

## Introducción

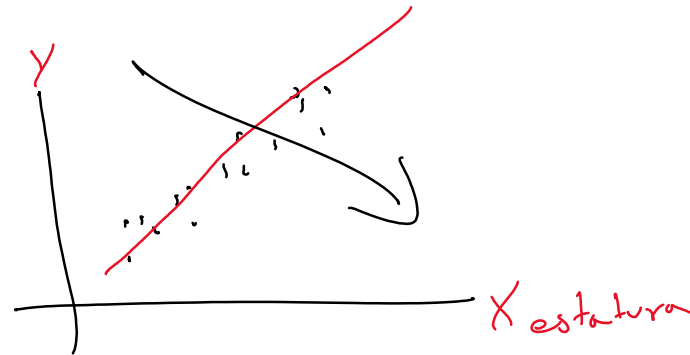
miércoles, 1 de febrero de 2023 7:15 a. m.

$$F_X(z)$$

$$P(X \leq z)$$

$$P_X(z) \rightarrow P(X=z)$$

$$\mu_1 > \mu_2$$



$$r = 0.9$$

- La realidad es que la mayoría de casos de estudio tienen varias variables

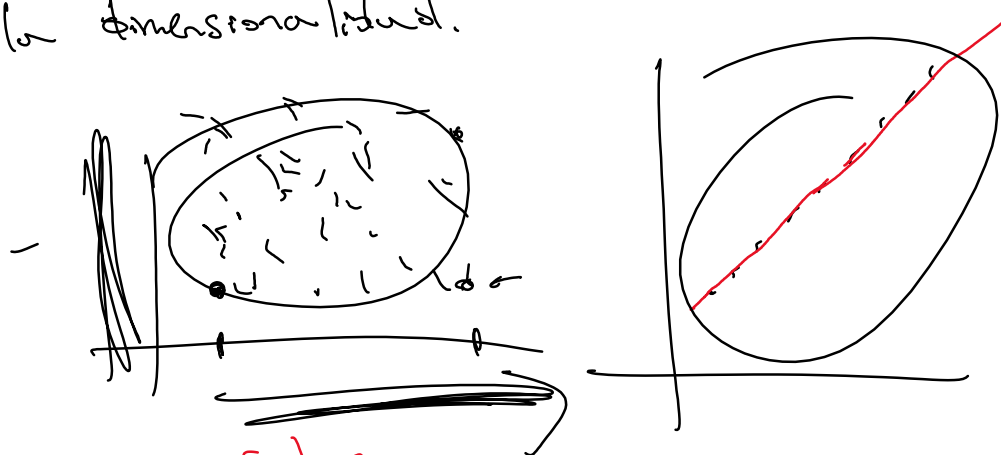
	Var 1	Var 2	...	Var K
Obs 1	—	—		—
...				
Obs n	—	—	—	—

... en  $P_{Y \in 1}$  y  $P_{Y \in 2}$

- Examinaremos los métodos ~~various~~ para estudiar estos casos.
- Utilizaremos R y R markdown.

Principales cosas que se hacen en investigación o en análisis multivariante son:

- Reducción de la dimensionalidad.



Obs 1

Obs n

Var 1 2 ... K

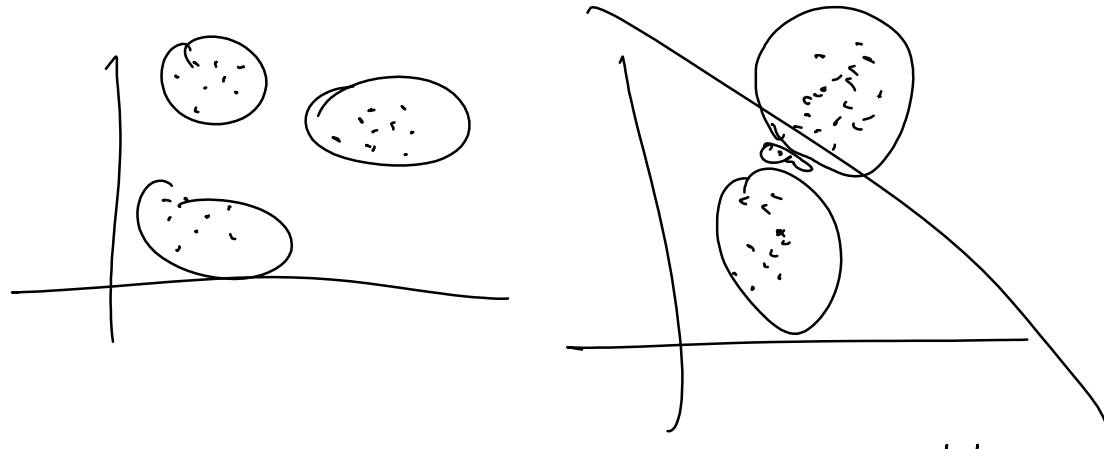
Obs 1

Obs

- Ordenamiento de datos
- Dependencia de variables
- Predicción

$$\underline{y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \dots$$

- Pruebas de hipótesis
- Clasificación.



Vamos a usar notación matricial  $X_{jk}$  es la medida de la  $k$ -ésima variable en la  $j$ -ésima observación

	Var 1	Var K	Var P
Item 1	$x_{11}$	$\vdots$	$x_{1p}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$
Item j	$\vdots$	$x_{jk}$	$\vdots$
$\vdots$	$\vdots$	$\vdots$	$\vdots$
Item n	$x_{n1}$	$x_{nk}$	$x_{np}$

$$X = \begin{bmatrix} x_{11} & \dots & x_{1k} & \dots & x_{1p} \\ \vdots & & \vdots & & \vdots \\ x_{n1} & & x_{nk} & & x_{np} \end{bmatrix}$$

Ejemplo: Variable 1: 42    55    48    58  
 Variable 2: 4        5        4        3

$$x_{11} = 42$$

$$x_{21} = \cancel{4} \quad 5$$

$$X = \begin{bmatrix} 42 & 4 \\ 55 & 5 \\ 48 & 4 \\ 58 & 3 \end{bmatrix}$$

Media  $\bar{x} = \frac{1}{n} \sum x_{ik}$

$$\mu_K = \frac{1}{n} \sum_{j=1}^n \dots$$

Varianza muestral

$$S_K^2 = \frac{1}{n} \sum_{j=1}^n (x_{jk} - \bar{x}_K)^2$$

Matriz de varianzas-covarianzas

$$S = \begin{bmatrix} S_{11} & S_{12} & \dots & S_{1p} \\ S_{21} & S_{22} & \dots & S_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ S_{p1} & S_{p2} & \dots & S_{pp} \end{bmatrix}$$

Simétrica.  
 $S_{21} = S_{12}$

$$S_K = \sqrt{S_K^2} \rightarrow \text{desviación estándar de } K$$

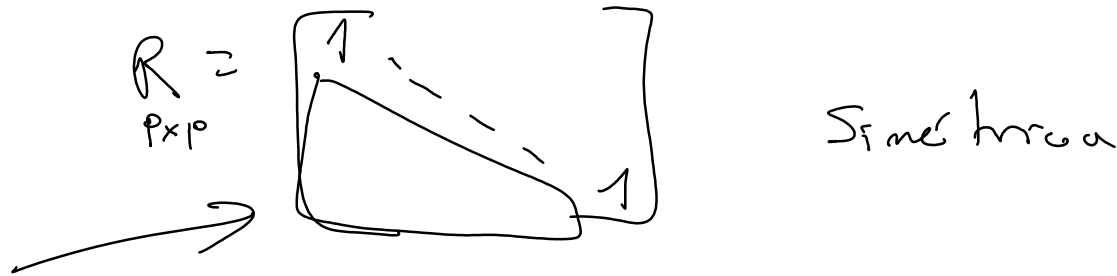
$$S_K^2 = S_{KK} \neq S_K$$

Varianza de K

Coefficiente de correlación muestral

$$r_{iK} = \frac{S_{iK}}{S_i S_K} = \frac{\sum (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_K)}{S_i S_K}$$

$$r_{jk} = \frac{\overline{S_{ji}} \overline{S_{kj}}}{\sqrt{\sum (x_{ji} - \bar{x}_i)^2} \sqrt{\sum (x_{jk} - \bar{x}_k)^2}}$$



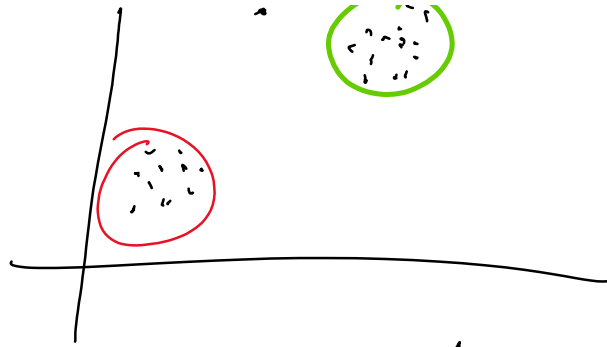
- Es una medida de la asociación lineal entre las variables que **NO** depende de las unidades de las variables.

Vector de medios

$$\bar{X} = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{bmatrix}$$

$p \times 1$

Para estudiar/analizar datos multivariados es muy útil utilizar el concepto de distancias



Cómo calcular la distancia?



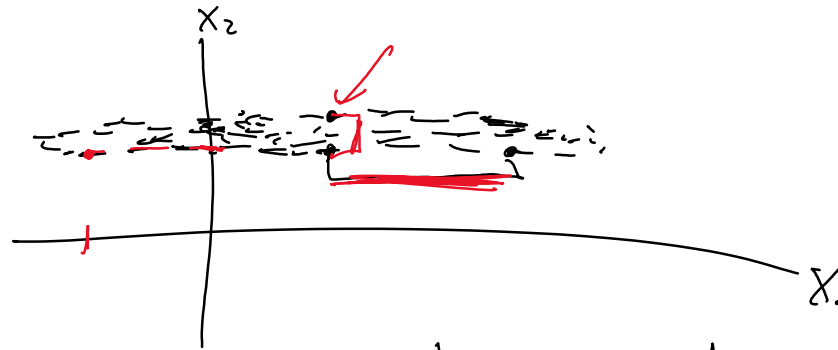
$$\text{Euclid}(O, F) = \sqrt{x_1^2 + x_2^2}$$

$$\text{Euclid}(O, F) = \sqrt{x_1^2 + x_2^2 + \dots + x_p^2}$$

Distancia euclídea entre dos puntos

$$F = (x_1, x_2, \dots, x_p) \quad Q = (y_1, y_2, \dots, y_p)$$

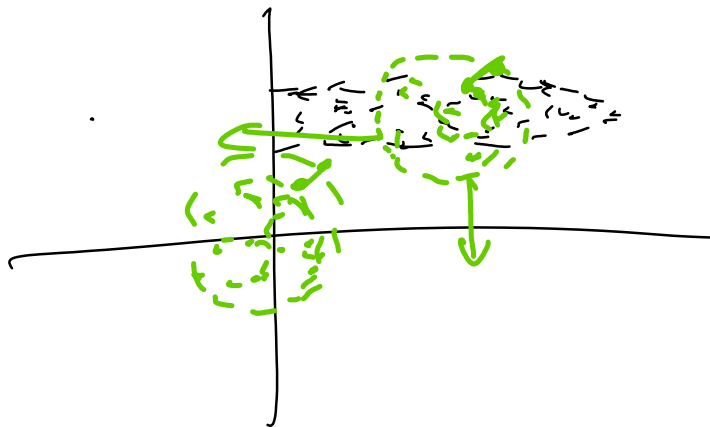
$$d(F, Q) = \sqrt{(x_1 - y_1)^2 + \dots + (x_p - y_p)^2}$$



Para la mayoría de aplicaciones, la distancia euclidiana no es apropiada

Calcularemos la distancia estadística

$$X_1^* = \frac{X_1 - \bar{X}_1}{\sqrt{S_{11}}} \quad X_2^* = \frac{X_2}{\sqrt{S_{22}}}$$



Distancia estadística de un punto  $F$  y el origen



$$d(O, F) = \sqrt{(x_1^*)^2 + (x_2)^2}$$

$$= \sqrt{\frac{x_1^2}{S_{11}} + \frac{x_2^2}{S_{22}}}$$

$$d(F, Q) = \sqrt{\frac{(x_1 - y_1)^2}{S_{11}} + \dots + \frac{(x_p - y_p)^2}{S_{pp}}}$$

$$= \sqrt{(x_1^* - y_1^*)^2 + \dots + (x_p^* - y_p^*)^2}$$

