

Intro

Wednesday, January 26, 2022 6:56 PM

- * Reduce de datos o simplifica estructura
- * Ordenar y agrupar datos
- * Investigar dependencia entre vars.
- * Técnicas de predic. → regresiones: lineal, multivar...
- * Construc. P.H. (pruebas Hip.)

↳ Nota: matricial

Item 1	Var 1	Var p
	x_{11}	x_{1p}
Item n	x_{n1}	x_{np}

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

→ Medidas q' resumen los datos: media, mediana, var.

• Media muestral

$$\bar{x}_k = \frac{1}{n} \sum_{i=1}^n x_{ik} \quad k=1, \dots, p$$

• Varianza muestral $s_k^2 = \frac{1}{n-1} \sum (x_{ik} - \bar{x}_k)^2$

$$s_k^2 = \frac{1}{n} \sum (x_{ik} - \bar{x}_k)^2 \quad k=1, \dots, p$$

$$s_k^2 = s_{kk}$$

• Desv. estándar: $\sqrt{s_{kk}} = s_k$ • n obs. de 2 vars. $\begin{bmatrix} x_{11} & \dots & x_{1p} \\ x_{12} & \dots & x_{n2} \end{bmatrix}$ • Covarianza: $s_{12} = \frac{1}{n} \sum (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2)$

↳ Asocia. lin. entre las vars.

↳ Promedio del producto entre las desviaciones de sus respectivas medias.

↳ Si hay tanto grandes como pequeños vals en ambas vars. se presentan de forma conjunta, s_{12} es positiva.↳ Si hay vals grandes en una y pequeños en otra, s_{12} es neg.↳ Si no hay asocia. entre ambas, s_{12} es aprox 0.

• Coef. Correla. muestral: Normalizo la cov para no depender de las unidades.

$$r_{12} = \frac{s_{12}}{\sqrt{s_{11} s_{22}}} = \frac{\sum (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2)}{\sqrt{\sum (x_{i1} - \bar{x}_1)^2 \sum (x_{i2} - \bar{x}_2)^2}} \rightarrow [-1, 1]$$

↳ Asocia. lineal entre 2 vars., q' no depende de unidades.

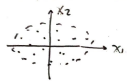
• $r=0$: No hay asocia.• $r=1$: correla. pos → una aumenta y la otra disminuye (en promedio). $p=(x_1, \dots, x_p)$ • Distancia Euclídeana $d(0, p) = \sqrt{x_1^2 + \dots + x_p^2}$ Todos los P a la misma dist. cuadrada c^2 de 0 están sobre una hipersf. n-1.

$$d(p, 0) = \sqrt{(x_1 - y_1)^2 + \dots + (x_p - y_p)^2}$$

↳ Podemos tener problemas por la dif. de unidades entre componentes, las cuales contribuyen igualm. al cálculo.

Cuando las coord. representan medidas q' pueden presentar variaciones aleatorias de dif. magnitudes es conveniente ponderar con mayor peso a las vars. con menos varia.

↳ Descomponer una distancia estadística q' considere las difs. en varia. y la presencia de correlac.

→ Hay n pares de medidas en x_1 y x_2 .

→ Tienden a 0 y varianzas indep.

→ Los P en x_1 tienen + variabilidad y los de x_2 → x_2 tiene más pondera.

$$\rightarrow x_1^* = x_1 / \sqrt{s_{11}}$$

$$\rightarrow x_2^* = x_2 / \sqrt{s_{22}}$$

$$d(0, P) = \sqrt{(x_1^*)^2 + (x_2^*)^2}$$

• Si la variabilidad en ambas coordenadas es igual, usar la distancia euclídeana.

• Todos los P (x_1, x_2) con dist cuadrada c^2 satisfacen q'

$$c^2 = \frac{x_1^2}{s_{11}} + \frac{x_2^2}{s_{22}} \text{ están en una elipse centrada en } 0.$$

• Distancia estadística entre 2 puntos. → Assume indep. en las vars.

$$d(p, q) = \sqrt{\frac{(x_1 - y_1)^2}{s_{11}} + \dots + \frac{(x_p - y_p)^2}{s_{pp}}}$$

No requiere la indep.

$$d(0, P) = \sqrt{\frac{x_1^2}{s_{11}} + \frac{x_2^2}{s_{22}}}$$

$$\begin{aligned} x_1 &= x_1 \cos \theta + x_2 \sin \theta \\ x_2 &= -x_1 \sin \theta + x_2 \cos \theta \end{aligned}$$

$$P = (x_1, x_2); Q = (y_1, y_2)$$

$$d(p, q) = \sqrt{a_{11}(x_1 - y_1)^2 + 2a_{12}(x_1 - y_1)(x_2 - y_2) + a_{22}(x_2 - y_2)^2}$$

$$d = \sqrt{a_{11}x_1^2 + 2a_{12}x_1x_2 + a_{22}x_2^2}$$

a_{ij} son func. de θ .

Ahora, con p vars.

$$d(p, q) = \sqrt{a_{11}(x_1 - y_1)^2 + \dots + a_{pp}(x_p - y_p)^2 + 2a_{12}(x_1 - y_1)(x_2 - y_2) + \dots + 2a_{p-1,p}(x_{p-1} - y_{p-1})(x_p - y_p)}$$