

# FUNDAMENTOS Y HERRAMIENTAS PARA LA MODELIZACIÓN DE PROCESOS TÉCNICOS-CIENTÍFICOS DE INVESTIGACIÓN

Bloque III: Modelización estadística.

**Regresión y correlación**



Roberto Espejo Mohedano e-mail: [malesmor@uco.es](mailto:malesmor@uco.es)

## Modelización.

- Necesidad del estudio simultaneo de varias variables aleatorias.
  - Analizar las posibles relaciones entre ellas. (Estudios de correlación).
  - Intentar establecer el modelo matemático que las relacione. (Estudios de regresión o ajuste).
- Interés por estimar una magnitud  $Y$  (variable dependiente) en función de una o varias  $X_1, X_2, \dots, X_k$  variables explicativas (variables independientes).
  - Imposibilidad de predicción exacta. Perturbación aleatoria  $\varepsilon$ .
  - Modelos lineales:
$$\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k$$
$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \varepsilon$$

## Modelos comunes

### Lineales

Lineal simple:  $\hat{Y} = \beta_0 + \beta_1 X$

Parabólico:  $\hat{Y} = \beta_0 + \beta_2 X + \beta_1 X^2$

Cúbico:  $\hat{Y} = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3$

Polinómico:  $\hat{Y} = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^2 + \dots + \beta_h X^h$

### No Lineales

Exponencial:  $\hat{Y} = \beta_0 k^{\beta_1 X}$

Potencial:  $\hat{Y} = \beta_0 X^{\beta_1}$

Hiperbólico:  $\hat{Y} = 1 / (\beta_0 + \beta_1 X)$

Logístico:  $\hat{Y} = 1 / (e^{-\beta_0 - \beta_1 X})$

Estimación de los coeficientes:  $\beta_0, \beta_1, \beta_2, \dots, \beta_k$

Problema de inferencia:

Contrastes sobre los coeficientes y sobre el ajuste del modelo.

## Correlación simple.

$$R_{xy} = \frac{S_{xy}}{S_x S_y} \in (-1;1)$$

$$R_{x_j x_i} = \frac{S_{x_j x_i}}{S_{x_j} S_{x_i}} \quad \forall i \neq j = 1, 2, \dots$$

Matriz de Correlación:  $\Gamma_{yx_1 x_2 \dots} = \begin{pmatrix} 1 & R_{yx_1} & R_{yx_2} & \dots \\ R_{yx_1} & 1 & R_{x_1 x_2} & \dots \\ R_{yx_2} & R_{x_1 x_2} & 1 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$

# Contrastes de incorrelación

Poblaciones  $X$  e  $Y$  numéricas.

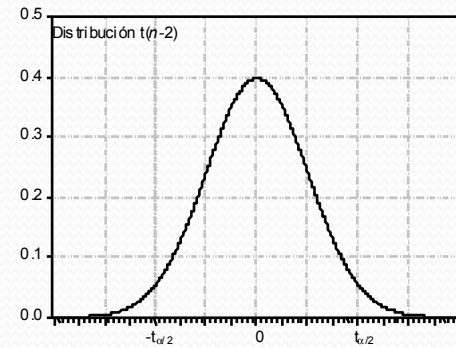
$$\text{Hipótesis: } \begin{cases} H_0: X \text{ está incorrelada con } Y \\ H_1: X \text{ está correlada con } Y \end{cases} \Rightarrow \begin{cases} H_0: \rho_{xy} = 0 \\ H_1: \rho_{xy} \neq 0 \end{cases}$$

$$\text{Estadístico: } T = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} \in t_{(n-2)} \quad r = r_{xy} = \frac{S_{xy}}{S_x S_y}$$

$$\text{Regla de decisión de nivel } \alpha: C_0 = (-t_{\alpha/2}, +t_{\alpha/2})$$

**Decisión:** A nivel  $\alpha$ : Si  $t \in C_0 \rightarrow$  se acepta  $H_0$   
Para todo  $\alpha \leq p \rightarrow$  se acepta  $H_0$   
Para todo  $\alpha > p \rightarrow$  se rechaza  $H_0$

$$\text{Probabilidad límite: } p = P(T > |t|)$$



## Contrastes de independencia

Poblaciones no numéricas.

$$X \in S_x = \{x_1^*, \dots, x_f^*\}$$

$$Y \in S_y = \{y_1^*, \dots, y_c^*\}$$

Hipótesis:  $\begin{cases} H_0 : X \text{ es independiente de } Y \\ H_1 : X \text{ esta relacionada con } Y \end{cases}$

Muestra:  $(x_t, y_t) \quad t = 1, \dots, n$  tabulada en forma de tabla de contingencia

$$\|n_{ij}\| \quad i = 1 \dots f; j = 1 \dots c$$

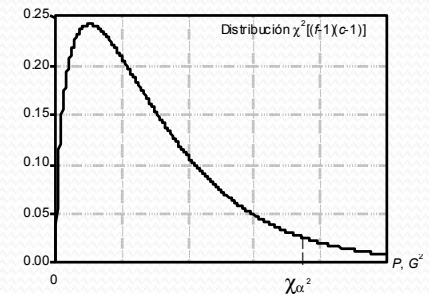
Estadísticos:  $G^2 = 2 \sum_{i=1}^f \sum_{j=1}^c n_{ij} \ln \frac{n_{ij}}{E_{ij}}$

$$P = \sum_{i=1}^f \sum_{j=1}^c \frac{(n_{ij} - E_{ij})^2}{E_{ij}} \quad P_c = \sum_{i=1}^f \sum_{j=1}^c \frac{(|n_{ij} - E_{ij}| - 0.5)^2}{E_{ij}}$$

Distribución muestral (asintótica):

$$G^2 / H_0, P / H_0, P_c / H_0 \in \chi^2((f-1), (c-1))$$

Condiciones asintóticas:  $E_{ij} = \frac{n_{i.} \cdot n_{.j}}{n} \geq 1 \quad \forall i, j$



Regla de decisión de nivel  $\alpha$ :

$$C_0 = (0, \chi_\alpha^2)$$

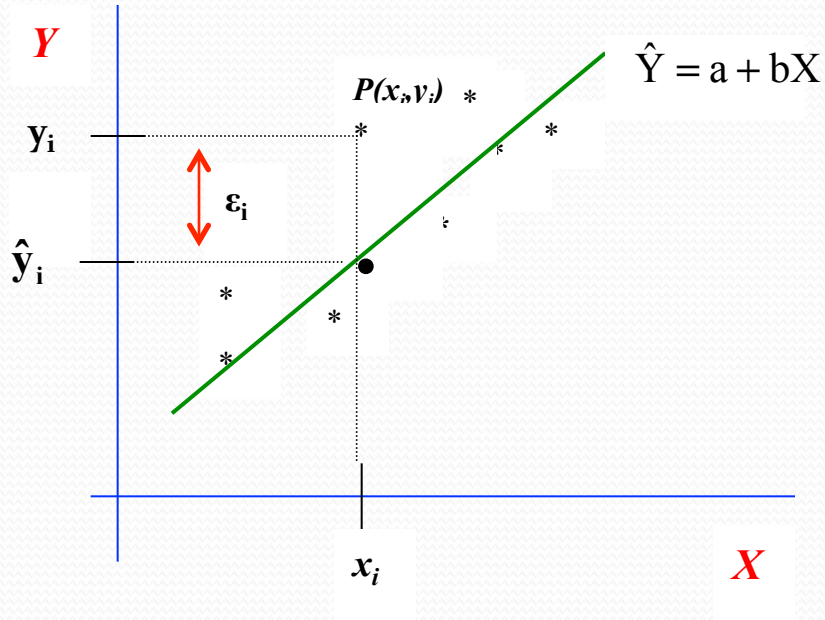
**Decisión:** A nivel  $\alpha$ : Si  $P, P_c, G^2 \in C_0 \rightarrow$  se acepta  $H_0$

Para todo  $\alpha \leq p \rightarrow$  se acepta  $H_0$

Para todo  $\alpha > p \rightarrow$  se rechaza  $H_0$

# Regresión Simple: Línea de Regresión

Recta de regresión de Y sobre X



$\hat{y}_i$  valor estimado de Y para  $X=x_i$

$\epsilon_i$  residuos:  $\epsilon_i = y_i - \hat{y}_i$

**Mínimos cuadrados**

$$H = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - a - bx_i)^2$$

$$b = \frac{S_{XY}}{S_X^2}$$

$$a = \bar{Y} - b\bar{X}$$

Los modelos de Regresión requieren Normalidad en la variable Dependiente

**Regresión Múltiple:**  $\hat{Y} = a + b_1 X_1 + b_2 X_2 + \dots + b_n X_n$

Matriz de Covarianzas

$$\Sigma_{y x_1 x_2 \dots} = \begin{pmatrix} S_y^2 & S_{yx_1} & S_{yx_2} & \dots \\ S_{x_1 y} & S_{x_1}^2 & S_{x_1 x_2} & \dots \\ S_{x_2 y} & S_{x_1 x_2} & S_{x_2}^2 & \dots \\ \dots & \dots & \dots & \dots \end{pmatrix}$$

Vector de Medias

$$\vec{\mu}_{y x_1 x_2 \dots} = \begin{pmatrix} \bar{Y} \\ \bar{X}_1 \\ \bar{X}_2 \\ \dots \end{pmatrix}$$

$$b_i = -\frac{A_{1,i+1}}{A_{1,1}} \quad \forall \quad i = 1, 2, \dots$$

$$a = \bar{Y} - b_1 \bar{X}_1 - b_2 \bar{X}_2 - \dots \quad \left( \bar{Y} = a + b_1 \bar{X}_1 + b_2 \bar{X}_2 + \dots \right)$$



**Correlación múltiple:**  $R_{yx_1x_2\ldots} = R_y = \frac{S_{y\hat{y}}}{S_y S_{\hat{y}}} \in [-1,1]$

Se verifica que:  $R_y^2 \geq R_{yx_i}^2 \quad \forall i = 1, 2, \dots$

**Coeficiente de determinación:**  $R^2 = R_{xy}^2 = \left( \frac{S_{xy}}{S_x S_y} \right)^2 = 1 - \frac{S_{\epsilon}^2}{S_y^2} \in [0; 1]$

"Tanto por 1 (ó %) de la varianza de  $Y$  explicada por el modelo."

Nos informa de hasta que punto el modelo se ajusta a los datos.

## Test de hipótesis sobre el modelo de regresión: Test F

Se basa en el **Teorema de descomposición de la varianza**  $S_y^2 = S_{\hat{y}}^2 + S_e^2$

o multiplicando por el número de datos (n), la suma de cuadrados totales se descompone en:  $S_y = S_{\hat{y}} + S_e$

Siendo:  $S_y = \sum_{i=1}^n (y_i - \bar{y})^2$  La Suma de cuadrados total

$S_{\hat{y}} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$  La Suma de cuadrados asociada a las variables explicativas (la regresión)

$S_e = \sum_{i=1}^n e_i^2$  La Suma de cuadrados residual, la no explicada por la regresión

Hipótesis a contrastar: 
$$\begin{cases} H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0 \\ H_1 : \text{algún o algunos } \beta_i \neq 0 \end{cases} \Rightarrow \begin{cases} H_0 : \text{Rechazo Goblamente el modelo} \\ H_1 : \text{Acepto Goblamente el modelo} \end{cases}$$

Estadístico del contraste:  $F = \frac{S_{\hat{y}} / k - 1}{S_e / (n - k - 1)}$  Siendo su distribución muestral es:  $F / H_0 \in F(k - 1; n - k - 1)$

Siendo k el número de variables explicativas (independientes) del modelo

**Regla de decisión** de nivel  $\alpha$ :  $C_0 = (0 ; F_\alpha)$

Para todo  $\alpha \leq p \rightarrow$  se acepta  $H_0$

Para todo  $\alpha > p \rightarrow$  se rechaza  $H_0$

## Test de hipótesis sobre el modelo de regresión: Test F

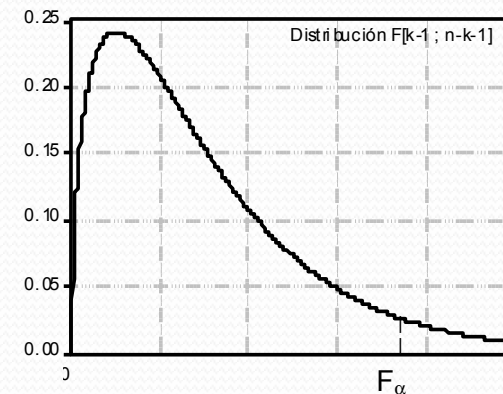
Los cálculos anteriores se resumen en la siguiente tabla de análisis de la Varianza

Fuente de Variación	Grados de Libertad	Sumas de Cuadrados	Medias de Cuadrados	Estadístico F
Regresión	k-1	$S_{\hat{y}} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$	$M_{\hat{y}} = \frac{S_{\hat{y}}}{k-1}$	$F = \frac{M_{\hat{y}}}{M_e}$
Residuo	n-k-1	$S_e = \sum_{i=1}^n e_i^2$	$M_e = \frac{S_e}{n-k-1}$	
Variación Total	n-1	$S_y = \sum_{i=1}^n (y_i - \bar{y})^2 = S_{\hat{y}} + S_e$		

**Decisión:** A nivel  $\alpha$ : Si  $F \in C_0 \rightarrow$  se acepta  $H_0$

Para todo  $\alpha \leq p \rightarrow$  se acepta  $H_0$

Para todo  $\alpha > p \rightarrow$  se rechaza  $H_0$



# Contrastes sobre los Coeficientes de Regresión

## Tests T

$$\text{Hipótesis: } \begin{cases} H_0 : \beta_j = 0 \\ H_1 : \beta_j \neq 0 \end{cases}$$

$$\text{Estadístico: } T = \frac{b_j}{\bar{s}_{b_j}} \quad \bar{s}_{b_j}^2 = \bar{s}_e^2 a_{jj} \quad \bar{s}_e^2 = \frac{1}{n-k-1} \sum e_i^2 \quad A = \|a_{ij}\| = (X'X)^{-1}$$

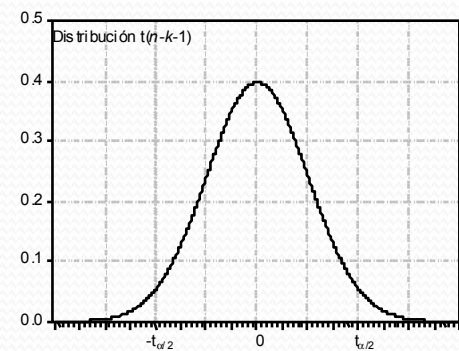
$$\text{Distribución muestral: Si } e_i \in N(0_n, \sigma_\varepsilon^2 I_n) \rightarrow T \in t(n-k-1)$$

$$\text{Regla de decisión de nivel } \alpha: C_0 = (-t_{\alpha/2}; t_{\alpha/2})$$

$$\text{Decisión: A nivel } \alpha: \text{ Si } t \in C_0 \rightarrow \text{se acepta } H_0$$

$$\text{Para todo } \alpha \leq p \rightarrow \text{se acepta } H_0$$

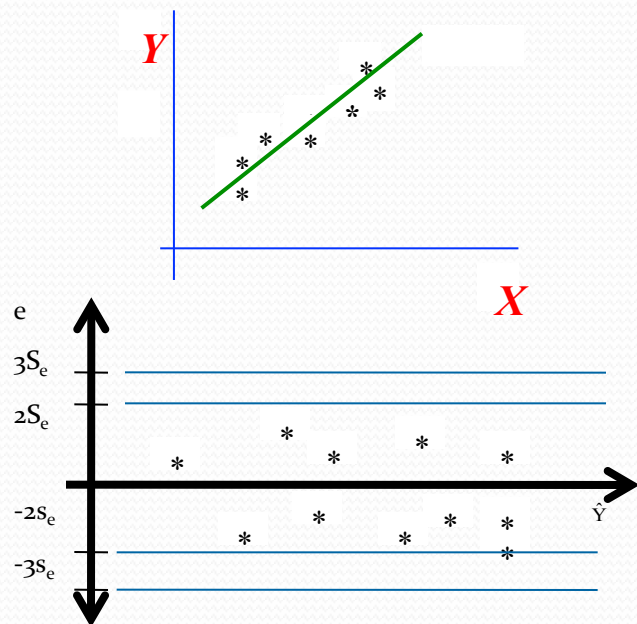
$$\text{Para todo } \alpha > p \rightarrow \text{se rechaza } H_0$$



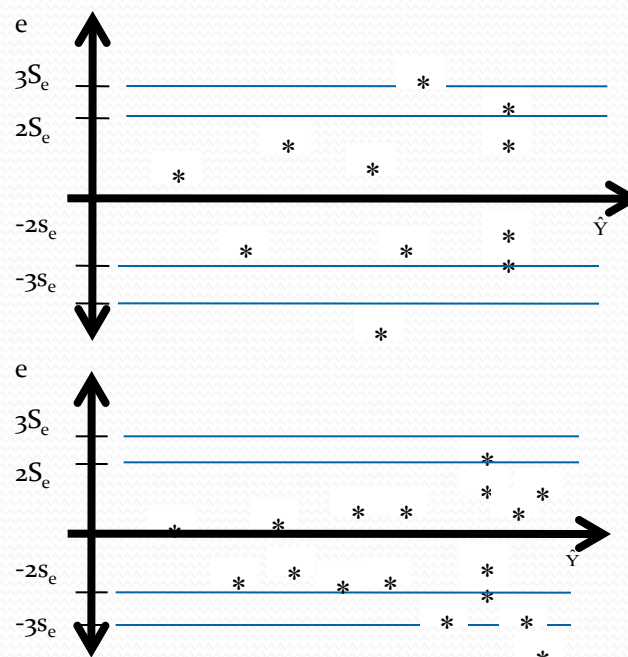
# Propiedades de los Residuos

1ª.-  $\bar{e} = 0$

2ª.-  $\text{Var}(\epsilon_i) = \text{cte}$  Varianza constante. (homocedásticos)



Homocedasticidad



Heterocedasticidad

# Propiedades de los Residuos

3ª.- Los residuos han de estar incorrelados entre si.

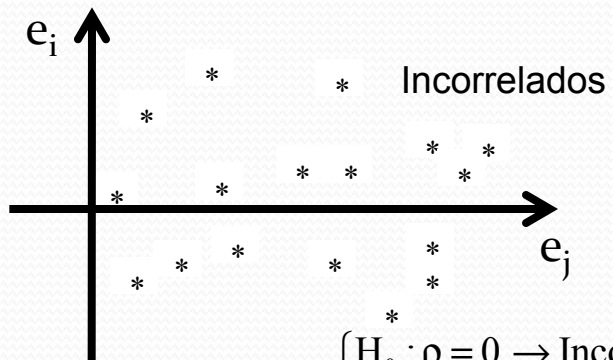
Estadístico de Durbin-Watson

$$DW = \frac{\sum_{i=2}^n (e_i - e_{i-1})^2}{\sum_{i=1}^n e_i^2} \in [0; 4]$$

DW  $\approx$  2, Incorrelados

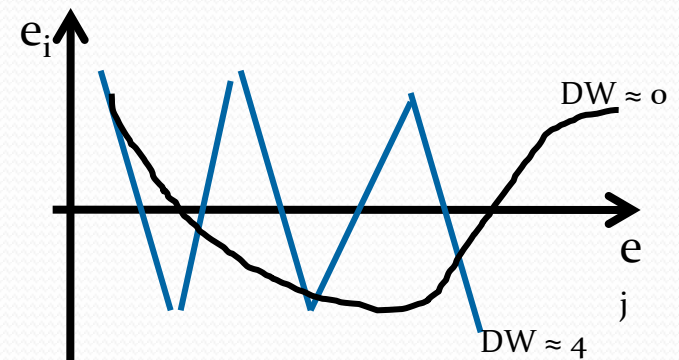
DW  $\approx$  4 Autocorrelación negativa

DW  $\approx$  0 Autocorrelación positiva



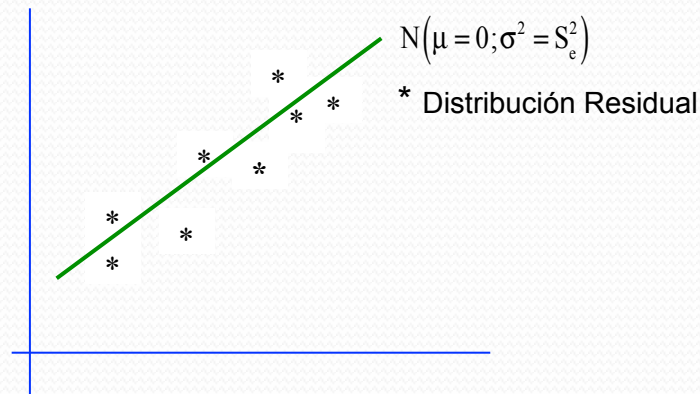
Incorrelados

Test de Hipótesis:

$$\begin{cases} H_0 : \rho = 0 \rightarrow \text{Incorrelados} \\ H_1 : \rho \begin{cases} > 0 \rightarrow \text{Autocorrelacion Positiva} \\ \neq 0 \rightarrow \text{Autocorrelados} \\ < 0 \rightarrow \text{Autocorrelacion Negativa} \end{cases} \end{cases}$$


4ª.- Normales

$$e \in N(\mu = 0; \sigma^2 = S_e^2)$$



## Transformaciones no lineales.

Tratan de corregir la falta de simetría de algunas distribuciones de datos. En ocasiones, estas transformaciones corrigen las desviaciones de normalidad.

La transformación mas utilizada es el logaritmo aunque existen otras mas simples como la raíz cuadrada o la inversa.

La transformación de Box-Cox (perteneciente a las transformaciones logarítmicas) también es muy usada y suele corregir las desviaciones de los supuestos necesarios para los modelos de regresión.

Dada una variable  $X$ , su transformada de Box-Cox se define como:

$$X^* = \left\{ \begin{array}{ll} \frac{(X+m)^\lambda - 1}{\lambda} & \text{si } \lambda \neq 0 \\ \ln(X+m) & \text{si } \lambda = 0 \end{array} \right\}$$

$m$  se escoge de tal forma que  $(X + m)$  sea positivo.

Si  $\lambda > 1$  se corrige la asimetría a la izquierda y si  $\lambda < 1$  la asimetría a la derecha.



# FUNDAMENTOS Y HERRAMIENTAS PARA LA MODELIZACIÓN DE PROCESOS TÉCNICOS-CIENTÍFICOS DE INVESTIGACIÓN

Bloque III: Modelización estadística.

**Regresión y correlación con R (Rcommander)**



Roberto Espejo Mohedano e-mail: [malesmor@uco.es](mailto:malesmor@uco.es)