

## Tema 2

### • Distribuciones condicionadas

$$X/Y=y_j \quad X/y_j$$

$$f_{ij}^j = \frac{n_{ij}}{n_{\cdot j}} = \frac{f_{ij}}{f_{\cdot j}}$$

$$f_{ij} = f_{i \cdot} \cdot f_{\cdot j}$$

$$Y/X=x_i$$

$$f_{ij}^i = \frac{n_{ij}}{n_{i \cdot}} = \frac{f_{ij}}{f_{i \cdot}}$$

$$f_{ij} = f_{ij}^i \cdot f_{i \cdot}$$

### • Relación entre variables

- El objetivo de analizar conjuntamente dos variables diferentes es establecer el tipo de relación existente entre ellas.

Hay 3 casos:

- Independientes: No hay relación alguna entre las variables.
- Dependencia funcional: El valor de una variable queda determinado conociendo el valor de la otra variable para esa misma observación a través de una función.
- Dependencia estadística: Una variable proporciona información sobre la otra pero la modalidad de una no queda determinada por la modalidad de la otra.

## - Independencia entre variables

$X$  es independiente de  $Y$  si:

$$X / Y=y_j \equiv X$$

Ejemplo:

$X \backslash Y$	$C_1$	$C_2$	$C_3$	$C_4$	
A	4	6	10	2	22
B	2	3	5	1	11

$Y/X=A$	$g_j^A$
$C_1$	2/11
$C_2$	3/11
$C_3$	5/11
$C_4$	1/11

$Y/X=B$	$g_j^B$
$C_1$	2/11
$C_2$	3/11
$C_3$	5/11
$C_4$	1/11

$Y_j$	$g_j$
$C_1$	2/11
$C_2$	3/11
$C_3$	5/11
$C_4$	1/11

Si  $X$  es independiente de  $Y$

$$g_j^i = g_i \quad \forall i \quad \forall j \quad g_j^i = \frac{n_{ij}}{n_{\cdot j}} : g_i = \frac{n_{i\cdot}}{N}$$

## - Dependencia funcional

$X \backslash Y$	$C_1$	$C_2$	$C_3$	$C_4$
$A_1$	4	6	0	0
$A_2$	0	0	5	7

22

$X/Y=C_3$  sólo toma el valor de  $A_2$   
 $X/Y=C_1$  " " " "  $A_1$   
 $X/Y=C_2$  " " " "  $A_1$   
 $X/Y=C_4$  " " " "  $A_2$

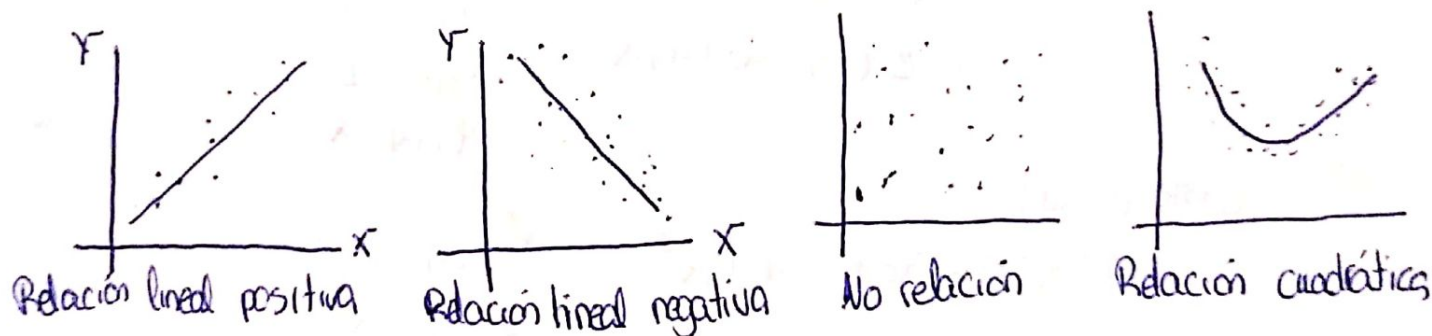
## - Dependencia estadística

- Ejemplo:
- Estatura y peso
  - Nacionalidad y renta
  - Familia por n° de hijos y n° de móviles

### + Pasos a seguir:

- 1) Nube de puntos
- 2) Buscar la línea o curva de regresión que mejor se ajuste a la nube de puntos.  $\rightarrow$  Regresión
- 3) Medir el grado de dependencia entre las variables.  $\rightarrow$  Correlación

Si todos los valores satisfacen la ecuación calculada, se dice que las variables ~~están~~ perfectamente correladas. La ecuación nos permite predecir valores desconocidos.

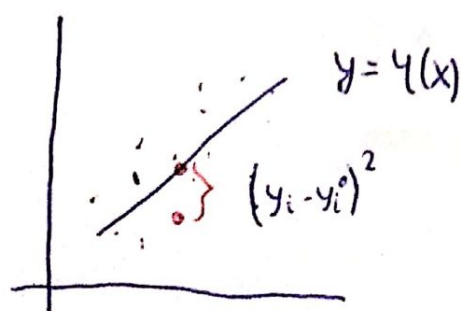


### • Regresión Lineal para tipo cuantitativas

#### • Método de mínimos cuadrados

Sean los datos  $\{(x_i, y_i)\}$  para dos v.e. cuantitativas  $X, Y$ .

#### 1. La Regresión de $Y/X$



El objetivo es minimizar

$$\sum_i (y_i - y_i^0)^2 = \sum_i (y_i - y(x_i))^2 = F(x)$$



$$\text{Min } F \rightarrow \text{Min } \sum_i e_i^2 \rightarrow \text{Minimizar el SSE}$$

$y_i^{\circ} = y_i^{\text{est}} = \varphi(x_i)$  es el valor de  $y$  estimado por la regresión para  $x_i$ .

$e_i = y_i - y_i^{\text{est}}$  es el error cometido por el ajuste para el  $i$ -ésimo dato.

El tipo de ajuste de mínimo cuadrados viene determinado por el tipo de función  $y = \varphi(x)$  elegido.

### + Caso Lineal Generalizado

$$\varphi(x) = a_0 \varphi_0(x) + a_1 \varphi_1(x) + a_2 \varphi_2(x) + \dots$$

Los más usados son:

- Ajuste lineal:  $y = \varphi(x) = a + bx$   $a_1 = a_1$   
 $\varphi(x) = a_0 + a_1 x$   $a_0 = a_0$   
 $\varphi_0(x) = 1$   
 $\varphi_1(x) = x$

- Ajuste parabólico:  
 $y = \varphi(x) = a_0 + a_1 x + a_2 x^2$   $\varphi_0(x) = 1$   
 $\varphi_1(x) = x$   
 $\varphi_2(x) = x^2$

- Otros ajustes:  
 $y = \varphi(x) = a_0 \cdot \cos(x) + a_1 \cdot \sin(x)$   $\varphi_0(x) = \cos x$   
 $\varphi_1(x) = \sin x$

### o Ajuste lineal. Propiedades

$$\begin{aligned} \min F &= \min \sum_i (y_i - y_i^{\text{est}})^2 = \min \sum_i (y_i - \varphi(x_i))^2 = \\ &= \min_{a,b} \sum (y_i - (a + bx_i))^2 \end{aligned}$$

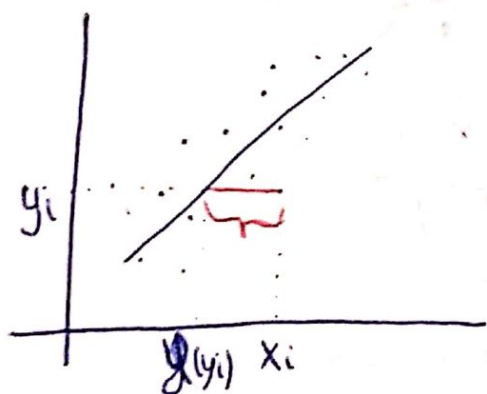
$$\frac{\partial F}{\partial a} = \sum_i 2(y_i - a - bx_i) \cdot (-1) = 0$$

$$\frac{\partial F}{\partial b} = \sum_i 2(y_i - a - bx_i) \cdot (-x_i) = 0$$

$$\left. \begin{aligned} -\sum_i y_i + a \sum_i 1 + b \sum_i x_i &= 0 \\ -\sum_i y_i x_i + a \sum_i x_i + b \sum_i x_i^2 &= 0 \end{aligned} \right\} \begin{aligned} a \cdot N + b \sum_i x_i &= \sum_i y_i \\ a \cdot \sum_i x_i + b \sum_i x_i^2 &= \sum_i x_i y_i \end{aligned}$$

$$\begin{pmatrix} N & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix}$$

+ Regresión de  $X/Y$



$$x = \phi'(y)$$

$$\begin{aligned} \text{Min } \sum_i (x_i - \phi'(y_i))^2 &= \\ = \text{Min } \sum_i (x_i - x_i^{\text{est}})^2 &= \\ = \text{Min } \sum_i (x_i - x_i^{\text{est}})^2 &= \end{aligned}$$

+ Ajuste lineal

$$\phi'(y) = a' + b'y$$

$$\begin{pmatrix} N & \sum_i y_i \\ \sum_i y_i & \sum_i y_i^2 \end{pmatrix} \begin{pmatrix} a' \\ b' \end{pmatrix} = \begin{pmatrix} \sum x_i \\ \sum x_i y_i \end{pmatrix}$$

Sist de Ec. normales

• Ajuste lineal: Propiedades

+  $Y/X$  5. Ecuaciones Normales

$$\left. \begin{aligned} Na + b \sum x_i &= \sum y_i \\ a \sum x_i + b \sum x_i^2 &= \sum x_i y_i \end{aligned} \right\} \left. \begin{aligned} a + b \frac{\sum x_i}{N} &= \frac{\sum y_i}{N} \\ a \frac{\sum x_i}{N} + b \frac{\sum x_i^2}{N} &= \frac{\sum x_i y_i}{N} \end{aligned} \right\}$$

$$\left. \begin{aligned} a + b \bar{x} &= \bar{y} \\ a \bar{x} + b m_{20} &= m_{11} \end{aligned} \right\} \textcircled{1} (\bar{x}, \bar{y}) \text{ pasa por la recta de regresión de } Y/X$$

$$\begin{aligned} a &= \bar{y} - b \bar{x} \rightarrow (\bar{y} - b \bar{x}) \bar{x} - b m_{20} = m_{11} \\ \bar{y} \bar{x} - b \bar{x}^2 + b m_{20} &= m_{11} \\ b (m_{20} - \bar{x}^2) &= m_{11} - \bar{y} \bar{x} \\ b &= \frac{m_{11} - \bar{y} \bar{x}}{m_{20} - \bar{x}^2} = \frac{\text{Cov}(X, Y)}{V^2 X} = b \end{aligned}$$

$$\left. \begin{aligned} Na' + b' \sum y_i &= \sum x_i \\ a' \sum y_i + b' \sum y_i^2 &= \sum x_i y_i \end{aligned} \right\} \left. \begin{aligned} a' + b' \bar{y} &= \bar{x} \\ b' &= \frac{\text{Cov}(X, Y)}{V^2 Y} \end{aligned} \right\}$$

②  $(\bar{x}, \bar{y})$  pasa por la recta de regresión de  $X/Y$   
 $(\bar{x}, \bar{y})$  es el punto de corte de las dos rectas de regresión

③  $\text{Signo } b = \text{Signo } b' = \text{Signo } \text{Cov}(X, Y)$

$$\sum_i e_i = 0 \quad e_i = y_i - y_i^o \quad y_i = a + b x_i$$

$$\sum_i e_i = \sum_i (y_i - y_i^o) = \sum_i (y_i - (a + b x_i)) = \sum_i y_i - a \cdot N - b \sum_i x_i = 0$$

1ª Ec. del sistema de ec. normales

$$\sum_i e_i x_i = 0$$

$$\sum_i (y_i - y_i^o) x_i = \sum_i (y_i - a - b x_i) x_i = \sum_i y_i x_i - a \sum_i x_i - b \sum_i x_i^2 = 0$$

2ª Ec. del S.E.N



• Descomposición de la varianza

$$\begin{aligned} \sigma^2_y &= \frac{1}{N} \sum_i (y_i - \bar{y})^2 = \frac{1}{N} \sum_i (y_i - y_i^* + y_i^* - \bar{y})^2 = \\ &= \frac{1}{N} \sum_i \left( (y_i - y_i^*) + (y_i^* - \bar{y}^*) \right)^2 = \frac{1}{N} \left[ \sum_i (y_i - y_i^*)^2 + \sum_i (y_i^* - \bar{y}^*)^2 + \right. \\ &\quad \left. 2 \sum_i (y_i - y_i^*) (y_i^* - \bar{y}^*) \right] = \frac{1}{N} \sum_i e_i^2 + \sigma_{y^*}^2 + \frac{1}{N} 2 \sum_i e_i \cdot (y_i^* - \bar{y}^*) = \\ &= MSE + \sigma_{y^*}^2 + \frac{1}{N} 2 \sum_i e_i (a + b x_i - (a + b \bar{x})) = MSE + \sigma_{y^*}^2 \end{aligned}$$

• Varianza residual o varianza de los errores

Sabíamos que:  $\sum_i e_i = 0 \Rightarrow \bar{e} = 0$

" "  $\sigma_{y^*}^2 = MSE + \sigma_{y^*}^2$

$$\sigma_{\text{resi}}^2 = \frac{1}{N} \sum_i (e_i - \bar{e})^2 = \frac{1}{N} \sum_i e_i^2 - \bar{e}^2 = MSE$$

$$\sigma_y^2 = \sigma_{\text{res}}^2 + \sigma_{y^*}^2$$

¿Quién es  $\sigma_{y^*}^2$ ?

$$\begin{aligned} \sigma_{y^*}^2 &= \frac{1}{N} \sum_i (y_i^* - \bar{y}^*)^2 = \frac{1}{N} \sum_i (y_i^* - \bar{y})^2 = \frac{1}{N} \sum_i (a + b x_i - (a + b \bar{x}))^2 = \\ &= \frac{1}{N} \sum_i (b (x_i - \bar{x}))^2 = b^2 \frac{1}{N} \sum_i (x_i - \bar{x})^2 = b^2 \sigma_x^2 \end{aligned}$$

Por tanto,  $\sigma_{\text{res}}^2 = \sigma_y^2 - b^2 \sigma_x^2$

## • Coefficiente de correlación lineal de Pearson

Mide el grado de relación lineal entre las variables. Se define como la media geométrica de los coeficientes de regresión de  $b$  y  $b'$

$$r = 1 = \sqrt{b \cdot b'}$$

signo de  $b$

signo de  $b'$

$$\mu_{12} = \text{Cov}(X, Y)$$

$$b = \frac{\text{Cov}(X, Y)}{\sigma^2_x}$$

$$b' = \frac{\text{Cov}(X, Y)}{\sigma^2_y}$$

$$r = \sqrt{\frac{\mu_{12}^2}{\sigma^2_x \sigma^2_y}} = \frac{\mu_{12}}{\sigma_x \sigma_y} \Rightarrow r = \frac{\mu_{12}}{\sigma_x \cdot \sigma_y}$$

$$r^2 \geq 0$$

-  $b$  y  $b'$  en función de  $r$ :

$$b = \frac{r \sigma_y}{\sigma_x} \quad b' = \frac{r \sigma_x}{\sigma_y}$$

- Propiedades de  $r^2$

Sabiendo que:  $\sigma^2_{res} = \sigma^2_y - b^2 \sigma^2_x$

como  $b = \frac{r \sigma_y}{\sigma_x} \Rightarrow \sigma^2_{res} = \sigma^2_y - \frac{r^2 \sigma^2_y}{\sigma_x^2} \cdot \sigma^2_x$

$$\frac{\sigma^2_{res}}{\sigma^2_y} = 1 - r^2 \Rightarrow r^2 = 1 - \frac{\sigma^2_{res}}{\sigma^2_y}$$

• Si  $r^2 = 1$  ( $r = 1$  ó  $r = -1$ )  
Correlación lineal es perfecta  $\begin{cases} r = 1 \rightarrow \text{directa} \\ r = -1 \rightarrow \text{inversa} \end{cases}$

• Si  $r = 0$   
Variables incorreladas

Cuando más próximo a 1 se encuentre el  $r^2$ , mejor será el ajuste lineal.



• ¿Cuándo es mejor una recta que otra?

① Depende de lo que queramos predecir

② Diremos que  $y/x$  es mejor que  $x/y$  si

$$\sigma^2_{\text{res } y/x} \leq \sigma^2_{\text{res } x/y}$$

$$\sigma_y^2 (1-r^2) \leq \sigma_x^2 (1-r^2)$$

$$\underline{\sigma_y^2 \leq \sigma_x^2}$$

• Relación entre  $b$  y  $b'$

$$y = a + bx$$

$$x = a' + b'y$$

$$y = \frac{x - a'}{b'}$$

$$m_{y/x} = b$$

$$m_{x/y} = \frac{1}{b'}$$

$$r^2 = |b \cdot b'| \leq 1 \iff r^2 = |b| \cdot |b'| \leq 1$$

$$|b| \leq \frac{1}{|b'|}$$

$$\underline{|m_{y/x}| \leq |m_{x/y}|}$$

## Modelo lineal generalizado

Ajustamos la nube de puntos a una función del tipo:

$$\varphi(x) = a_0 \varphi_0(x) + a_1 \varphi_1(x) + a_2 \varphi_2(x) + \dots$$

El objetivo:

- minimizar  $\sum_i (y_i - y_i^*)^2$  ( $x$  indep)

- minimizar  $\sum_i (y_i - \varphi(x_i))^2$  ya que  $y_i^* = y_i^{est} = \varphi(x_i)$

- minimizar  $\sum e_i^2$

• Considero los vectores:

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} \quad x = \begin{pmatrix} x_1 \\ \vdots \\ x_N \end{pmatrix}$$

$$y^* = \varphi(x) = \begin{pmatrix} y_1^* \\ y_2^* \\ \vdots \\ y_N^* \end{pmatrix} = \begin{pmatrix} \varphi(x_1) \\ \varphi(x_2) \\ \vdots \\ \varphi(x_N) \end{pmatrix}$$

$$M = \begin{pmatrix} \varphi_0(x_1) & \varphi_1(x_1) & \dots \\ \varphi_0(x_2) & \varphi_1(x_2) & \dots \\ \vdots & \vdots & \ddots \\ \varphi_0(x_N) & \varphi_1(x_N) & \dots \end{pmatrix}$$

$$A = \begin{pmatrix} a_0 \\ a_1 \\ a_2 \\ \vdots \end{pmatrix}$$

$$e = y - y^* = \begin{pmatrix} y_1 - y_1^* \\ y_2 - y_2^* \\ \vdots \\ y_N - y_N^* \end{pmatrix}$$

①  $\varphi(x_1) = a_0 \varphi_0(x_1) + a_1 \varphi_1(x_1) + a_2 \varphi_2(x_1) + \dots$

$\varphi(x_2) = a_0 \varphi_0(x_2) + a_1 \varphi_1(x_2) + a_2 \varphi_2(x_2) + \dots$

$\vdots$   
 $\varphi(x_N) = a_0 \varphi_0(x_N) + a_1 \varphi_1(x_N) + a_2 \varphi_2(x_N) + \dots$

$$\varphi(x) = M \cdot A$$

②  $(y - y^*)^t (y - y^*) = (y_1 - y_1^* \dots y_N - y_N^*) \begin{pmatrix} y_1 - y_1^* \\ \vdots \\ y_N - y_N^* \end{pmatrix} = \sum_i (y_i - y_i^*)^2$

$$\min \sum_i e_i^2$$

$$\min_{a_1, a_2, a_3, \dots} \sum (y_i - y_i^*)^2$$

$$\min_A (y - y^*)^t (y - y^*)$$

$$\min_A (y - \phi(x))^t (y - \phi(x))$$

$$\min_A (y - MA)^t (y - MA)$$

$$\min_A (y^t - A^t M^t) (y - MA)$$

$$\min_A y^t y - y^t MA - A^t M^t y + A^t M^t MA$$

$$J(A) = y^t y - 2 M^t y A + A^t M^t MA$$

$$\frac{\partial J}{\partial A} = 0 \Rightarrow -2 M^t y + 2 M^t MA = 0$$

$$M^t MA = M^t y \rightarrow \text{Sistema de Ec. normales}$$

$$\phi(x) = a + bx \quad \phi_0(x) = 1 \quad \phi_1(x) = x$$

$$M = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_N \end{pmatrix} \quad A = \begin{pmatrix} a \\ b \end{pmatrix} \quad y = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix}$$

$$M^t MA = \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_N \end{pmatrix} \cdot \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_N \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} N & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix}$$

$$M^t y = \begin{pmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_N \end{pmatrix} \cdot \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} = \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix}$$



- Ajuste mediante un plano

$$z = \varphi(x, y) = a + bx + cy$$

$$z_i = a + bx_i + cy_i$$

$$M = \begin{pmatrix} 1 & x_1 & y_1 \\ 1 & x_2 & y_2 \\ \vdots & \vdots & \vdots \\ 1 & x_N & y_N \end{pmatrix} \quad A = \begin{pmatrix} a \\ b \\ c \end{pmatrix}$$

- Ecuaciones normales

$$M^t M A = M^t Z \rightarrow \text{dependiente}$$

$$\begin{pmatrix} N & \sum x_i & \sum y_i \\ \sum x_i & \sum x_i^2 & \sum x_i y_i \\ \sum y_i & \sum x_i y_i & \sum y_i^2 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} \sum z_i \\ \sum x_i z_i \\ \sum y_i z_i \end{pmatrix}$$

- Otros ajustes no lineales

- $y = ae^{bx} \rightarrow$  Ajuste exponencial

- Tomamos logaritmos para linealizar el modelo

$$\ln y = \ln a e^{bx}$$

$$\ln y = \ln a + \ln e^{bx}$$

$$\ln y = \ln a + bx$$

- Llamamos  $\ln y = \hat{y}$  a una nueva variable  
 $\ln a = \hat{a}$

$$\underline{\hat{y} = \hat{a} + bx}$$