

ANÁLISIS DE COMPONENTES PRINCIPALES (PCA)

Como estudiar datos x_{g1}, \dots, x_{gn} , cuya dimensión puede ser grande, no es una tarea inmediata, nos gustaría una forma de representar x_{g1}, \dots, x_{gn} en dimensiones menores. En estadística hemos aprendido que un "buen" resumen de los datos es un promedio, pero una forma de hacerlo más "flexible" sería darle un peso a cada componente de x_{gi} , la cual midiera una "importancia" en el promedio.

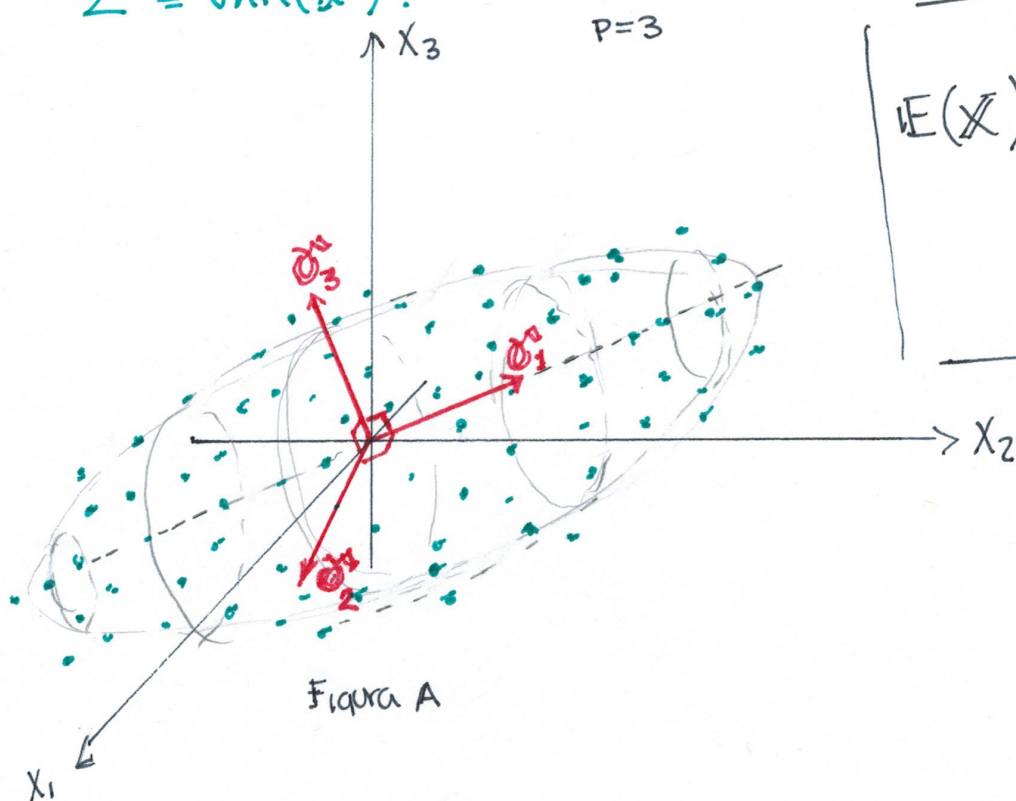
x_{g1}, \dots, x_{gn} realizaciones del vect.
aleatorio X

$$\alpha'X = \sum_{i=1}^P \alpha_i X_i, \quad \text{donde } \sum_{i=1}^P \alpha_i^2 = 1.$$

La forma en que X varía es una característica de interés, por lo cual nos planteamos el problema de encontrar $\alpha \in \mathbb{R}^P$, tal que $\sum_{j=1}^P \alpha_j^2 = 1$ y $\alpha'X$ tiene varianza máxima. En otras palabras hay que encontrar $\alpha \in \mathbb{R}^P$ tal que $\alpha'X$ alcance el máximo

$$\begin{aligned} & \max \{ \text{VAR}(\alpha'X) : \alpha \in \mathbb{R}^P \text{ y } \|\alpha\|=1 \} \\ & = \max \{ \alpha' \text{VAR}(X) \alpha : \alpha \in \mathbb{R}^P \text{ y } \|\alpha\|=1 \} \end{aligned}$$

Si regresamos al Teorema C, específicamente al caso en que $B = I_p$, es claro que una solución al problema de optimización anterior, es tomar $\alpha = \varphi_1$, donde φ_1 es el vector propio correspondiente al más grande valor propio λ_1 de la matriz de covarianzas $\Sigma = \text{VAR}(X)$.



$$E(X) = \begin{pmatrix} E(X_1) \\ E(X_2) \\ E(X_3) \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

Si los datos tuvieran un comportamiento "ideal" ⁽¹⁾ como en el dibujo, la dirección en \mathbb{R}^3 en la que existe mayor variabilidad está

(1) ver página 15

(1)

Asumiendo x_{g1}, \dots, x_{gn} (los datos) son una muestra aleatoria del vector aleatorio \mathbb{X} tal que $\mathbb{X} \sim F$, donde la función de distribución multivariada F es tal que:

$$\mathbb{E}(\mathbb{X}) = \mathbf{0} \quad \text{y} \quad \text{VAR}(\mathbb{X}) = \Sigma > 0.$$

Además $F: \mathbb{R}^3 \rightarrow [0,1]$ con soporte como en el dibujo de la Figura A.

$$\text{Soporte}(F) = \{x \in \mathbb{R}^3: f(x) > 0\}$$

f = densidad de F .

dada por \mathbf{v}_1

Pregunta: ¿Porqué es así? (*)

Por el Teorema C la elección $\alpha = \mathbf{v}_1$ garantiza que $\text{VAR}(\mathbf{v}_1' \mathbf{X}) = \mathbf{v}_1' \text{VAR}(\mathbf{X}) \mathbf{v}_1$ alcanza el máximo, el cual vale λ_1 . ($\text{VAR}(\mathbf{v}_1' \mathbf{X}) = \lambda_1$).

(PC1) ... $Y_1 \equiv \mathbf{v}_1' \mathbf{X} \leftarrow$ La primera componente principal
Es una combinación lineal
"estandarizada" (es decir convexa:
 $\sum_{j=1}^P v_{1j}^2 = 1$; $\mathbf{v}_1 = (v_{11}, \dots, v_{1P})$)
de las componentes de \mathbf{X} .
Está construida para tener
varianza máxima.

(PC2) ... $Y_2 \equiv \mathbf{v}_2' \mathbf{X} \leftarrow$ La segunda componente principal
donde \mathbf{v}_2 es el vector propio de Σ asociado a
 $\lambda_2 (\leq \lambda_1)$, entonces $\Sigma \mathbf{v}_2 = \lambda_2 \mathbf{v}_2$ por lo
que $\mathbf{v}_2' \Sigma \mathbf{v}_2 = \lambda_2$ ($\text{VAR}(\mathbf{v}_2' \mathbf{X}) = \lambda_2 \leq \lambda_1$).

(PC3) ... $Y_3 \equiv \mathbf{v}_3' \mathbf{X} \leftarrow$ La tercera componente principal
etc...