

Tema 2: Estadística descriptiva: Dos variables.

Representación de tablas bidimensionales.

C \ E	0	1	$n_{i\cdot}$	
0	2	4	6	Distribución de C
1	4	8	12	
2	2	0	2	
$n_{\cdot j}$	8	12	20 = N	

Distribución de E

$n_{i\cdot}$: 2 veces aparece el dato (x_i, y_i) en mi muestra

Distribuciones marginales: Distribución de E y distribución de C.

$$f_{ij} = \frac{n_{ij}}{N}$$

N : tamaño de la muestra

$$N = \sum_{i=1}^2 \sum_{j=1}^3 n_{ij}$$

Tablas de frecuencia de C y E por separadas

C | $n_{i\cdot}$

0 | 6

1 | 12

2 | 2

E | $n_{\cdot j}$

0 | 8

1 | 12

$\Sigma = 20$

$\Sigma = 20$

INDEPENDENCIA: Relación entre variables

INDEPENDENCIA

X \ Y	C ₁	C ₂	C ₃	C ₄	
A	4	6	10	2	22
B	2	3	5	1	11
	6	9	15	3	33

X Y = C ₁	n_i	f_i
A	4	$4/6 = 2/3$
B	2	$2/6 = 1/3$
	6	

X Y = C ₂	n_i	f_i
A	6	$6/9 = 2/3$
B	3	$3/9 = 1/3$
	9	

X Y = C ₃	n_i	f_i
A	10	$10/15 = 2/3$
B	5	$5/15 = 1/3$
	15	

X Y = C ₄	n_i	f_i
A	2	$2/3$
B	1	$1/3$
	3	

Def 1 \rightarrow X es independiente de Y si $f_i^j =$

Def 2 \rightarrow Y es independiente de X si $f_i^j = f_{\cdot j} \quad \forall i, j$

Def 3 \rightarrow X es independiente de Y (y viceversa) si $f_{ij} = f_{i\cdot} \times f_{\cdot j}$

Propiedad: X es independiente de Y \Leftrightarrow Y es independiente de X.

DEPENDENCIA FUNCIONAL

X = "número de días en un año"

Y = "número de horas en un año"

$$Y = 24X \quad \left(\begin{array}{l} \text{Una en función} \\ \text{de la otra} \end{array} \right)$$

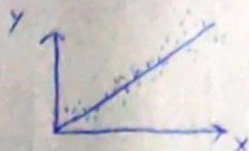
(Error = 0)

DEPENDENCIA ESTADÍSTICA

x = "número de monitores que hay en el aula de informática"

y = "número de computadores que hay en el aula de informática"

$$y = x + E \quad (E = \text{error})$$



DEPENDENCIA FUNCIONAL

x = "n° de días en un año" // y = "n° de horas en un año"

$$y = 24x \quad (\text{Error} = 0)$$

REGRESION Y CORRELACION

TIPOS DE AJUSTES

Método general para el ajuste de mínimos cuadrados

Supongamos que tenemos $\{(x_i, y_i)\}_{i=1}^N$ para dos variables cuantitativas x e y

↓ Tenemos una tabla simple de datos (¿Cómo puedo obtener una tabla simple de datos a partir de una tabla bidimensional?)

$x \backslash y$	1	2
3	0	1
4	9	4

tabla bidimensional

$$\{(x_i, y_i)\}_{i=1,2}^{j=1,2}$$

$x \backslash y$	1	2
3	1	0
4	2	1
9	1	4
8	2	4

tabla simple

$$\{(x_i, y_i)\}_{i=1}^4$$

Lo que buscamos es una función p que mejor "explique" los datos: Es decir, dado un dato (x_i, y_i) voy a tratar de aproximar $y_i \approx y_i^*$, donde:

$$y_i^* = p(x_i) \quad \text{Introduciendo notación}$$

•) $y_i^* = p(x_i)$ es el valor de " y " estimado por la regresión de x

•) $p_i = y_i - y_i^*$ es el error cometido por el ajuste por el dato i -ésimo (x_i, y_i)

¿Cuál será la idea para encontrar p ?

→ Minimizar p_i para todo $i = 1, 2, \dots, N$

Definimos la variable $e = y_i - y_i^*$ ($p_i = y_i - y_i^*$)

Lo que haremos es minimizar el "error cuadrático medio" de E

$$ECM_e = \sum_{i=1}^N p_i^2 = \sum_{i=1}^N (y_i - y_i^*)^2$$

→ Como minimizar el error anterior: ORGANIZO LOS DATOS: $\{(x_i, y_i)\}_{i=1}^N$

$$y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}; \quad y^* = \begin{pmatrix} y_1^* \\ y_2^* \\ \vdots \\ y_n^* \end{pmatrix} = \begin{pmatrix} p(x_1) \\ p(x_2) \\ \vdots \\ p(x_n) \end{pmatrix} = p(\bar{x})$$

Vamos a suponer que p es de la siguiente forma:

$$p(x) = a_0 p_0(x) + a_1 p_1(x) + \dots + a_5 p_5(x)$$

• Ajuste lineal $y = a_0 + a_1 x$

Se corresponde a $S=1$

$$\begin{cases} \varphi_0(x) = 1 \\ \varphi_1(x) = x \end{cases}$$

• Ajuste parabólico $y = a_0 + a_1 x + a_2 x^2$

$S=2$

$$\begin{cases} \varphi_0(x) = 1 \\ \varphi_1(x) = x \\ \varphi_2(x) = x^2 \end{cases}$$

• Ajuste trigonométrico $y = a_0 \cdot \cos(x) + a_1 \cdot \sin(x)$

$S=1$

$$\begin{cases} \varphi_0(x) = \cos(x) \\ \varphi_1(x) = \sin(x) \end{cases}$$

→ Este problema de regresión estará resuelto cuando encuentre a_0, a_1, \dots, a_s que minimizan el ECM anterior definido

$$y^* = \begin{pmatrix} \varphi_0 a_0(x_1) + \varphi_1 a_1(x_1) + \varphi_2 a_2(x_1) + \dots + \varphi_s a_s(x_1) \\ \varphi_0 a_0(x_2) + \varphi_1 a_1(x_2) + \dots + \varphi_s a_s(x_2) \\ \vdots \\ \varphi_0 a_0(x_n) + \varphi_1 a_1(x_n) + \dots + \varphi_s a_s(x_n) \end{pmatrix}$$

$$y^* = \underbrace{\begin{pmatrix} \varphi_0 a_0(x_1) + \varphi_1 a_1(x_1) + \dots + \varphi_s a_s(x_1) \\ \varphi_0 a_0(x_2) + \varphi_1 a_1(x_2) + \dots + \varphi_s a_s(x_2) \\ \vdots \\ \varphi_0 a_0(x_n) + \varphi_1 a_1(x_n) + \dots + \varphi_s a_s(x_n) \end{pmatrix}}_M \cdot \underbrace{\begin{pmatrix} a_0 \\ a_1 \\ \vdots \\ a_s \end{pmatrix}}_A = M \cdot A$$

Vamos a minimizar el ECM

$$\min_A \left\{ \sum_{i=1}^n e_i^2 \right\} = \min_A \left\{ \sum_{i=1}^n (y_i - y_i^*)^2 \right\} = \min_A \left\{ (y - y^*)^T (y - y^*) \right\} = \min_A \left\{ y^T y - y^T M A - (M A)^T y + (M A)^T (M A) \right\}$$

El mínimo se alcanza:

$$\text{Sea } J(A) = y^T y - y^T M A - (M A)^T y + (M A)^T (M A)$$

$$\frac{\partial J}{\partial A} = 0 \rightarrow$$

Eso ocurre cuando:

$$(M^T M) A = M^T y$$

Ecuaciones
normales

CASO PARTICULAR

Ajuste lineal $\rightarrow y = a_0 + a_1 x$ la llamamos recta de regresión Y/X

Encontrar A resolviendo el sistema lineal

$$(M^T M)A = M^T Y$$

$$\begin{cases} p_0(x) = 1 \\ p_1(x) = x \end{cases}$$

$$M = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \quad M^T = \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{pmatrix} \quad M^T M = \begin{pmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \end{pmatrix} \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} = \begin{pmatrix} N & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix}$$

Sistema lineal

$$M^T Y = \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix} \quad \begin{pmatrix} N & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \end{pmatrix} = \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix}$$

Ejemplo dispositivo 40

$$(M^T M)A = M^T Y$$

$$M = \begin{pmatrix} p_0(x_1) & p_1(x_1) & \dots & p_k(x_1) \\ p_0(x_2) & p_1(x_2) & \dots & p_k(x_2) \\ \vdots & \vdots & \dots & \vdots \\ p_0(x_n) & p_1(x_n) & \dots & p_k(x_n) \end{pmatrix}$$

X	2002	2003	2004	2005	...	2011
Y	93215	926308	963513	981403		1039183

Recta de regresión Y/X

1) Tabla de entrada

2) $N = 10$

3) Ecuaciones normales: $(M^T M)A = M^T Y$

4) Construimos M

$$M = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_{10} \end{pmatrix}, \quad A = \begin{pmatrix} a_0 \\ a_1 \end{pmatrix}$$

5) Cálculo $M^T M$

$$M^T M = \begin{pmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_{10} \end{pmatrix} \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_{10} \end{pmatrix} = \begin{pmatrix} 10 & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix} = \begin{pmatrix} 10 & 20065 \\ 20065 & 40260505 \end{pmatrix}$$

6) Cálculo $M^T Y$

$$M^T Y = \begin{pmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_{10} \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_{10} \end{pmatrix} = \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix} = \begin{pmatrix} 10236140 \\ 2'154529648 \cdot 10^3 \end{pmatrix}$$

7) Resolvemos el sistema lineal

$$\begin{pmatrix} 10 & 20065 \\ 20065 & 40260505 \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \end{pmatrix} = \begin{pmatrix} 10236140 \\ 2'154529648 \cdot 10^3 \end{pmatrix}$$

8) Obtenemos la recta

$$Y(x) = a_0 + a_1 x$$

Este "modelo matemático" explica la variable Y en función de la variable X

→ ¿Podría encontrar otro "modelo matemático" lineal que explique x en función de y ?
 Si, a eso la llamamos recta de regresión de $x|y$

$$x|y = a_0 + a_1 y$$

El método general ya visto, puede aplicarse también.

•) Modelo: $x = a_0 + a_1 y$

•) $(M^t M)M = M^t X \rightarrow$ Pongo lo que quiero explicar

•) Construyo M . Como en este caso el modelo es $x = a_0 + a_1 y$, $\begin{cases} p_0(y) = 1 \\ p_1(y) = y \end{cases} \quad M = \begin{pmatrix} 1 & y_1 \\ \vdots & \vdots \\ 1 & y_n \end{pmatrix}$

$$M^t M = \begin{pmatrix} 1 & \dots & 1 \\ y_1 & \dots & y_n \end{pmatrix} \begin{pmatrix} 1 & y_1 \\ \vdots & \vdots \\ 1 & y_n \end{pmatrix} = \begin{pmatrix} N & \sum y_i \\ \sum y_i & \sum y_i^2 \end{pmatrix} \quad M^t X = \begin{pmatrix} 1 & \dots & 1 \\ y_1 & \dots & y_n \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} \sum x_i \\ \sum x_i y_i \end{pmatrix}$$

•) Resolver

$$\begin{pmatrix} N & \sum y_i \\ \sum y_i & \sum y_i^2 \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \end{pmatrix} = \begin{pmatrix} \sum x_i \\ \sum x_i y_i \end{pmatrix}$$

$$\begin{aligned} \rightarrow \text{Recta } y|x &= \frac{y - \bar{y}}{\sigma_y} = r \frac{(x - \bar{x})}{\sigma_x} \\ \rightarrow \text{Recta } x|y &= \frac{x - \bar{x}}{\sigma_x} = r \frac{(y - \bar{y})}{\sigma_y} \end{aligned} \quad \left\{ \begin{aligned} r &= \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} \end{aligned} \right.$$

" r " → "Coeficiente de Correlación de Pearson", mide la magnitud y dirección de las rectas, es decir, da una idea de cómo es la pendiente de las rectas (tanto $y|x$ como $x|y$)

- ¿Qué recta explica mejor los datos?

La recta que menor error cuadrático medio tenga con respecto a los datos

ECM de $y|x$

$$y = \bar{y} + \frac{r \sigma_y}{\sigma_x} (x - \bar{x})$$

$$\text{El error cuadrático medio es: } \frac{1}{N} \sum (y_i - \bar{y})^2 = \dots = \frac{1}{N} \sum \left(y_i - \left(\bar{y} + \frac{\sigma_y}{\sigma_x} r (x_i - \bar{x}) \right) \right)^2 = \dots = \sigma_y^2 (1 - r^2)$$

Del mismo modo, el ECM de $x|y$ es $\sigma_x^2 (1 - r^2)$

$$\rightarrow \begin{cases} \text{Será mejor } y|x & \text{si } \sigma_y^2 (1 - r^2) < \sigma_x^2 (1 - r^2) \\ \text{Será mejor } x|y & \text{si } \sigma_x^2 (1 - r^2) < \sigma_y^2 (1 - r^2) \end{cases}$$

Observaciones

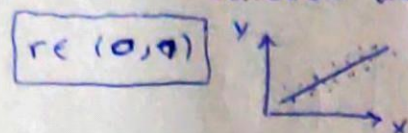
• $r \in [-1, 1]$

$$\frac{1}{N} \sum (x_i - \bar{x})^2 = \sigma_x^2 (1 - r^2) \geq 0$$

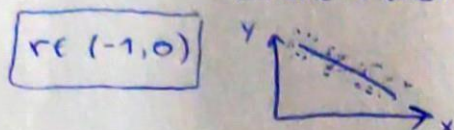
$\Rightarrow 1 - r^2 \geq 0; r^2 \leq 1 \rightarrow r \in [-1, 1]$

*

• Si $r > 0 \rightarrow$ Correlación lineal directa



• Si $r < 0 \rightarrow$ Correlación lineal inversa



• Si $r = 0$

\rightarrow recta $y/x \Rightarrow y = \bar{y} \quad \left(\frac{y - \bar{y}}{\sigma_y} = r \frac{(x - \bar{x})}{\sigma_x} \right)$ Variables incorreladas

\rightarrow recta $x/y \Rightarrow x = \bar{x}$

• Si $r = 1$ o $r = -1 \Rightarrow$ El error cuadrático es 0 \Rightarrow Correlación lineal perfecta.

\rightarrow Definición: Dada una nube de puntos llamamos vector residuo a:

$$E: e_i = y_i - y_i^* \quad (y_i^* \text{ valor estimado, es decir } y_i^* = p(x_i))$$

\rightarrow Definición: Llamaremos varianza residual a:

$$V_r = \frac{1}{N} \sum (e_i - \bar{e})^2 = \frac{1}{N} \sum e_i^2 - \bar{e}^2$$

\rightarrow Definición: Llamaremos coeficiente de determinación a:

$$R^2 = 1 - \frac{V_r}{V_y} \rightarrow y \text{ es la variable que quiero explicar}$$

Observaciones:

• El coeficiente R^2 coincide con el coeficiente de correlación de Pearson al cuadrado para el caso en que tenga las rectas x/y o y/x : $R^2 = r^2$

• Si $R^2 = 1 \rightarrow$ Ajuste perfecto

• Si $R^2 = 0 \rightarrow$ Variables incorreladas

Ej 2 Ajuste de un plano

x	y	z
0	0	0
1	1	2
0	1	1

Ajuste mediante un plano

$$z = \hat{y}(x, y) = a_0 + a_1 x + a_2 y$$

N=3

$$M^T M = \begin{pmatrix} 1 & 1 & 1 \\ x_1 & x_2 & x_3 \\ y_1 & y_2 & y_3 \end{pmatrix} \begin{pmatrix} 1 & x_1 & y_1 \\ 1 & x_2 & y_2 \\ 1 & x_3 & y_3 \end{pmatrix} = \begin{pmatrix} N & \sum x_i & \sum y_i \\ \sum x_i & \sum x_i^2 & \sum x_i y_i \\ \sum y_i & \sum x_i y_i & \sum y_i^2 \end{pmatrix} = \begin{pmatrix} 3 & 1 & 2 \\ 1 & 1 & 1 \\ 2 & 1 & 2 \end{pmatrix}$$

$$M^T Z = \begin{pmatrix} 1 & 1 & 1 \\ x_1 & x_2 & x_3 \\ y_1 & y_2 & y_3 \end{pmatrix} \begin{pmatrix} z_1 \\ z_2 \\ z_3 \end{pmatrix} = \begin{pmatrix} \sum z_i \\ \sum z_i x_i \\ \sum z_i y_i \end{pmatrix} = \begin{pmatrix} 3 \\ 2 \\ 3 \end{pmatrix} \quad \begin{array}{l} \text{Sol. al sistema lineal} \\ a_0 = 0 \quad a_1 = 1 \quad a_2 = 1 \end{array}$$

$$z = \hat{y}(x, y) = 0 + x + y = x + y \quad z^* = x + y = \hat{y}(x, y)$$

$$R^2 = 1 - \frac{V_r}{V_t} \quad \begin{array}{l} \text{varianza de la} \\ \text{variable que queremos explicar} \end{array} \Rightarrow R^2 = 1 - 0 = 1$$

$$V_r = \frac{1}{N} \sum (E - \bar{E})^2; \quad E = z - z^*; \quad E = \sum_{\text{error}} z - \hat{y}(x, y) = z - x - y;$$

$$\bar{E} = \overline{(z - x - y)} = \bar{z} - \bar{x} - \bar{y} = \frac{3}{3} - \frac{1}{3} - \frac{2}{3} = 0$$

$$V_r = \frac{1}{N} \sum_{i=1}^N E_i^2 = \frac{1}{3} (E_1^2 + E_2^2 + E_3^2) = \frac{1}{3} ((z_1 - (x_1 + y_1))^2 + (z_2 - (x_2 + y_2))^2 + (z_3 - (x_3 + y_3))^2) = 0;$$

Como $R^2 = 1$ el ajuste es perfecto

Observados: z_i

Estimados: $z_i = \hat{y}(x_i, y_i) \quad (x_i, y_i, \hat{y}(x_i, y_i))$