

Una muestra de 1 obsn. multivari. es el conjunto de mediciones de 'p' variables distintas en el mismo sujeto ó ensayo. Si se obtienen 'n' observaciones, los datos se resumen en la matriz

$$\underline{X}_{(n \times p)} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \ddots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix} \leftarrow \text{1 obs. 4 v. con } p \text{ componentes}\right.$$

### 1) Representación Geométrica de la muestra

Re-escribamos  $\underline{X}$  así:  $\underline{X}_{(n \times p)} = [\underline{Y}_1 | \underline{Y}_2 | \dots | \underline{Y}_p]$ , donde

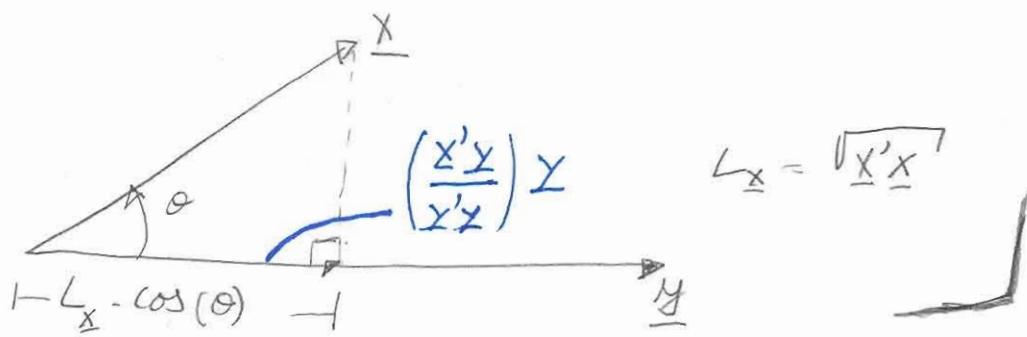
$\underline{Y}_i = [x_{1i} \ x_{2i} \ \dots \ x_{ni}]$  son las 'n' mediciones hechas para los 'n' individuos, respecto de la ' $i$ '-ésima variable

Aparentado: La proyección (o sombra) del vector  $\underline{X}$  sobre el vector  $\underline{Y}$

$$\text{es: } \frac{(\underline{X}' \underline{X})}{\underline{Y}' \underline{Y}} \underline{Y} = \frac{(\underline{X}' \underline{Y})}{\underline{Y}' \underline{Y}} \frac{1}{\|\underline{Y}\|} \underline{Y}, \text{ donde } L_{\underline{Y}} = \sqrt{\underline{Y}' \underline{Y}}$$

$$\cos(\theta) = \frac{\underline{X}' \underline{Y}}{\|\underline{X}\| \|\underline{Y}\|}$$

(ver p. 53 cap. 2)



Consideren el vector  $\frac{1}{\sqrt{m}} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}_{(m \times 1)} = \frac{1}{\sqrt{m}} \mathbf{1}$ . Entonces,

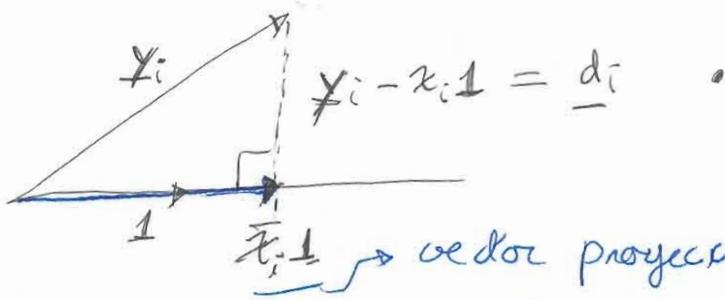
$\frac{1}{\sqrt{m}} \mathbf{1} \cdot \mathbf{1} = 1$ . La proyección del vector  $\underline{x_i}$  en el vector

unitario  $\left(\frac{1}{\sqrt{m}}\right) \mathbf{1}$  es igual a:  $\underline{x_i}' \left(\frac{1}{\sqrt{m}} \mathbf{1}\right) \cdot \frac{1}{\sqrt{m}} \mathbf{1}$ , y

$$\underline{x_i}' \left(\frac{1}{\sqrt{m}} \mathbf{1}\right) \frac{1}{\sqrt{m}} \mathbf{1} = \frac{x_{1i} + x_{2i} + \dots + x_{ni}}{m} \cdot \mathbf{1} = \bar{x}_i \mathbf{1}.$$

Esto es, la

media en la muestra de la variable ' $i'$ ' ( $\bar{x}_i = \underline{x}' \mathbf{1} / m$ ), es la cantidad por la cual hay que multiplicar el vector  $\mathbf{1}$  para obtener la proyección de  $\underline{x_i}$  sobre la línea determinada por  $\mathbf{1}$ :



Entonces,  $\underline{d_i} = \begin{bmatrix} x_{1i} - \bar{x}_i \\ x_{2i} - \bar{x}_i \\ \vdots \\ x_{ni} - \bar{x}_i \end{bmatrix}$  es entonces el vector que contiene las desviaciones de la media  $\bar{x}_i$  para cada observación, para la variable ' $i'$ .

La varianza de la muestra para la variable 'i' es proporcional al (largo)<sup>2</sup> de  $\underline{d}_i$ : (3)

$$(\text{largo})^2 \text{ de } \underline{d}_i : \quad (\underline{L}_{\underline{d}_i})^2 = \underline{d}_i \cdot \underline{d}_i = \sum_{j=1}^m (x_{ji} - \bar{x}_i)^2 \quad \text{ec. (**)}$$

Asimismo, el largo de dichas desviaciones es proporcional a la desviación estandar de la muestra (para la variable 'i').

Para dos vectores de desviaciones  $\underline{d}_i$  y  $\underline{d}_k$ ,

$$\underline{d}_i \cdot \underline{d}_k = \sum_{j=1}^m (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k) \quad \text{ec. (***)}$$

Si  $\theta_{ik}$  denota el ángulo entre  $\underline{d}_i$  y  $\underline{d}_k$ , entonces

$$\underline{d}_i \cdot \underline{d}_k = L_{\underline{d}_i} L_{\underline{d}_k} \cos(\theta_{ik}), \quad \text{y usando las ecs. (*) y (**) :}$$

$$\sum_{j=1}^m (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k) = \underbrace{\sqrt{\sum_{j=1}^m (x_{ji} - \bar{x}_i)^2}}_{= S_{ik}} \cdot \underbrace{\sqrt{\sum_{j=1}^m (x_{jk} - \bar{x}_k)^2}}_{= S_{kk}} \cdot \cos(\theta_{ik})$$

Por lo que :

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

El coseno del ángulo entre  $\underline{d}_i$  y  $\underline{d}_k$  es la correlación entre la muestra 'i' y 'k'. Entonces, si dos vectores desviación tienen casi la misma orientación (i.e.  $\theta_{ik} \approx 0$ ), entonces la correlación va a ser casi 1. Si dos vectores son perpendiculares, la correlación de la muestra sera 0 y si están orientados de manera opuesta, -1.

## 2) Muestras aleatorias y momentos de la muestra.

(4)

Antes del muestreo, asumimos que cada uno de los vectores  $\underline{X}_j$  que tendrá mediciones en 'p' variables es un vector aleatorio, y que juntas, los 'n' vectores conforman una matriz aleatoria.

$$\underline{X}_{(n \times p)} = \begin{bmatrix} \underline{X}'_1 \\ \underline{X}'_2 \\ \vdots \\ \underline{X}'_n \end{bmatrix} \quad \text{eq (***)}$$

Def.: Muestra aleatoria: Si los vectores  $\underline{X}'_1, \underline{X}'_2, \dots, \underline{X}'_n$  en (\*\*\*)  
representan observaciones independientes de una misma distribución conjunta con densidad  $f(\underline{X}) = f(x_1, x_2, \dots, x_p)$ , entonces  
se dice que  $\underline{X}_1, \underline{X}_2, \dots, \underline{X}_n$  forman una muestra aleatoria de  $f(\underline{X})$ . La distribución conjunta de la muestra está dada por

$$f(\underline{X}_1)f(\underline{X}_2)\dots f(\underline{X}_n), \text{ donde } f(\underline{X}_j) = f(x_{j1}, x_{j2}, \dots, x_{jp}).$$

Resultado 2.1: Sean  $\underline{X}_1, \underline{X}_2, \dots, \underline{X}_n$  una muestra aleatoria  
de una distribución conjunta con vector media  $\underline{\mu}$  y matriz de  
var-cov.  $\Sigma$ . Entonces  $\bar{\underline{X}}$  es un vector estimador inseguido  
de  $\underline{\mu}$  y su matriz de var-cov es  $\frac{1}{n}\Sigma$ . Además,  $\frac{n}{n-1}S_n$  es  
un estimador inseguido de  $\Sigma$ , i.e.:

$$E\left[\left(\frac{n}{n-1}\right)S_n\right] = \Sigma$$

$$\text{Demo/ } E(\bar{X}) = E\left(\frac{1}{m} \underline{X}_1 + \frac{1}{m} \underline{X}_2 + \dots + \frac{1}{m} \underline{X}_m\right) \quad (5)$$

$$= \frac{1}{m} E\underline{X}_1 + \dots + \frac{1}{m} E\underline{X}_m = \sum_{j=1}^m \frac{1}{m} E(X_j)$$

$$\stackrel{iid}{=} \frac{1}{m} m \cdot \mu = \mu //$$

Ahora bien:

$$(\bar{X} - \mu)(\bar{X} - \mu)' = \left( \frac{1}{m} \sum_{j=1}^m (\underline{X}_j - \mu) \right) \left( \frac{1}{m} \sum_{e=1}^m (\underline{X}_e - \mu) \right)'$$

$$\text{(¿Porque?)} \Rightarrow = \frac{1}{m^2} \sum_{j=1}^m \sum_{e=1}^m (\underline{X}_j - \mu)(\underline{X}_e - \mu).$$

Responda!

$$\text{Entonces } E\left[ (\bar{X} - \mu)(\bar{X} - \mu)' \right] = \frac{1}{m^2} \sum_{j=1}^m \sum_{e=1}^m E[(\underline{X}_j - \mu)(\underline{X}_e - \mu)']. \quad (\text{ec. 4})$$

$$\text{Para todo } j \neq e, \quad E[(\underline{X}_j - \mu)(\underline{X}_e - \mu)'] = 0 \quad (\text{¿Porque?}),$$

$$\text{entonces } \text{Cov}(\bar{X}) = E\left[ (\bar{X} - \mu)(\bar{X} - \mu)' \right] \stackrel{(\text{ec. 4})}{=} \frac{1}{m^2} \sum_{j=1}^m E[(\underline{X}_j - \mu)(\underline{X}_j - \mu)'].$$

$$\underline{X}_j \stackrel{iid}{=} \frac{1}{m} \times m \cdot \underline{\Sigma} = \frac{1}{m} \sum.$$

(6)

Ahora calculemos  $E[S_m]$ . Para esto, notemos que:

$$\begin{aligned}
 \sum_{j=1}^m (\underline{x}_j - \bar{x})(\underline{x}_j - \bar{x})' &= \sum_{j=1}^m (\underline{x}_j - \bar{x}) \underline{x}_j' + \left( \sum_{j=1}^m (\underline{x}_j - \bar{x}) \right) (-\bar{x})' \\
 &= \sum_{j=1}^m \underline{x}_j \underline{x}_j' - \sum_{j=1}^m \bar{x} \underline{x}_j' \\
 &= \sum_{j=1}^m \underline{x}_j \underline{x}_j' - \bar{x} \sum_{j=1}^m \underline{x}_j' \\
 &= \sum_{j=1}^m \underline{x}_j \underline{x}_j' - \bar{x} m \bar{x}' = \sum_{j=1}^m \underline{x}_j \underline{x}_j' - m \bar{x} \bar{x}' 
 \end{aligned}$$

$$\therefore E\left[\sum_{j=1}^m (\underline{x}_j - \bar{x})(\underline{x}_j - \bar{x})'\right] = \sum_{j=1}^m E(\underline{x}_j \underline{x}_j') - m E(\bar{x} \bar{x}') \quad (\text{ec. } \otimes)$$

Ahora bien, para todo vector aleatorio  $\underline{V}$  con media  $\mu$  y varianza

$\mu_V$  y  $\Sigma_V$  respectivamente,  $E(VV') = \Sigma_V + \mu_V \mu_V'$ , por lo tanto:

$$\begin{array}{lcl}
 \text{L.D. de ec. } \otimes &=& \sum_{j=1}^m (\Sigma + \mu \mu') - m \left[ \frac{1}{m} \Sigma + \mu \mu' \right] \\
 \text{(lado derecho)} & &
 \end{array}$$

$$\begin{aligned}
 \underline{x}_j \text{ iid} &\Rightarrow m \Sigma + m \cancel{\mu \mu'} - \cancel{\Sigma} - m \cancel{\mu \mu'} \\
 &= (m-1) \Sigma. \quad \text{Como } S_m = \frac{1}{m} \left( \sum_{j=1}^m \underline{x}_j \underline{x}_j' - m \bar{x} \bar{x}' \right)
 \end{aligned}$$

resulta que

$$E(S_m) = \frac{(m-1)}{m} \Sigma$$

□

(7)

Matriz Var-Cov (insurgada):

$$S = \left( \frac{m}{m-1} \right) S_m = \frac{1}{m-1} \sum_{j=1}^m (\underline{x}_j - \bar{\underline{x}})(\underline{x}_j - \bar{\underline{x}})'$$

En lo que queda del curso,  $S$  reemplazará a  $S_m$  en la mayoría de las pruebas estadísticas que veremos.

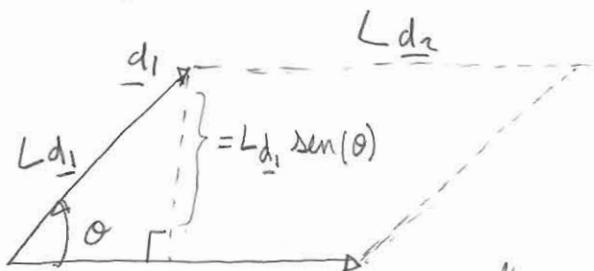
### 3) Variancia Generalizada:

Con una sola variable, la varianza de la muestra es suficiente para describir la variabilidad en los datos. Cuando se observan  $p$  variables, se puede representar la variabilidad total de la muestra con un sólo número: el determinante de  $S$ . A este se le llama "Variancia generalizada".

#### 3.1 Variancia generalizada cuando $p=2$ :

Sean  $S = \begin{bmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{bmatrix}$  y  $\begin{cases} \underline{d}_1 = \underline{x}_1 - \bar{\underline{x}}_1 \\ \underline{d}_2 = \underline{x}_2 - \bar{\underline{x}}_2 \end{cases}$ ,

Entonces,



$$\begin{cases} L_{\underline{d}_1} = \sqrt{\sum_{j=1}^m (x_{j1} - \bar{x}_{1j})^2} \\ \vdots = \sqrt{(m-1) S_{11}} \end{cases} \quad \begin{cases} L_{\underline{d}_2} = \sqrt{\sum_{j=1}^m (x_{j2} - \bar{x}_{2j})^2} \\ \vdots = \sqrt{(m-1) S_{22}} \end{cases}$$

$$\theta = \text{Área paralelogramo} = L_{\underline{d}_1} \times (L_{\underline{d}_2} \cdot \sin(\theta))$$

$$A = L_d_2 L_d_1 \sin(\phi) = L_d_2 L_d_1 \sqrt{1 - \cos^2(\phi)}$$

$$= \sqrt{(m-1)S_{22}} \sqrt{(m-1)S_{11}} \sqrt{1 - n_{12}^2} \quad (\text{porque } \cos(\phi) = n_{12}) \\ = (m-1) \sqrt{S_{22}} \sqrt{S_{11}} \sqrt{1 - n_{12}^2}$$

Teniendo en cuenta que  $|S| = S_{11} S_{22} - S_{12}^2$

$$= S_{11} S_{22} - (\sqrt{S_{11}} \sqrt{S_{22}} n_{12})^2 \\ = S_{11} S_{22} (1 - n_{12}^2)$$

entonces  $|S| = \alpha^2 / (m-1)^2$  y en  $p$  dimensiones:

$$\text{Varianza general: } \sigma^2 = |S| \cdot (m-1)^{-p} (\text{Volumen})^2$$

este es el significado  
físico de los parámetros  
de covariancia.

Considerando la interpretación de una forma cuadrática entre  $\underline{x}$ , y sustituyendo la matriz  $A$  en la ecuación para  $S^{-1}$  y  $\underline{x}$  como el punto fijo, entre las coordenadas  $\underline{x}' = [x_1, x_2, \dots, x_p]$  de los puntos q' están a una distancia este. 'c' de  $\underline{x}$  satisfacen

$$(\underline{x} - \underline{\delta})^T S^{-1} (\underline{x} - \underline{\delta}) = c^2$$

(9)

La hyper-ellipse definida por  $(\underline{x} - \bar{\underline{x}})' S^{-1} (\underline{x} - \bar{\underline{x}}) = c^2$  está relacionada a la varianza generalizada por la ecuación

$$\text{Volumen de } \left\{ \underline{x} : (\underline{x} - \bar{\underline{x}})' S^{-1} (\underline{x} - \bar{\underline{x}}) \leq c^2 \right\} = \frac{2 \pi^{p/2} \cdot c^p \cdot |S|^{1/2}}{p \Gamma(p/2)}$$