



# **GEOMETRÍA MUESTRAL**



## Observación Multivariada

Una observación multivariada es una colección de mediciones sobre 'p' variables medida sobre el mismo objeto o ensayo:

$$\mathbf{X}_{n \times p} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}$$

Diagram illustrating the structure of a multivariate observation matrix  $\mathbf{X}_{n \times p}$ . The matrix is shown with rows representing observations and columns representing variables. The first row is labeled "1ra observación" and the last row is labeled "n-ésima observación". The columns are labeled  $y_1$ ,  $y_2$ , and  $y_p$  at the bottom, with arrows pointing to the corresponding columns in the matrix.



## NOTACIÓN

$X_{n \times p}$  : Matriz de datos para una muestra de tamaño 'n' para 'p' variables

$X_{ij}$  : Un dato para la variable 'j' del individuo o unidad muestral 'i'.

$X'_i$  : i-ésima fila.

$X_{(j)}$  : j-ésima columna, correspondiente a la j-ésima variable.



## NOTACIÓN (...continuación)

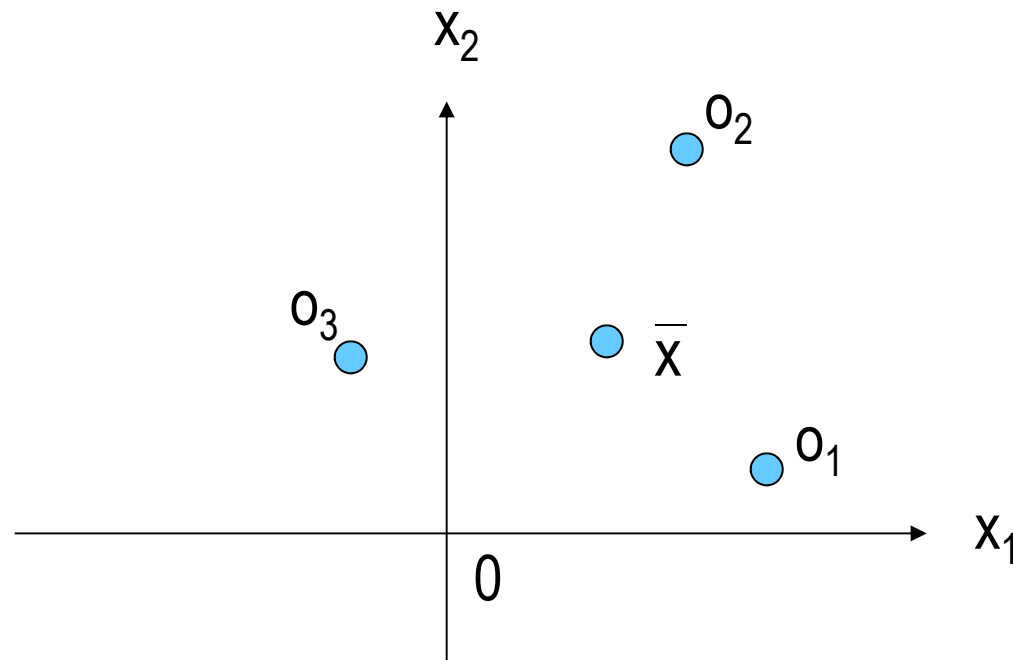
$$X_{n \times p} = \begin{pmatrix} X'_1 \\ X'_2 \\ \vdots \\ X'_n \end{pmatrix} = \begin{pmatrix} X_{(1)}, X_{(2)}, \dots, X_{(p)} \end{pmatrix}$$

$$\text{donde, } X_i = \begin{pmatrix} X_{i1} \\ X_{i2} \\ \vdots \\ X_{ip} \end{pmatrix} \quad (i=1, \dots, n) \quad X_{(j)} = \begin{pmatrix} X_{1j} \\ X_{2j} \\ \vdots \\ X_{nj} \end{pmatrix} \quad (j=1, \dots, p)$$



## Caso $n=3$ y $p=2$

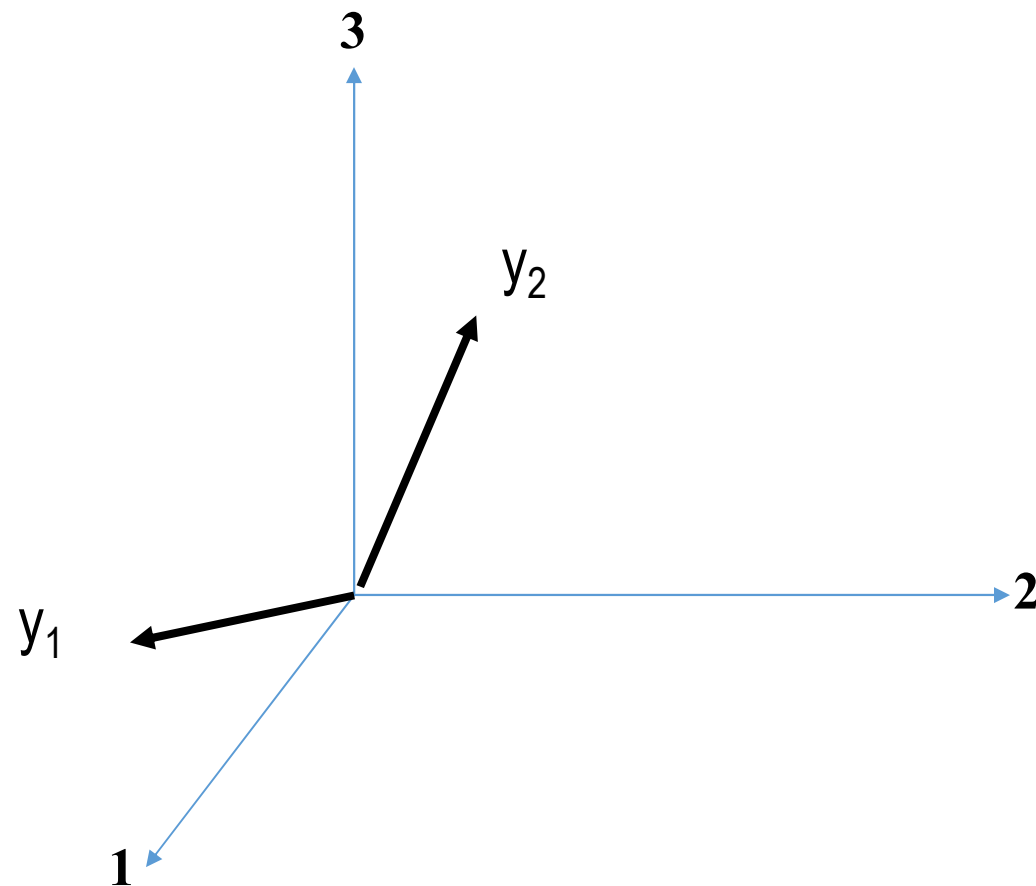
Ploteo de puntos (observaciones) para una matriz  $X$





## Caso $n=3$ y $p=2$

### Gráfico de variables para la matriz $X$

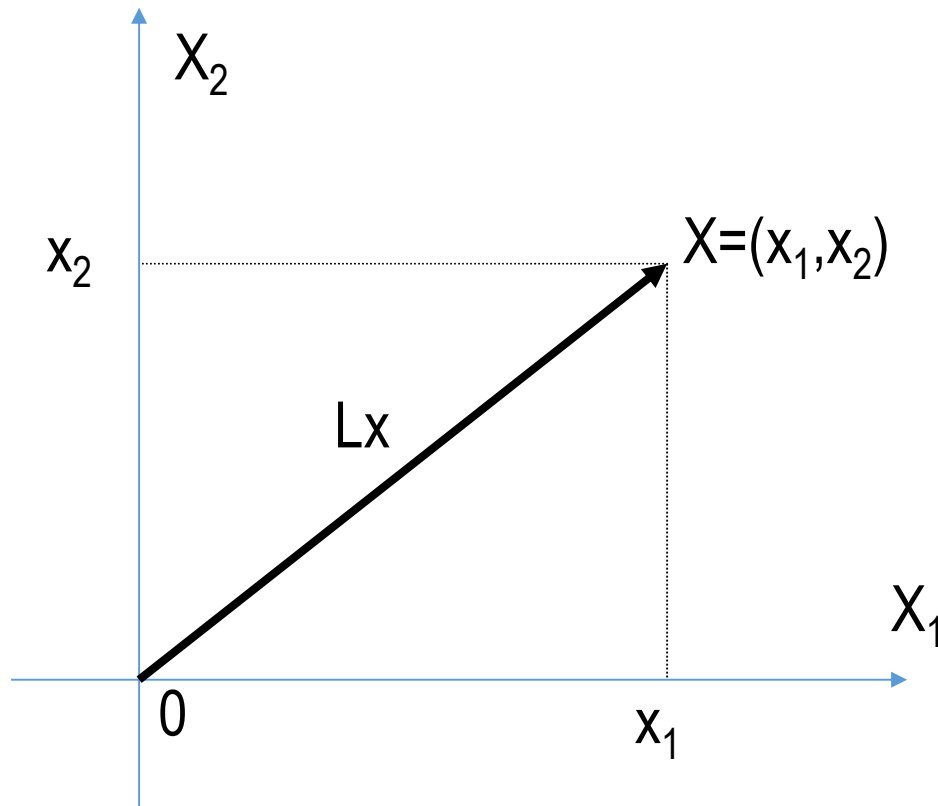




# VECTORES ALEATORIOS



## Longitud o Norma de un Vector



$$L_x = \sqrt{X_1^2 + X_2^2}$$

$$L_x = \sqrt{X'X}$$

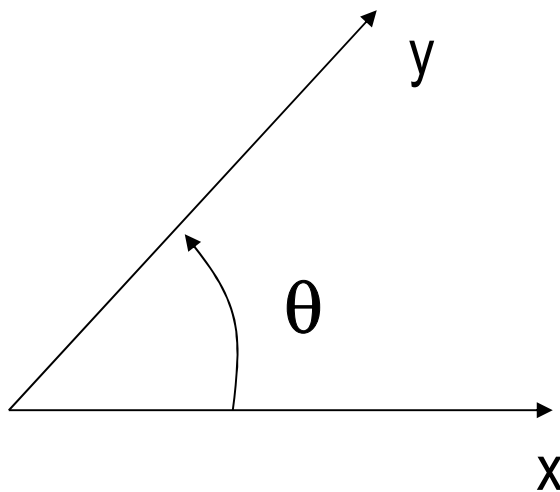
Caso General:

$$X = (x_1, x_2, \dots, x_p)$$





## Ángulo entre Dos Vectores



$$\cos \theta = \frac{x_1 y_1 + x_2 y_2 + \dots + x_p y_p}{L_x L_y}$$

$$\cos \theta = \frac{X'Y}{\sqrt{X'X} \sqrt{Y'Y}}$$



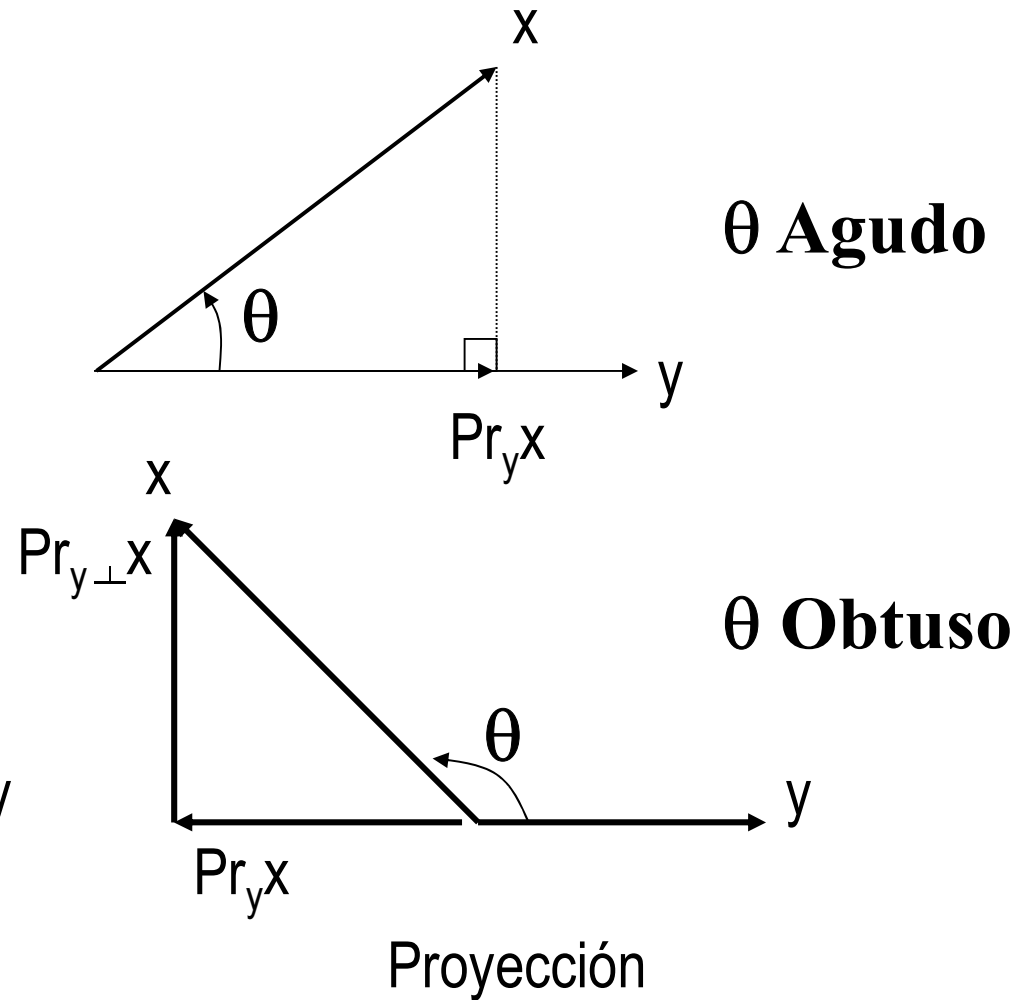
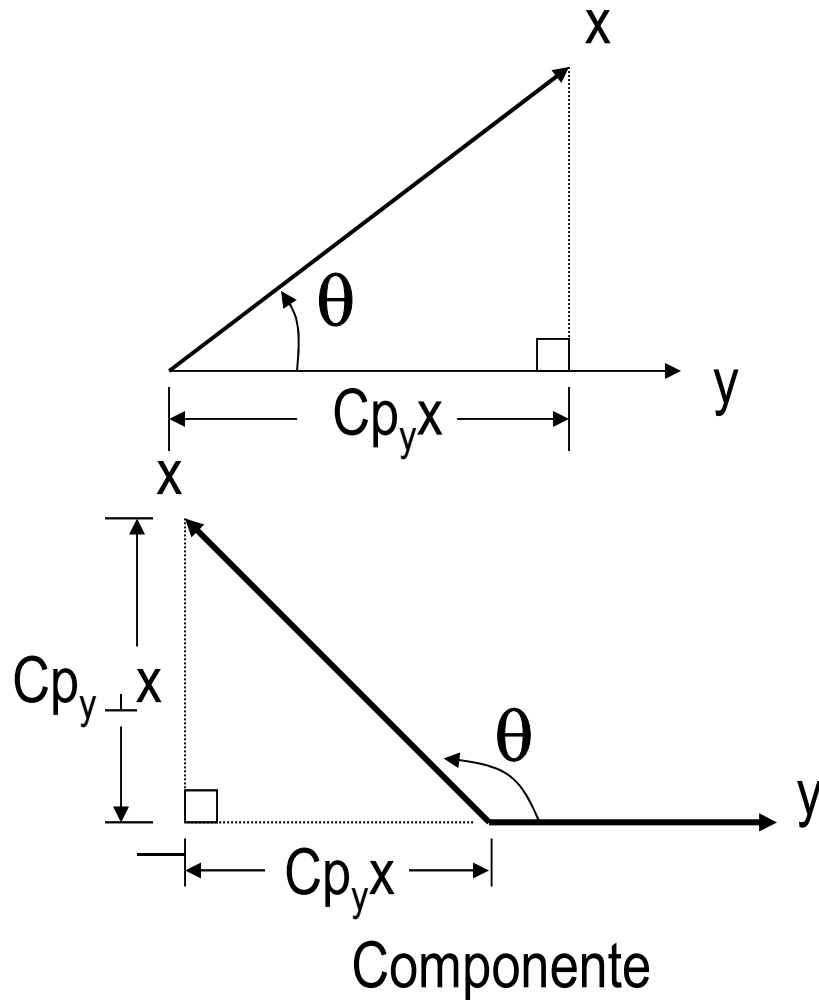
## Componente y Proyección de X sobre Y

La Proyección Ortogonal de X sobre Y ( $Pr_y x$ ) es el vector cuya dirección está dada por la dirección y sentido del vector Y, y la longitud por la Componente de X sobre Y ( $Cp_y x$ )

$$Pr_y x = \frac{x^l y}{L^2 y} \cdot y = \frac{x^l y}{L_y} \cdot \frac{y}{L_y} = Cp_y x \cdot \frac{y}{L_y}$$

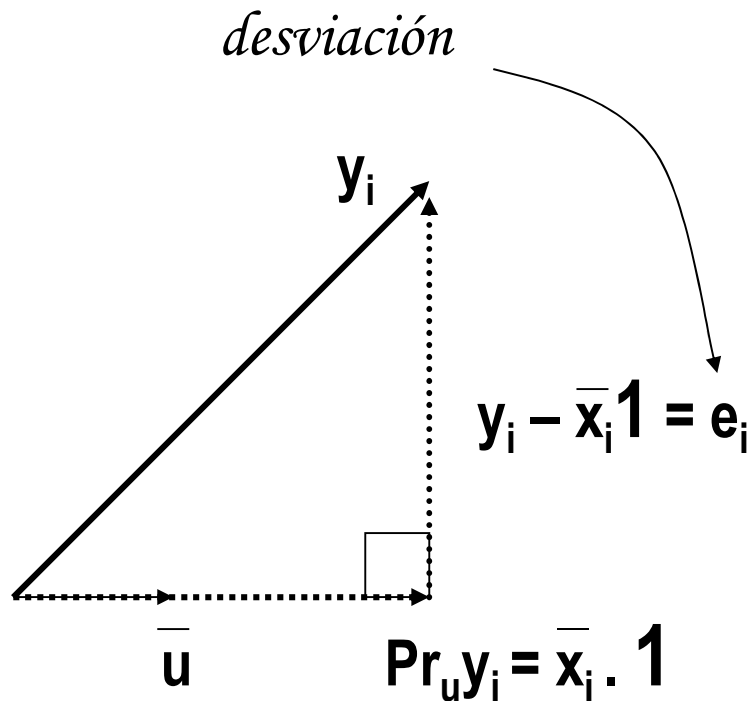


## CASOS





## Media Muestral



Sea

$$y_i = (x_{i1}, x_{i2}, \dots, x_{in}) \text{ y } \mathbf{1} = (1, 1, \dots, 1)_{n \times 1}$$

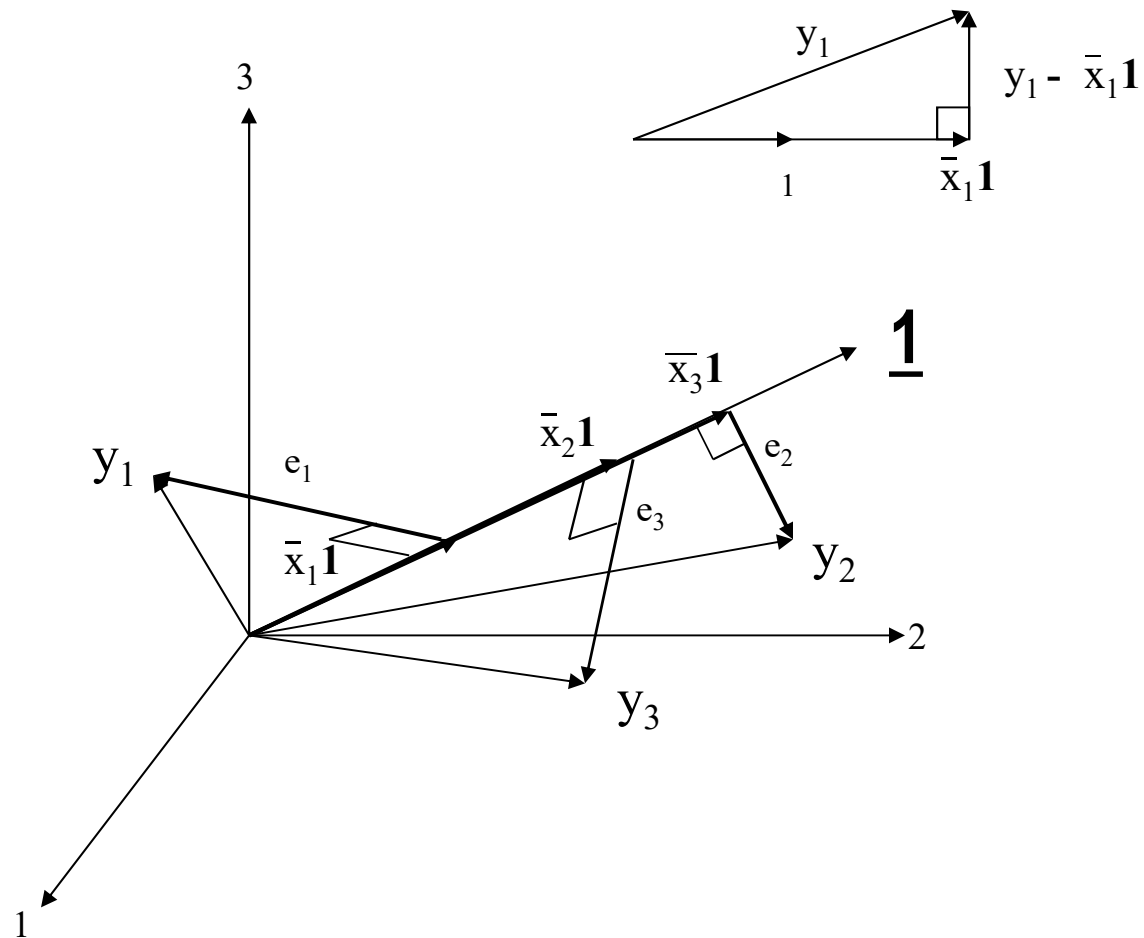
$$\text{Pr}_u y_i = \frac{y_i^T \bar{u}}{|\bar{u}|^2} \cdot \bar{u} \text{ donde } \bar{u} = \frac{\mathbf{1} \cdot \mathbf{1}}{\sqrt{n}}$$

$$\text{Pr}_u y_i = y_i^T \bar{u} \cdot \bar{u} = \bar{x}_i \cdot \mathbf{1}$$

*componente  
de medias*



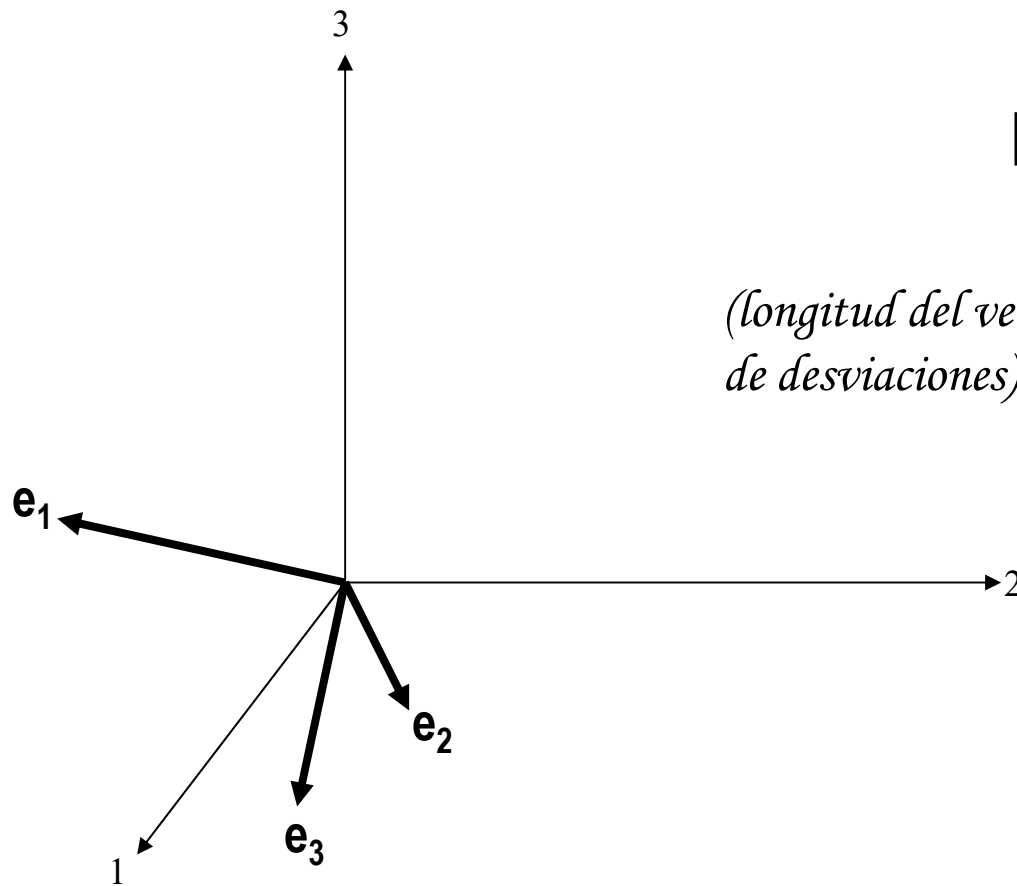
## Caso para $p=3$ y $n=3$





Caso para  $p=3$  y  $n=3$

## Graficando las Desviaciones



$$L^2_{ei} = e_i^T e_i = \sum_j (x_{ij} - \bar{x}_i)^2$$

*(longitud del vector  
de desviaciones)<sup>2</sup> = (suma de cuadrados de  
las desviaciones)*

**Longitud más grande del  
vector de desviaciones  
representa mayor  
variabilidad que vectores  
más pequeños**



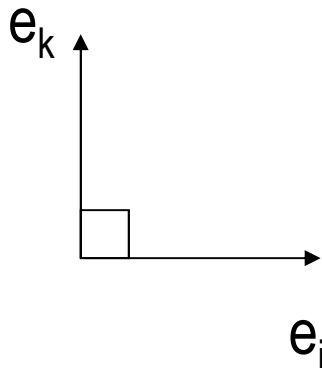
## Coefficiente de Correlación

El coseno del ángulo formado por dos vectores de desviaciones es el coeficiente de correlación

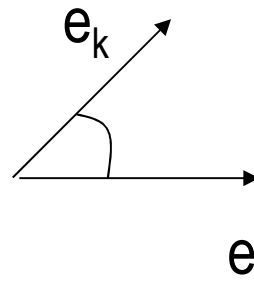
$$\mathbf{e}_i^T \mathbf{e}_k = L_{ei} L_{ek} \cos(\theta_{ik})$$

$$\cos(\theta_{ik}) = \frac{\mathbf{e}_i^T \mathbf{e}_k}{L_{ei} L_{ek}}$$

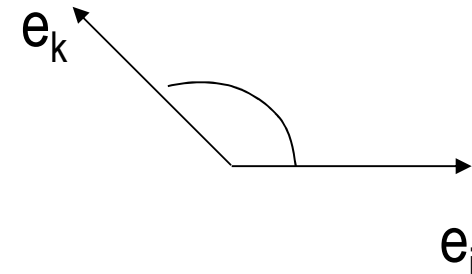
$$\cos(\theta_{ik}) = \frac{S_{ik}}{\sqrt{S_{ii}} \sqrt{S_{kk}}}$$



$$\cos(\theta_{ik}) = 0$$



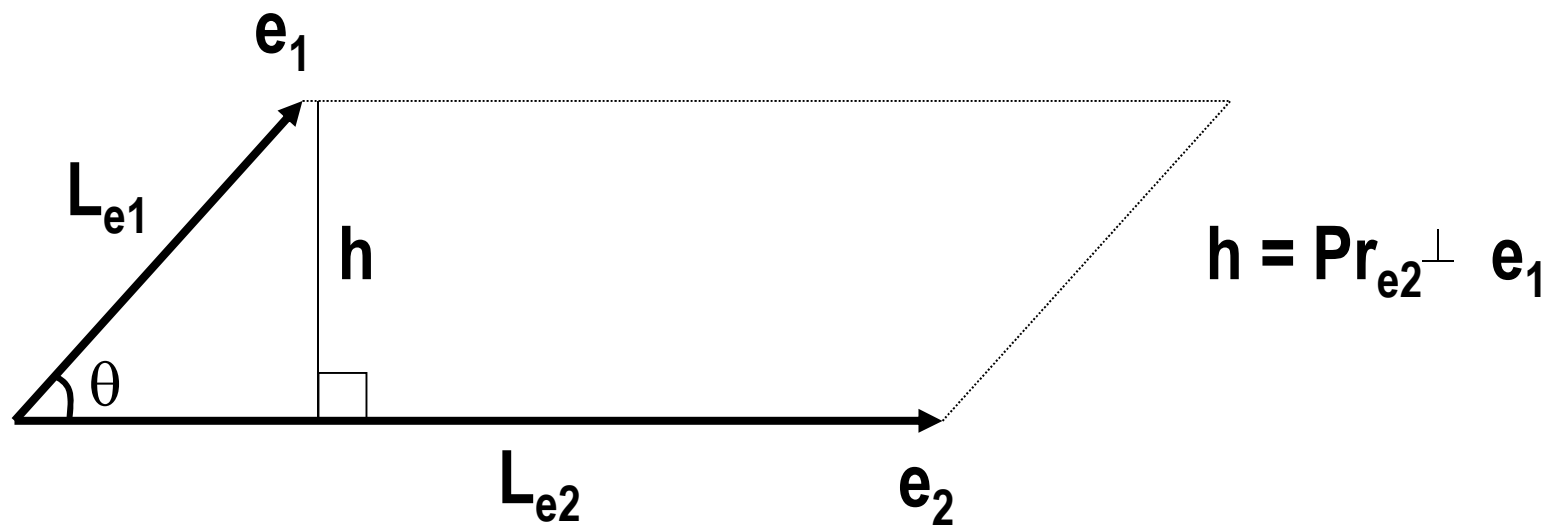
$$0 < \cos(\theta_{ik}) < 1$$



$$\cos(\theta_{ik}) < 0$$



## Varianza Generalizada



$$h = L_{e1} \sin(\theta)$$

$$\text{Area} = L_{e1} L_{e2} \sqrt{1 - \cos^2(\theta)}$$





$$L_{e1} = \sqrt{\sum_{j=1}^n (x_{1j} - \bar{x}_1)^2} = \sqrt{(n-1)s_{11}}$$

$$L_{e2} = \sqrt{\sum_{j=1}^n (x_{2j} - \bar{x}_2)^2} = \sqrt{(n-1)s_{22}}$$

$$\cos(\theta) = r_{12}$$

Luego,

$$\text{Area} = (n-1) \sqrt{s_{11}} \sqrt{s_{22}} \sqrt{1 - r_{12}^2}$$

$$|S| = (\text{Area})^2 / (n-1)^2$$

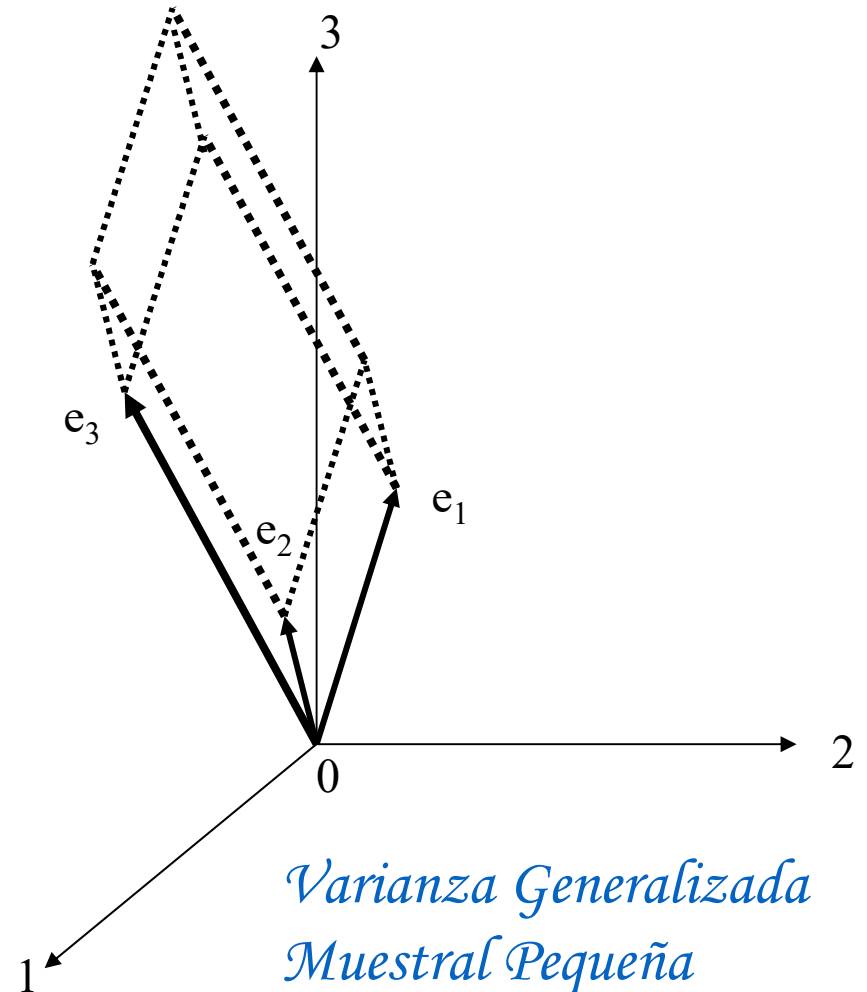
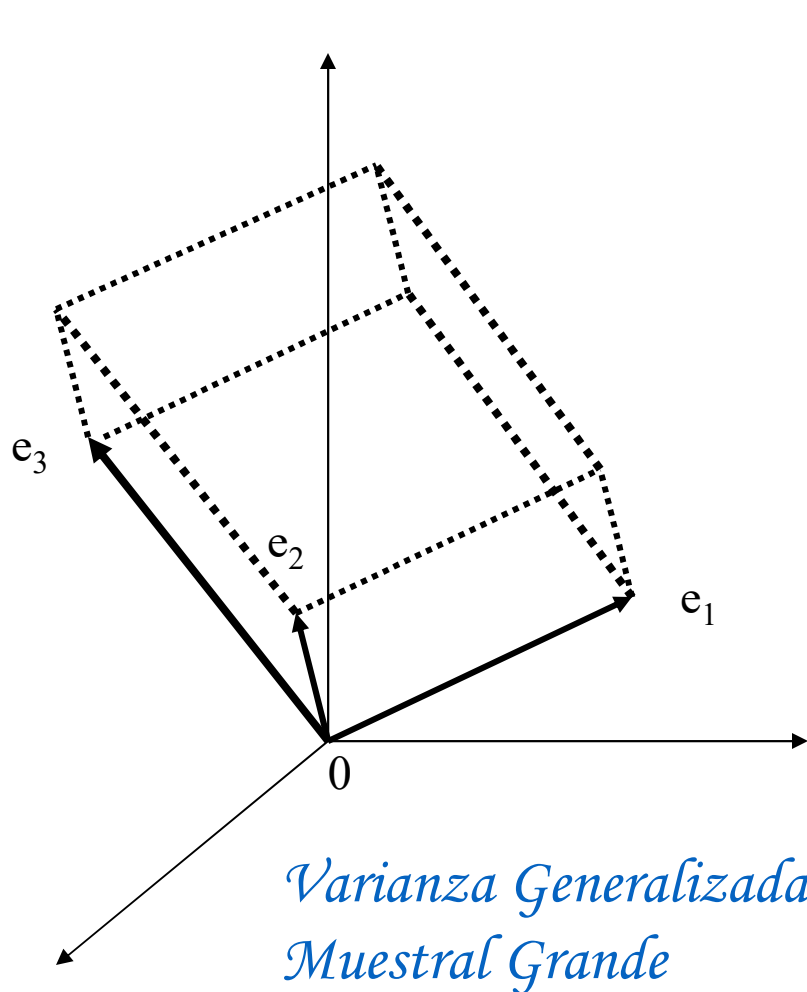


Generalizando para un espacio  $n$ -dimensional y  $p$  vectores de desviaciones,

$$|S| = (n-1)^{-p} (\text{Volumen})^2$$

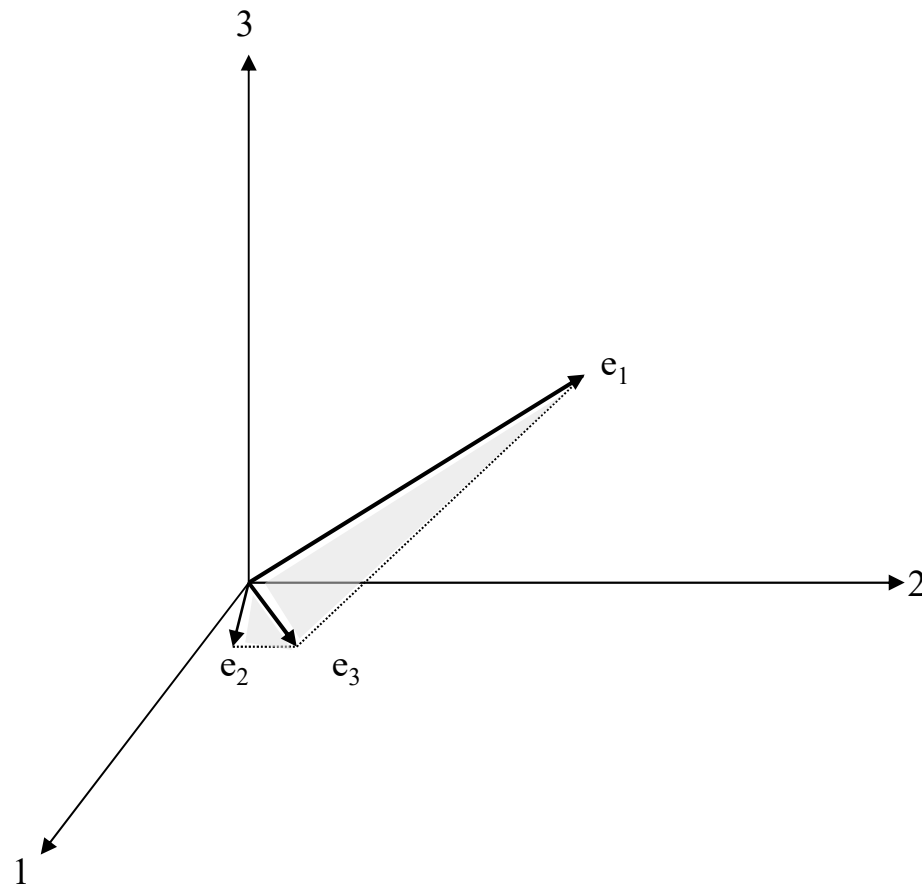


Caso:  $p=3$



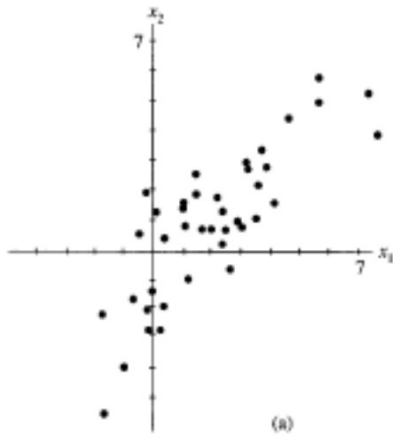


## Caso en que el volumen tri-dimensional es cero $|S|=0$





## Interpretación de la varianza generalizada

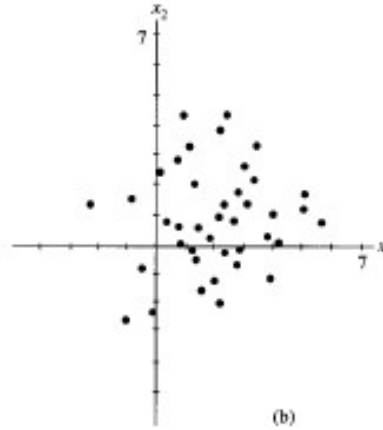


$$S = \begin{pmatrix} 5 & 4 \\ 4 & 5 \end{pmatrix}$$

$$r = 0.8$$

$$\bar{x}' = [2, 1]$$

$$|S| = 9$$

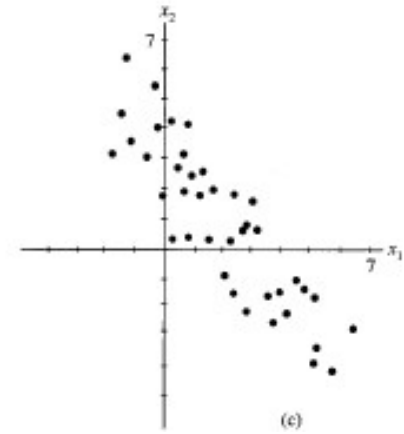


$$S = \begin{pmatrix} 3 & 0 \\ 0 & 3 \end{pmatrix}$$

$$r = 0$$

$$\bar{x}' = [2, 1]$$

$$|S| = 9$$



$$S = \begin{pmatrix} 5 & -4 \\ -4 & 5 \end{pmatrix}$$

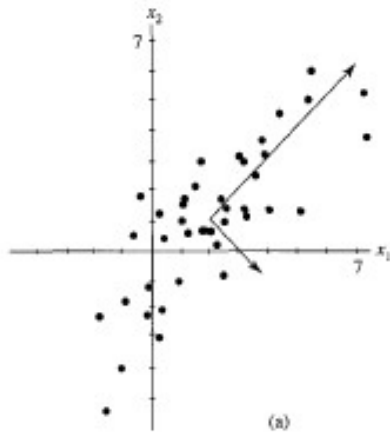
$$r = -0.8$$

$$\bar{x}' = [2, 1]$$

$$|S| = 9$$



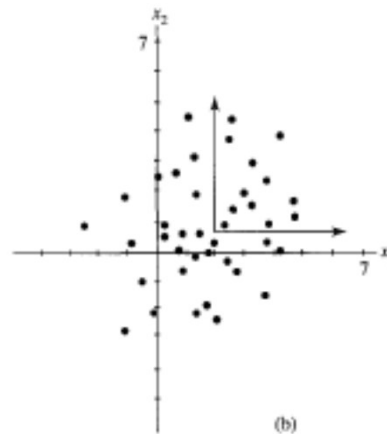
## Diferentes estructuras de correlación no son detectadas por $|S|$



$$\lambda_1 = 9 \quad e_1 = \left[ 1/\sqrt{2}, 1/\sqrt{2} \right]$$

$$\lambda_2 = 1 \quad e_2 = \left[ 1/\sqrt{2}, -1/\sqrt{2} \right]$$

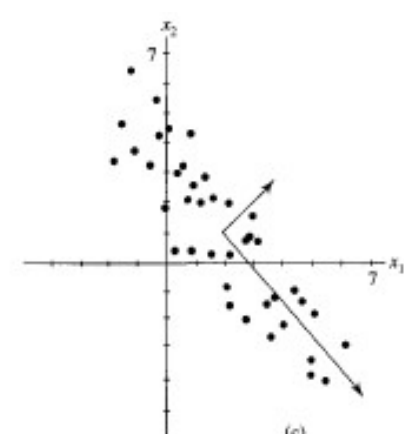
$$\bar{x}' = \begin{bmatrix} 2, 1 \end{bmatrix}$$



$$\lambda_1 = 3 \quad e_1 = \begin{bmatrix} 1, 0 \end{bmatrix}$$

$$\lambda_2 = 3 \quad e_2 = \begin{bmatrix} 0, 1 \end{bmatrix}$$

$$\bar{x}' = \begin{bmatrix} 2, 1 \end{bmatrix}$$



$$\lambda_1 = 9 \quad e_1 = \left[ 1/\sqrt{2}, -1/\sqrt{2} \right]$$

$$\lambda_2 = 1 \quad e_2 = \left[ 1/\sqrt{2}, 1/\sqrt{2} \right]$$

$$\bar{x}' = \begin{bmatrix} 2, 1 \end{bmatrix}$$



**Ejemplo 3.3 (\*) *Descomposición de un vector en su media y desviaciones centrales***

$$X = \begin{pmatrix} 4 & 1 \\ -1 & 3 \\ 3 & 5 \end{pmatrix}$$

En este caso,

$$\bar{x}_1 = (4 - 1 + 3)/3 = 2 \text{ and } \bar{x}_2 = (1 + 3 + 5)/3 = 3, \text{ so}$$

$$\bar{x}_1 \mathbf{1} = 2 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix} \quad \bar{x}_2 \mathbf{1} = 3 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 3 \\ 3 \\ 3 \end{bmatrix}$$

Consecuentemente,

$$\mathbf{d}_1 = \mathbf{y}_1 - \bar{x}_1 \mathbf{1} = \begin{bmatrix} 4 \\ -1 \\ 3 \end{bmatrix} - \begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix} = \begin{bmatrix} 2 \\ -3 \\ 1 \end{bmatrix}$$

(\*) Johnson y Wichern (2002) Applied Multivariate Statistical Analysis



## FSM71 - ESTADÍSTICA MULTIVARIADA

### 1. GEOMETRÍA DE LA MUESTRA Y MUESTREO MULTIVARIADO

y,

$$\mathbf{d}_2 = \mathbf{y}_2 - \bar{\mathbf{x}}_2 \mathbf{1} = \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix} - \begin{bmatrix} 3 \\ 3 \\ 3 \end{bmatrix} = \begin{bmatrix} -2 \\ 0 \\ 2 \end{bmatrix}$$

Se comprueba que  $\bar{\mathbf{x}}_1 \mathbf{1}$  y  $\mathbf{d}_1$  son perpendiculares

$$(\bar{\mathbf{x}}_1 \mathbf{1})'(\mathbf{y}_1 - \bar{\mathbf{x}}_1 \mathbf{1}) = [2 \quad 2 \quad 2] \begin{bmatrix} 2 \\ -3 \\ 1 \end{bmatrix} = 4 - 6 + 2 = 0$$

El mismo resultado se obtiene para  $\bar{\mathbf{x}}_2 \mathbf{1}$  y  $\mathbf{d}_2$ .

La descomposición de  $\mathbf{y}_1$  e  $\mathbf{y}_2$  es:

$$\begin{aligned} \mathbf{y}_1 &= \begin{bmatrix} 4 \\ -1 \\ 3 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix} + \begin{bmatrix} 2 \\ -3 \\ 1 \end{bmatrix} \\ \mathbf{y}_2 &= \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix} = \begin{bmatrix} 3 \\ 3 \\ 3 \end{bmatrix} + \begin{bmatrix} -2 \\ 0 \\ 2 \end{bmatrix} \end{aligned}$$





## FSM71 - ESTADÍSTICA MULTIVARIADA

### 1. GEOMETRÍA DE LA MUESTRA Y MUESTREO MULTIVARIADO

Para cualquiera dos vectores de desviaciones  $\mathbf{d}_i$  y  $\mathbf{d}_k$  :

$$\mathbf{d}_i' \mathbf{d}_k = \sum_{j=1}^n (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k)$$

Si  $\theta_{ik}$  es el ángulo formado por los vectores  $\bar{\mathbf{d}}_i$  y  $\bar{\mathbf{d}}_k$  :

$$\mathbf{d}_i' \mathbf{d}_k = L_{\mathbf{d}_i} L_{\mathbf{d}_k} \cos(\theta_{ik})$$

O, equivalentemente con las observaciones muestrales:

$$\sum_{j=1}^n (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k) = \sqrt{\sum_{j=1}^n (x_{ji} - \bar{x}_i)^2} \sqrt{\sum_{j=1}^n (x_{jk} - \bar{x}_k)^2} \cos(\theta_{ik})$$

El coseno del ángulo es el coeficiente de correlación muestral entre  $\mathbf{d}_i$  y  $\mathbf{d}_k$ :

$$r_{ik} = \frac{s_{ik}}{\sqrt{s_{ii}} \sqrt{s_{kk}}} = \cos(\theta_{ik})$$



Calculando la matriz varianza covarianza **S<sub>n</sub>** y de correlaciones **R**

Se sabe que:

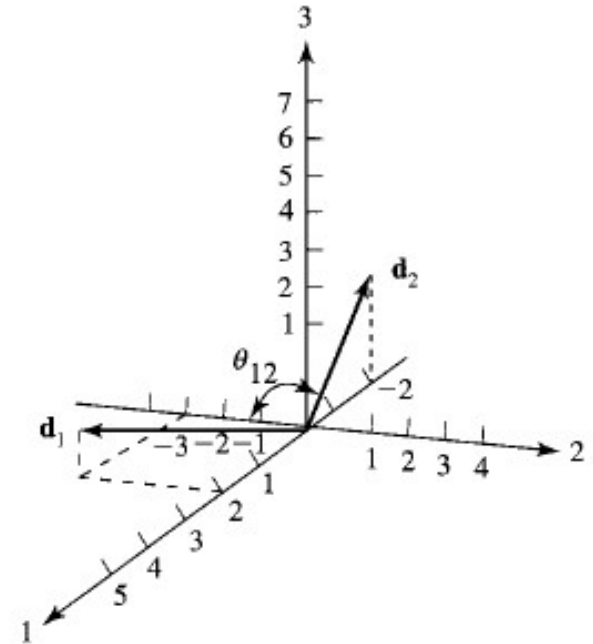
$$\mathbf{d}_1 = \begin{bmatrix} 2 \\ -3 \\ 1 \end{bmatrix} \quad \text{and} \quad \mathbf{d}_2 = \begin{bmatrix} -2 \\ 0 \\ 2 \end{bmatrix}$$

Veamos los componentes de la matriz varianza-covarianza:

$$\mathbf{d}_1' \mathbf{d}_1 = [2 \quad -3 \quad 1] \begin{bmatrix} 2 \\ -3 \\ 1 \end{bmatrix} = 14 = 3s_{11} \quad \longrightarrow \quad s_{11} = 14 / 3$$

$$\mathbf{d}_2' \mathbf{d}_2 = [-2 \quad 0 \quad 2] \begin{bmatrix} -2 \\ 0 \\ 2 \end{bmatrix} = 8 = 3s_{22} \quad \longrightarrow \quad s_{22} = 8 / 3$$

$$\mathbf{d}_1' \mathbf{d}_2 = [2 \quad -3 \quad 1] \begin{bmatrix} -2 \\ 0 \\ 2 \end{bmatrix} = -2 = 3s_{12} \quad \longrightarrow \quad s_{12} = -2 / 3$$





## FSM71 - ESTADÍSTICA MULTIVARIADA

### 1. GEOMETRÍA DE LA MUESTRA Y MUESTREO MULTIVARIADO

Consecuentemente,

$$r_{12} = \frac{s_{12}}{\sqrt{s_{11}} \sqrt{s_{22}}} = \frac{-\frac{2}{3}}{\sqrt{\frac{14}{3}} \sqrt{\frac{8}{3}}} = -.189 \quad \longrightarrow \quad \theta_{12} = 100.65^\circ$$

La matriz varianza-covarianza y de correlación muestral serían:

$$\mathbf{S}_n = \begin{bmatrix} \frac{14}{3} & -\frac{2}{3} \\ -\frac{2}{3} & \frac{8}{3} \end{bmatrix}, \quad \mathbf{R} = \begin{bmatrix} 1 & -.189 \\ -.189 & 1 \end{bmatrix}$$



## INTERPRETACIÓN GEOMÉTRICA DE LA MUESTRA

1. La proyección de la columna  $x(j)$  de la matriz  $X$  sobre el vector equiangular  $\mathbf{1}$  es el vector  $\bar{x}_j \mathbf{1}$ . Éste tiene norma  $\sqrt{n} |\bar{x}_j|$ . Así, la  $j$ -ésima media muestral,  $\bar{x}_j$ , se relaciona con la longitud de la proyección de  $x(j)$  sobre  $\mathbf{1}$  o la componente de  $x(j)$  sobre  $\mathbf{1}$ .
2. La información contenida en  $\mathbf{S}_n$  se obtuvo de los vectores de desviaciones  $\mathbf{d}_j = \mathbf{y}_j - \bar{x}_{(j)} \mathbf{1} = \begin{pmatrix} x_{1j} - \bar{x}_{(j)}, x_{2j} - \bar{x}_{(j)}, \dots, x_{nj} - \bar{x}_{(j)} \end{pmatrix}$ . El cuadrado de la longitud  $\mathbf{d}_j$  es  $ns_{jj}$  y el producto interno entre  $\mathbf{d}_i$  y  $\mathbf{d}_j$  es  $ns_{ij}$ .
3. La correlación muestral  $r_{ik}$  es el coseno del ángulo formado por  $\mathbf{d}_i$  y  $\mathbf{d}_k$ .