO TIENE UN COMPORTAMIENTO ALEATORIO //

(3) ε_i SON VARIABLES ALEATORIAS : $E[\varepsilon_i] = 0$

$$\text{Var}[\varepsilon_i] = \sigma^2$$

(4) ε_i SON INDEPENDIENTES E IDÉNTICAMENTE DISTRIBUIDOS (i.i.d)

(5) ε_i SON NORMALES $\varepsilon_i \sim N(0, \sigma^2)$



ESTIMACIÓN POR MÍNIMOS CUADRADOS:

⇒ PROPONEMOS $\hat{\beta}_0, \hat{\beta}_1$ Y CALCULAR

$$\min_{\hat{\beta}_0, \hat{\beta}_1 \in \mathbb{R}} \sum_{i=1}^n (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2$$

RESIDUALES ε_i

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hat{\beta}_1 = \frac{r_{xy} \cdot s_y}{s_x}$$

ESTIMACIÓN POR MÁXIMA VEROSIMILITUD

Y A CONOCER LA DISTRIBUCIÓN Y VAMOS A BUSCAR PARÁMETROS QUE LLEVEN A LA DISTRIBUCIÓN MÁS PARECIDA A NUESTROS DATOS

VEROSIMILITUD:

$$l(\beta_0, \beta_1, \sigma^2) = f(y_1, y_2, \dots, y_n) = f(y_1) f(y_2) \dots f(y_n)$$

↙
FUNCIÓN DE
DENSIDAD

INDEPENDENCIA

$$= \prod_{i=1}^n f(y_i) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\left(\frac{y_i - \mu_i}{\sigma^2}\right)^2}$$

NO SON IDÉNTICAMENTE
DISTRIBUIDAS

LOG VEROSIMILITUD

$$L = \ln(l(\beta_0, \beta_1, \sigma^2))$$

¿QUÉ PARÁMETROS APROXIMA?

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hat{\sigma}_{MV}^2 = \frac{\sum_{i=1}^n \left(y_i - [\hat{\beta}_0 + \hat{\beta}_1 x_i] \right)^2}{n}$$

e_i
RESIDUALES

$\hat{\sigma}_{MV}^2$ NO ES INSESGADA $E[\hat{\sigma}_{MV}^2] \neq \sigma^2$

\Rightarrow RESTRICCIONES DE LOS RESIDUALES

$$\left. \begin{aligned} \sum_{i=1}^n e_i &= 0 \\ \sum_{i=1}^n x_i e_i &= 0 \end{aligned} \right\} \begin{aligned} &2 \text{ RESTRICCIONES} \\ &\text{SI CONOZCO } n-2 \text{ RESIDUALES} \\ &\text{PUEDO CONOCER LOS } 2 \text{ QUE FALTAN} \end{aligned}$$

$n-2$: GRADOS DE LIBERTAD DE RESIDUALES

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n e_i^2}{n-2} = S^2 \rightarrow S^2 \text{ ES INSESGADO}$$

$$E[\hat{\sigma}^2] = \sigma^2$$

VARIANZAS

$$\begin{aligned} \text{ESTIMADOS} \quad \text{Var}(\hat{\beta}_1) &= \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \rightarrow \hat{\text{Var}}(\hat{\beta}_1) = \frac{\hat{\sigma}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ \text{Var}(\hat{\beta}_0) &= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) \rightarrow \hat{\text{Var}}(\hat{\beta}_0) = \hat{\sigma}^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) \\ \text{Var}(\hat{\sigma}^2) &= \frac{2\sigma^4}{n-2} \rightarrow \hat{\text{Var}}(\hat{\sigma}^2) = \frac{2\hat{\sigma}^4}{n-2} \end{aligned}$$

INTERVALOS

$$\beta_j \in \left[\hat{\beta}_j - S_e(\hat{\beta}_j) \cdot t_{(n-2)}^{(\alpha/2)}, \hat{\beta}_j + S_e(\hat{\beta}_j) \cdot t_{(n-2)}^{(\alpha/2)} \right]$$

PRUEBA DE HIPOTESIS

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon$$

$\varepsilon = 0$

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

REGRESIÓN LINEAL MÚLTIPLE

$$Y_i = \beta_0 + \beta_1 X_i + \dots + \beta_n X_n + \epsilon_i \rightarrow \text{MINIMIZAR}$$

β_0 → PESO
 $\beta_1 X_i$ → ESTATURA
 $\beta_n X_n$ → EDAD

SUBAJUSTE → MENOS VARIABLES

SOBREAJUSTE → MÁS VARIABLES

→ FLEXIBLE

→ NO LOGRE EXPLICAR Y

→ NO AÑADIR PREDICTORES INSIGNIFICANTES

→ SUFICIENTEMENTE SENCILLO

MÁS VARIABLES → R^2 AUMENTA O SE QUEDA IGUAL

R^2 AJUSTADA → ENTRE MÁS GRANDE MEJOR

$$R^2_a = 1 - \frac{(n-1)}{n-k-1} (1 - R^2)$$

R^2 : COEF. DETERMINACIÓN
 n : # DATOS

k : # PREDICTORES

CRITERIOS INFORMACIÓN

$$AIC \text{ (AKAIKE)} = 2(k+1) - 2 \ln(L)$$

$$BIC \text{ (BAYESIANO)} = \ln(L)(k+1) - 2 \ln(L)$$

ENTRE MÁS PEQUEÑO MEJOR

MÁS RESTRICTIVO

MENOS RESTRICTIVO

POCOS DATOS

