

Regresión y correlación lineal II

ESCUELA POLITÉCNICA
SUPERIOR DE CÓRDOBA

Universidad de Córdoba

DEPARTAMENTO DE ESTADÍSTICA



Regresión parabólica.

$$\hat{y} = a + b_1x + b_2x^2$$

Usando el método de mínimos cuadrados, minimizamos la expresión:

$$H = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - a - b_1x_i - b_2x_i^2)^2$$

para ello derivamos parcialmente e igualamos a cero:

$$\left. \begin{aligned} \frac{\partial H}{\partial a} &= -2 \sum_{i=1}^n (y_i - a - b_1x_i - b_2x_i^2) = 0 \\ \frac{\partial H}{\partial b_1} &= -2 \sum_{i=1}^n (y_i - a - b_1x_i - b_2x_i^2)x_i = 0 \\ \frac{\partial H}{\partial b_2} &= -2 \sum_{i=1}^n (y_i - a - b_1x_i - b_2x_i^2)x_i^2 = 0 \end{aligned} \right\} \text{ Ecuaciones Normales}$$

Resolviendo el sistema, se obtienen los parámetros: a , b_1 y b_2 .



Podemos simplificar realizando el cambio $Z=X^2$, con lo que:

$$\hat{y} = a + b_1x + b_2z$$

Y de nuevo aplicando mínimos cuadrados:

$$\left. \begin{aligned} \frac{\partial H}{\partial a} &= -2 \sum_{i=1}^n (y_i - a - b_1x_i - b_2z_i) = 0 \\ \frac{\partial H}{\partial b_1} &= -2 \sum_{i=1}^n (y_i - a - b_1x_i - b_2z_i)x_i = 0 \\ \frac{\partial H}{\partial b_2} &= -2 \sum_{i=1}^n (y_i - a - b_1x_i - b_2z_i)z_i = 0 \end{aligned} \right\} \text{ Ecuaciones Normales}$$

resolviendo el sistema, se obtienen los valores de los parámetros.

Conclusiones:

→ $\bar{\varepsilon} = 0$

→ $a = \bar{y} - b_1\bar{x} - b_2\bar{z}$

$(\bar{y} = a + b_1\bar{x} + b_2\bar{z})$

$$\rightarrow COV(\varepsilon X) = 0 \quad \rightarrow COV(\varepsilon Z) = 0$$

$$\rightarrow b_i = -\frac{A_{1,i+1}}{A_{1,1}} \text{ para } i = 1, 2. \text{ Donde } A_{1,1} \text{ es el adjunto al elemento } a_{1,1} \text{ de la matriz de covarianzas } \Sigma.$$

$$\Sigma = \begin{pmatrix} S_y^2 & S_{yx} & S_{yz} \\ S_{xy} & S_x^2 & S_{xz} \\ S_{zy} & S_{zx} & S_z^2 \end{pmatrix} \text{ y } A_{1,i+1} \text{ el adjunto al elemento } a_{1,i+1}$$

Al menor complementario del elemento $a_{1,1}$ de Σ , es decir:

$$\Sigma_{\vec{x}} = \begin{pmatrix} S_x^2 & S_{xz} \\ S_{zx} & S_z^2 \end{pmatrix} \text{ se le denomina matriz de covarianzas de las } \vec{x} \text{ (v. indep.)}$$

Este proceso se puede generalizar al caso **lineal múltiple**, dado el hiperplano de regresión:

$$\hat{y} = a + b_1 x_1 + b_2 x_2 + \cdots + b_n x_n$$

Regresión exponencial y potencial.

Exponencial: $\hat{y} = ak^{bx}, k > 0$

Si tomamos logaritmos: $\log_k \hat{y} = \log_k a + bx \log_k k \Rightarrow \hat{y}^* = a^* + b^* x$

donde: $\left. \begin{array}{l} \hat{y}^* = \log_k \hat{y} \\ a^* = \log_k a \\ b^* = b \log_k k = b \end{array} \right\}$ con lo que se ha transformado el problema en uno de tipo lineal. Para obtener finalmente a , basta tomar antilogaritmos.

Potencial: $\hat{y} = ax^b$

tomando logaritmos de nuevo: $\log \hat{y} = \log a + b \log x \Rightarrow \hat{y}^* = a^* + bx^*$

donde: $\left. \begin{array}{l} \hat{y}^* = \log \hat{y} \\ a^* = \log a \\ x^* = \log x \end{array} \right\}$ de nuevo se ha transformado el problema en uno de tipo lineal.

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

Varianza residual.

$$V(\varepsilon) = S_\varepsilon^2 = \frac{1}{n} \sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon})^2 = \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2$$

Por otra parte:

$$S_\varepsilon^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{1}{n} \sum_{i=1}^n (y_i - a - bx_i)^2 = \dots = S_y^2 \left(1 - \frac{S_{xy}^2}{S_x^2 S_y^2} \right) = S_y^2 (1 - R_{xy}^2)$$

Coefficiente de determinación: $R^2 = R_{xy}^2 = \left(\frac{S_{xy}}{S_x S_y} \right)^2 = 1 - \frac{S_\varepsilon^2}{S_y^2} \in [0, 1]$

Siendo: $S_\varepsilon^2 = \frac{\det \Sigma}{\det \Sigma_{\bar{x}}}$ "Tanto por 1 (ó %) de la varianza de Y explicada por el modelo."

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



Correlación parcial.

Objetivo: Eliminar los efectos distorsionantes de terceras variables sobre las relaciones lineales entre variables.

Fundamento: Dados los modelos:

$$x = a + bz + \varepsilon_x \quad ; \quad y = a' + b'z + \varepsilon_y$$

ε_x representa la parte de X que no es capaz de explicar Z (luego Z no interviene), y de forma análoga ε_y .

Por tanto se define **el coeficiente de correlación entre X e Y , parcial Z** como:

$$\rho_{xy,z} = \rho_{\varepsilon_x \varepsilon_y}$$

Operativamente se puede obtener como:

$$\rho_{xy,z} = \frac{\rho_{xy} - \rho_{xz}\rho_{yz}}{\sqrt{1 - \rho_{xz}^2} \sqrt{1 - \rho_{yz}^2}} = -\frac{A_{12}}{\sqrt{A_{11}A_{22}}}$$



Conclusiones:

$$\rho_{xy,z}^2 > \rho_{xy}^2 \Rightarrow \begin{cases} \text{La variable } Z \text{ oculta o amortigua la} \\ \text{dependencia entre } X \text{ e } Y \end{cases}$$

$$\rho_{xy,z}^2 \cong 0 \Rightarrow \begin{cases} \text{La interdependencia entre } X \text{ e } Y \text{ se debe} \\ \text{casi exclusivamente al efecto de } Z \end{cases}$$

$$\rho_{xy,z}^2 < \rho_{xy}^2 \Rightarrow \begin{cases} \text{La interdependencia entre } X \text{ e } Y \text{ se debe} \\ \text{parcialmente a la influencia de } Z \end{cases}$$

ESCUELA POLITÉCNICA
SUPERIOR DE CÓRDOBA

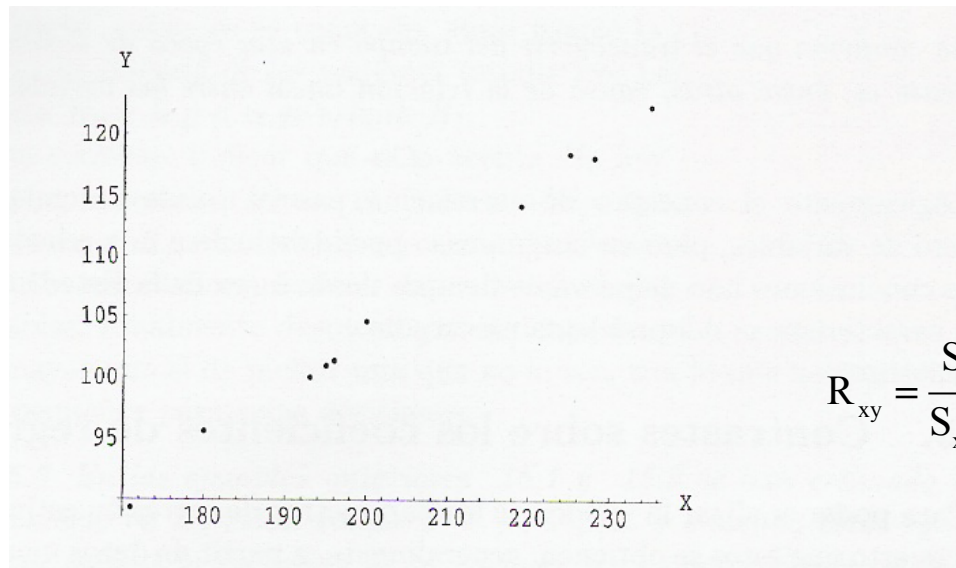
Universidad de Córdoba

DEPARTAMENTO DE ESTADÍSTICA





| | | | | | | | | | | |
|----------|------|------|-----|-----|-------|-------|-------|-------|-------|-------|
| X | 171 | 180 | 193 | 195 | 196 | 200 | 219 | 225 | 228 | 235 |
| Y | 89.3 | 95.6 | 100 | 101 | 101.4 | 104.6 | 114.3 | 118.5 | 118.2 | 122.4 |



$$\bar{X} = 204.2 ; \bar{S}_x = 21.47246$$

$$\bar{Y} = 106.53 ; \bar{S}_y = 11.11206$$

$$S_{xy} = 237.91556$$

$$R_{xy} = \frac{S_{xy}}{S_x S_y} = \frac{237.91556}{(21.47246)(11.11206)} = \underline{0.9971}$$

| | | | | | | | | | | |
|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|-----------|
| X | 171 | 180 | 193 | 195 | 196 | 200 | 219 | 225 | 228 | 235 |
| Y | 89.3 | 95.6 | 100 | 101 | 101.4 | 104.6 | 114.3 | 118.5 | 118.2 | 122.4 |
| Z | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| | 78 | 79 | 80 | 81 | 82 | 83 | 84 | 85 | 86 | 87 |

X \equiv "Nº de coches aparcados en una determinada facultad"

Y \equiv "Producción industrial de Japón"

Z \equiv "Periodo de años 1978-1987"

El vector de medias y la matriz de covarianzas son:

$$\vec{\mu}_{xyz} = \begin{pmatrix} \mu_x \\ \mu_y \\ \mu_z \end{pmatrix} = \begin{pmatrix} 204.20 \\ 106.21 \\ 5.50 \end{pmatrix}$$


$$\Sigma_{XYZ} = \begin{pmatrix} S_{xx} = S_x^2 & S_{xy} & S_{xz} \\ S_{yx} & S_{yy} = S_y^2 & S_{yz} \\ S_{zx} & S_{zy} & S_{zz} = S_z^2 \end{pmatrix} = \begin{pmatrix} 461.0666 & 237.9155 & 63.7777 \\ 237.9155 & 123.4778 & 32.8722 \\ 63.7777 & 32.8722 & 9.16666 \end{pmatrix}$$

Coeficiente de correlación parcial: $R_{xy,z} = -\frac{A_{12}}{\sqrt{A_{11} \times A_{22}}}$

$$A_{12} = - \begin{vmatrix} 237.9155 & 32.8722 \\ 63.7777 & 9.16666 \end{vmatrix} = -84.3772$$

$$A_{22} = \begin{vmatrix} 461.0666 & 63.7777 \\ 63.7777 & 9.16666 \end{vmatrix} = 158.8457$$

$$A_{11} = \begin{vmatrix} 123.4778 & 32.8722 \\ 32.8722 & 9.16666 \end{vmatrix} = 51.29747$$



$$R_{xy,z} = -\frac{A_{12}}{\sqrt{A_{11} \times A_{22}}} = -\frac{-84.3772}{\sqrt{51.29747 \times 158.8457}} = \underline{0.93473}$$

$$\left. \begin{array}{l} R_{xy}^2 = 0.9971^2 = 0.9942 \\ R_{xy,z}^2 = 0.93473^2 = 0.87372 \end{array} \right\} ; R_{xy,z}^2 < R_{xy}^2 \Rightarrow$$

La interdependencia entre X e Y se debe parcialmente a la influencia de Z.

Aún eliminando los efectos de Z, la relación entre X e Y sigue siendo muy alta desde un punto de vista estadístico, pero esto no significa que tenga sentido fuera de este ámbito, como es obvio para el ejemplo que nos ocupa.

Continuaremos con el ejemplo determinando el modelo: $\hat{x} = a + b_1y + b_2z$

(Recordemos que la matriz de covarianzas \sum_{XYZ} hace referencia a este modelo.)

$$\Sigma_{XYZ} = \begin{pmatrix} 461.0666 & 237.9155 & 63.7777 \\ 237.9155 & 123.4778 & 32.8722 \\ 63.7777 & 32.8722 & 9.16666 \end{pmatrix} \Rightarrow \det \Sigma_{XYZ} = 112.2411$$

$$\Sigma_{\bar{X}} = \begin{pmatrix} 237.9155 & 123.4778 \\ 63.7777 & 32.8722 \end{pmatrix} \Rightarrow \det \Sigma_{\bar{X}} = 51.2974$$

$$b_i = -\frac{A_{1,i+1}}{A_{1,1}} \quad \text{para } i = 1, 2 \Rightarrow \begin{cases} b_1 = -\frac{A_{12}}{A_{11}} = -\frac{-84.3772}{51.29747} = 1.6448 \\ b_2 = -\frac{A_{13}}{A_{11}} = -\frac{-54.3241}{51.29747} = 1.059 \end{cases}$$

$$A_{13} = -\begin{vmatrix} 237.9155 & 123.4778 \\ 63.7777 & 32.8722 \end{vmatrix} = -54.3241$$

$$a = \bar{y} - b_1 \bar{x} - b_2 \bar{z} = 23.3145 \Rightarrow \hat{x} = 23.3145 + 1.6448 y + 1.059 z$$

→ Para $y=130$ y $Z=11$ $\hat{x} = 23.3145 + 1.6448 (130) + 1.059 (11) = 248.7875 \equiv 248$

$$S_{\varepsilon}^2 = \frac{\det \Sigma}{\det \Sigma_{\bar{x}}} = 2.188 \quad R^2 = R_{xy}^2 = \left(\frac{S_{xy}}{S_x S_y} \right)^2 = 1 - \frac{S_{\varepsilon}^2}{S_y^2} = 0.9822$$

Regresión y correlación lineal II

ESCUELA POLITÉCNICA
SUPERIOR DE CÓRDOBA

Universidad de Córdoba

DEPARTAMENTO DE ESTADÍSTICA

