



ANÁLISIS DE COMPONENTES PRINCIPALES (ACP)

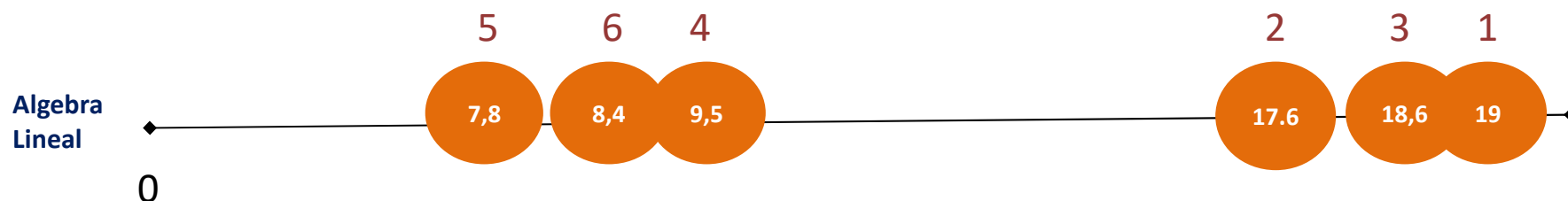


Definición

Es una técnica matemática que permite reducir la dimensionalidad centrándose en las varianzas de un conjunto de datos multivariados obtenidos de una población cuya distribución de probabilidades no necesita ser conocida.

CALIFICACION FINAL DE ESTUDIANTES SEGÚN CURSOS

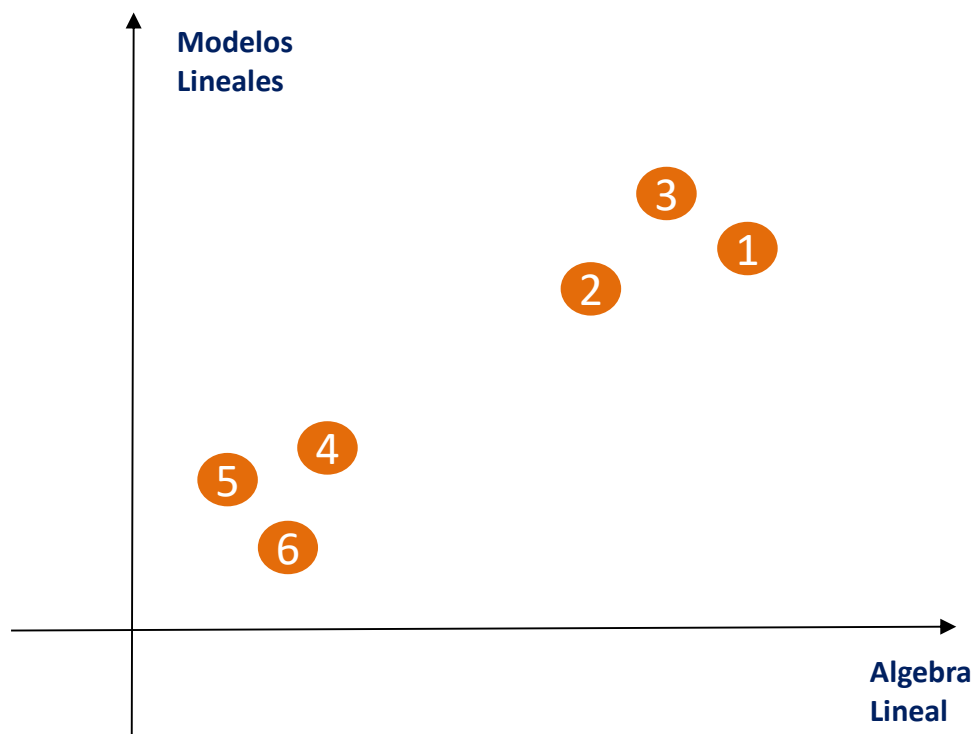
GRAFICO EN
UNA DIMENSIÓN



	Estudiante 1	Estudiante 2	Estudiante 3	Estudiante 4	Estudiante 5	Estudiante 6
Algebra Lineal	19,0	17,6	18,6	9,5	7,8	8,4

CALIFICACION FINAL DE ESTUDIANTES SEGÚN CURSOS

GRAFICO EN
DOS DIMENSIONES

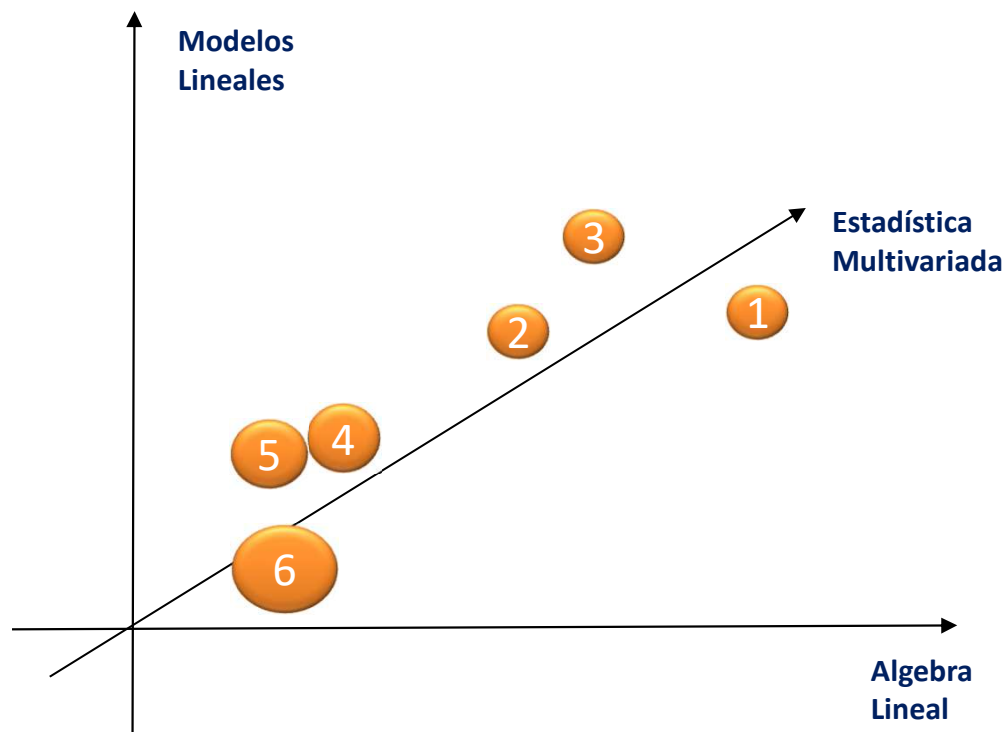


	Estudiante 1	Estudiante 2	Estudiante 3	Estudiante 4	Estudiante 5	Estudiante 6
Algebra Lineal	19,0	17,6	18,6	9,5	7,8	8,4
Modelos Lineales	19,2	15,8	19,6	8,7	7,2	5,0



CALIFICACION FINAL DE ESTUDIANTES SEGÚN CURSOS

BOSQUEJO EN
TRES DIMENSIONES



	Estudiante 1	Estudiante 2	Estudiante 3	Estudiante 4	Estudiante 5	Estudiante 6
Algebra Lineal	19,0	17,6	18,6	9,5	7,8	8,4
Modelos Lineales	19,2	15,8	19,6	8,7	7,2	5,0
Estadística Multivariada	17,2	11,0	15,1	5,0	4,4	3,0

¿EN CUATRO DIMENSIONES?

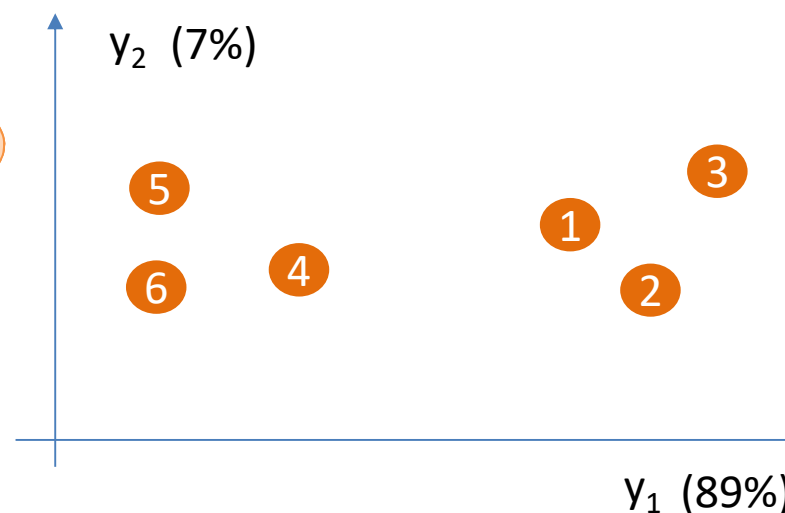


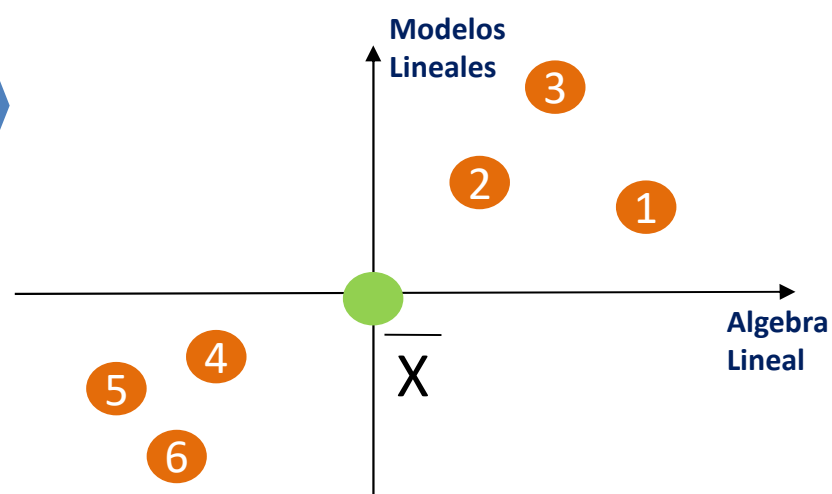
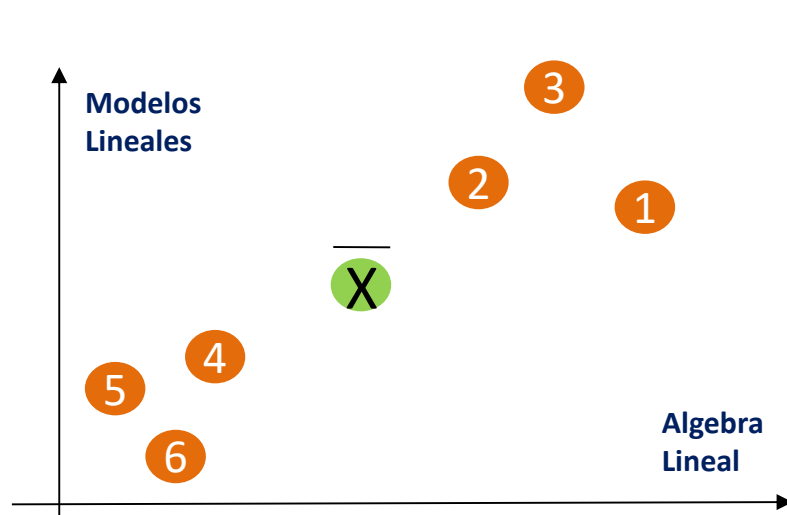
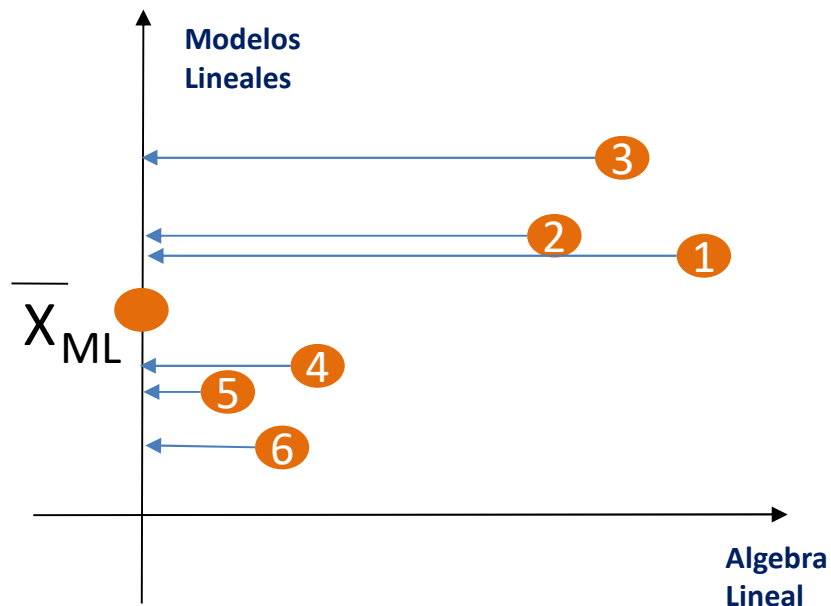
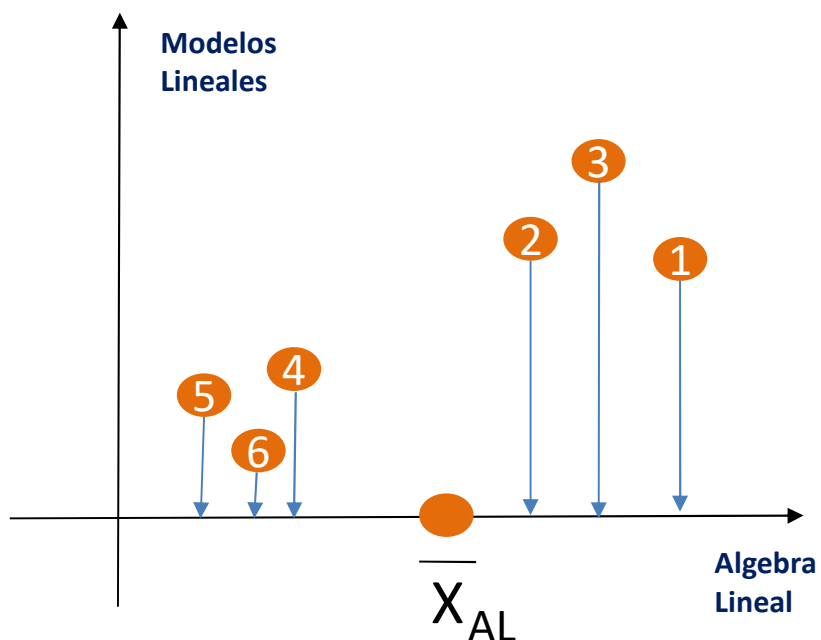
	Estudiante 1	Estudiante 2	Estudiante 3	Estudiante 4	Estudiante 5	Estudiante 6
Algebra Lineal	19,0	17,6	18,6	9,5	7,8	8,4
Modelos Lineales	19,2	15,8	19,6	8,7	7,2	5,0
Estadística Multivariada	17,2	11,0	15,1	5,0	4.4	3,0
Técnicas de redacción	12,1	15,0	13,3	6,0	8.3	15,0

¿Como se puede representar en
más de tres dimensiones
gráficamente?

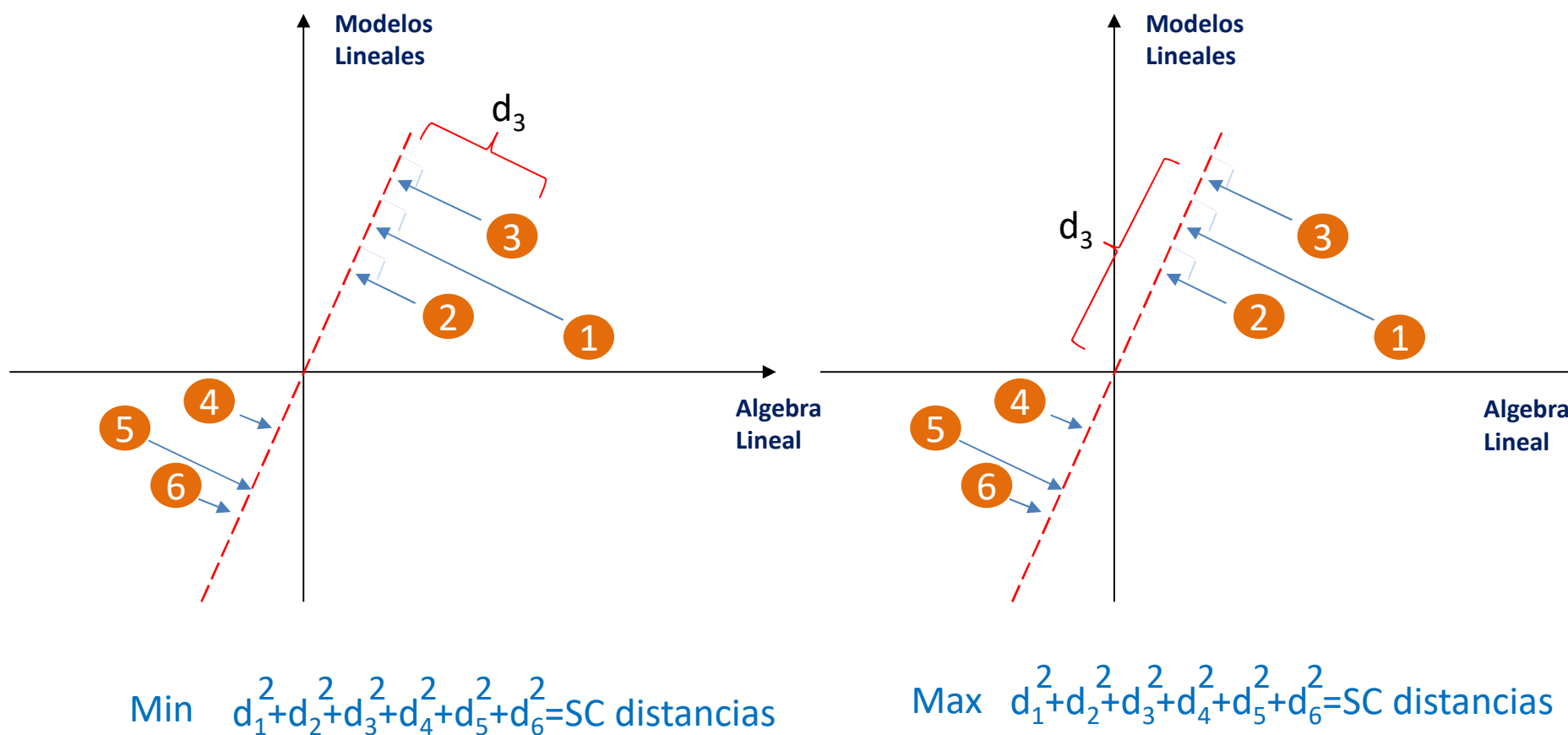
Calificación Final de Estudiantes según Cursos

	Estudiante 1	Estudiante 2	Estudiante 3	Estudiante 4	Estudiante 5	Estudiante 6
Algebra Lineal	19,0	17,6	18,6	9,5	7,8	8,4
Modelos Lineales	19,2	15,8	19,6	8,7	7,2	5,0
Estadística Multivariada	17,2	11,0	15,1	5,0	4,4	3,0
Técnicas de redacción	12,1	15,0	13,3	6,0	8,3	15,0

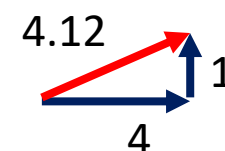




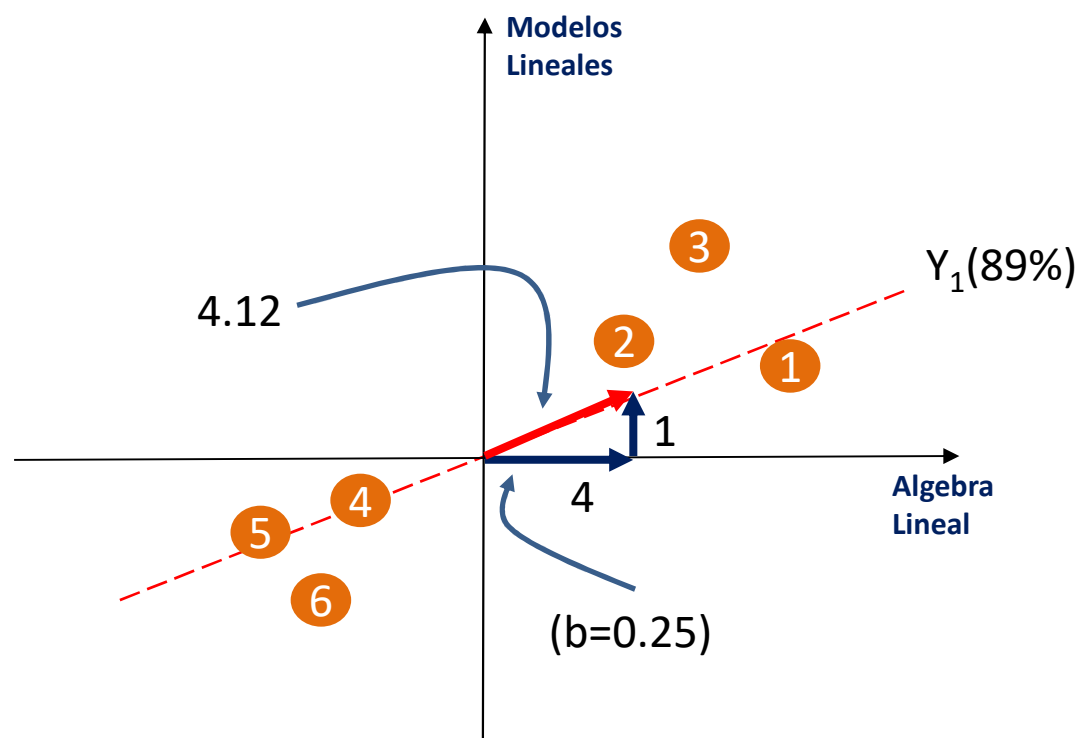
¿Como trabaja el análisis de componentes principales?



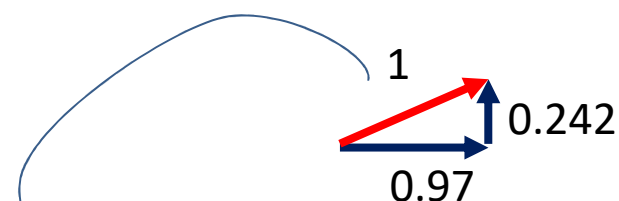
Combinación lineal = 4 de AL por 1 de ML



Relación pitagórica



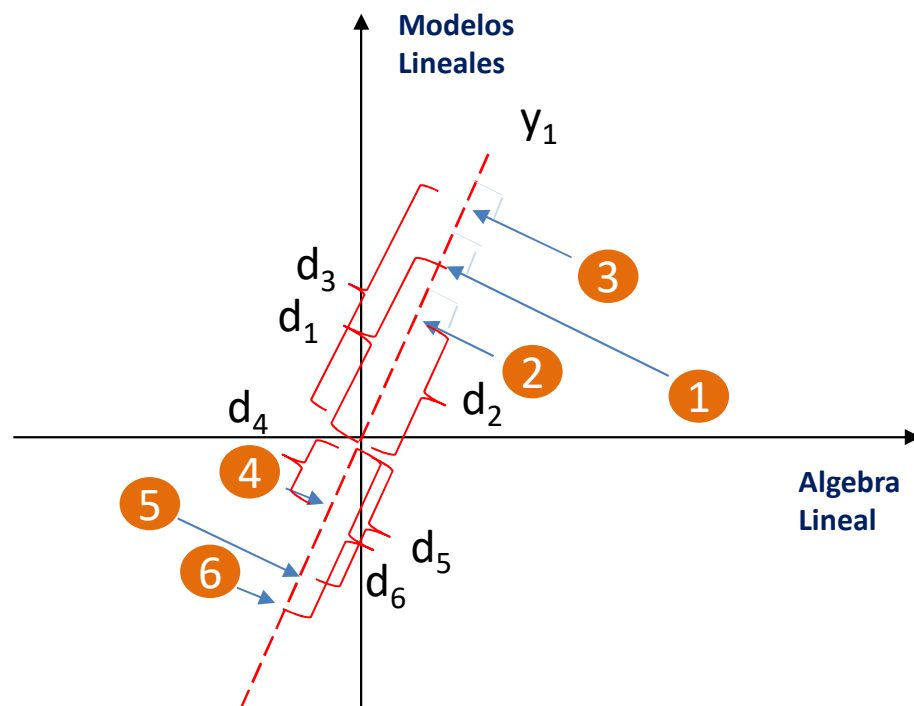
$$\frac{4.12}{4.12} = \sqrt{\left(\frac{1}{4.12}\right)^2 + \left(\frac{4}{4.12}\right)^2}$$



0.242
0.970

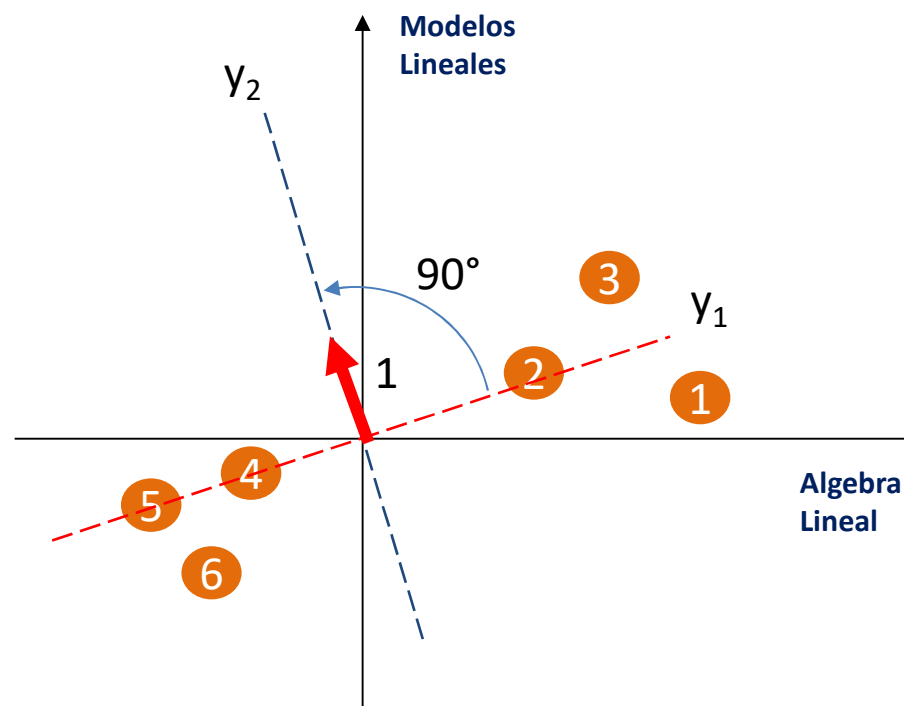
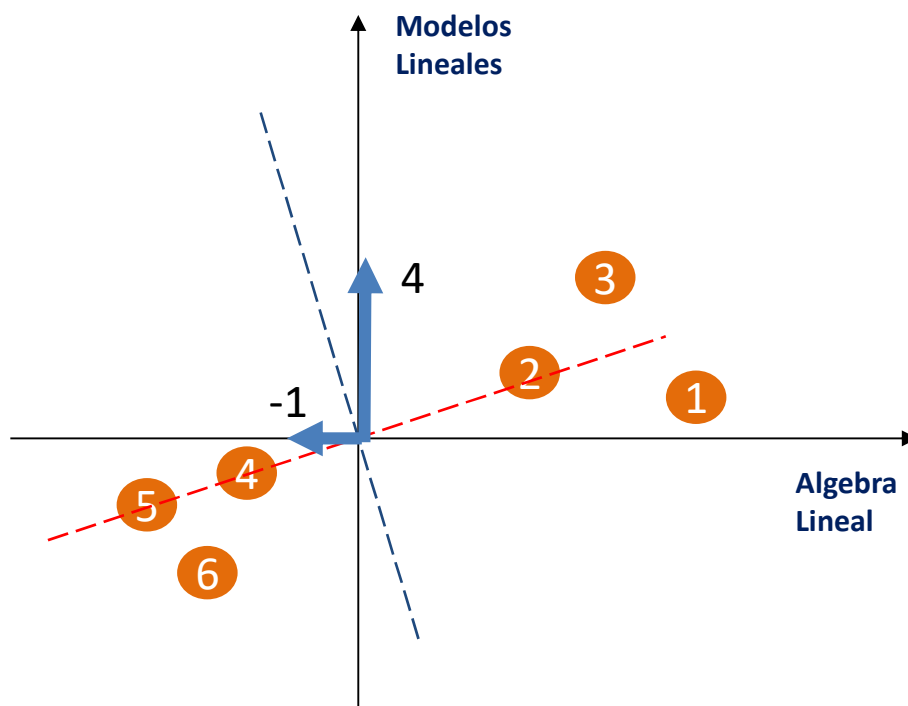
Vector
característico,
vector singular o
“eigenvector”

Combinación lineal = **0.97** puntos de AL por **0.242** de ML

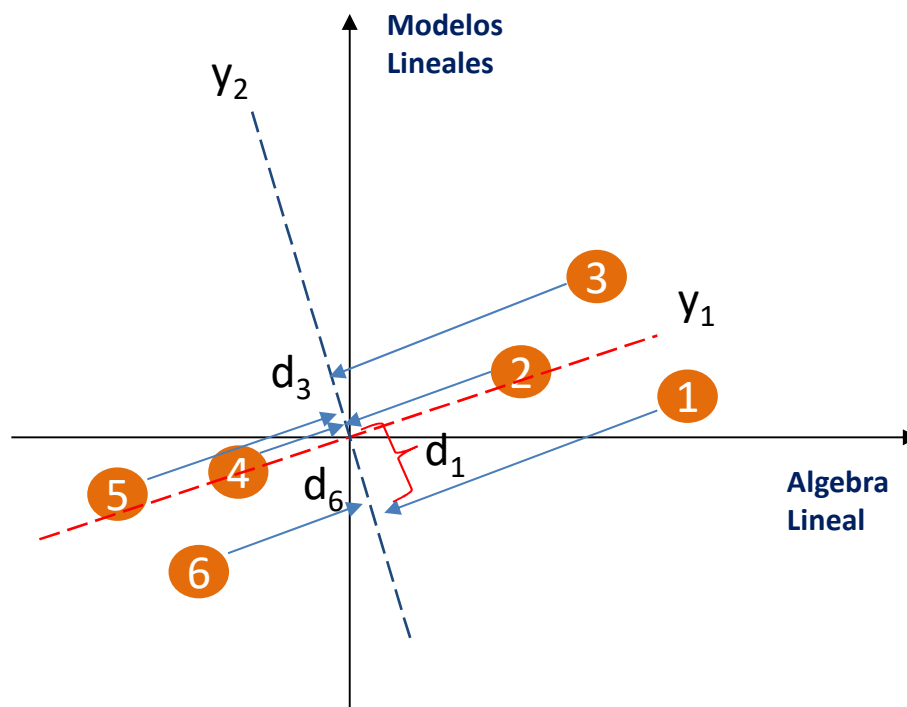


$$d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2 + d_6^2 = \text{SC distancias} = (\text{valor característico de } y_1)^2$$

Combinación lineal =
0.97 puntos de ML por **-0.242** de AL



$$\begin{bmatrix} -0.242 \\ 0.970 \end{bmatrix}$$
 Vector característico,
 vector singular o
 “eigenvector” de y_2



$$d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2 + d_6^2 = \text{SC distancias} = (\text{valor característico de } y_2)^2$$

Componentes Principales utilizando la descomposición de los valores característicos

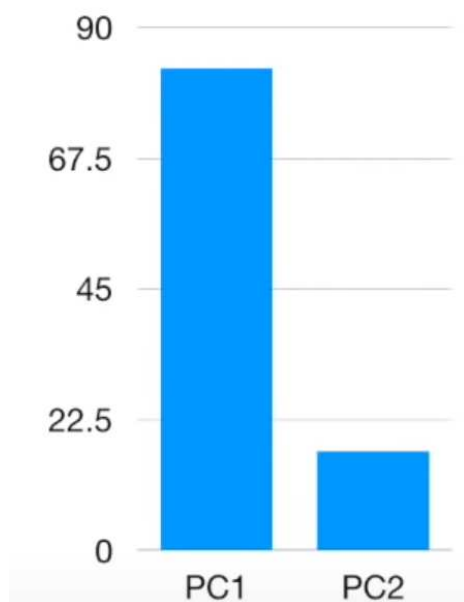
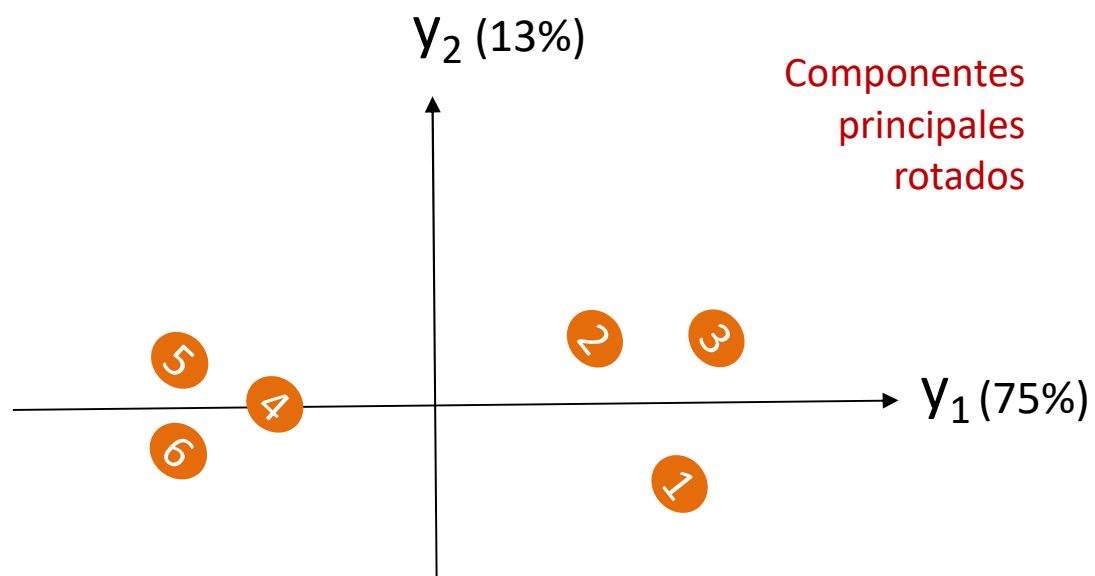


Gráfico de Sedimentación o
Descomposición Espectral





UNIVERSIDAD NACIONAL
DE INGENIERÍA

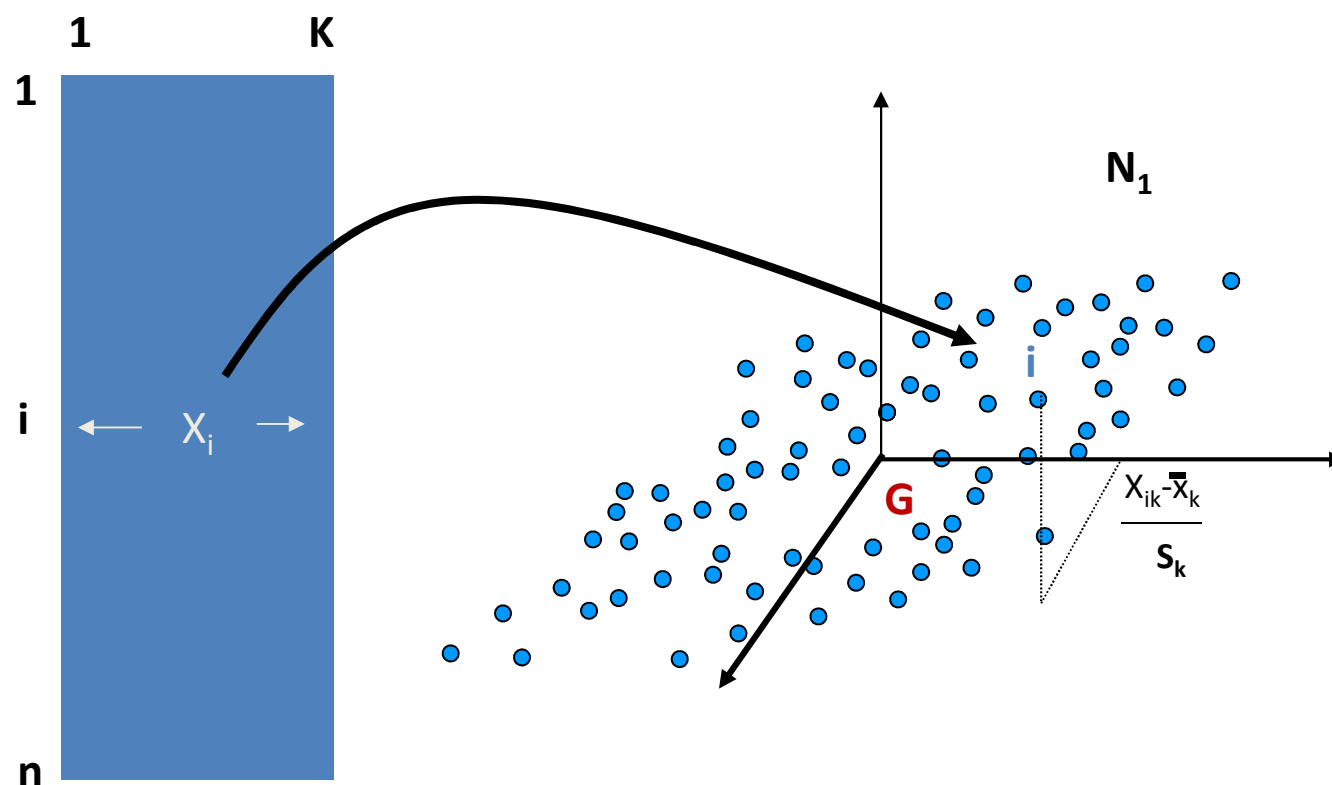
Escuela Profesional de Ingeniería Estadística – FIEECS

ESTADÍSTICA MULTIVARIADA – Análisis de Componentes Principales

Prof. Luis Huamanchumo de la Cuba

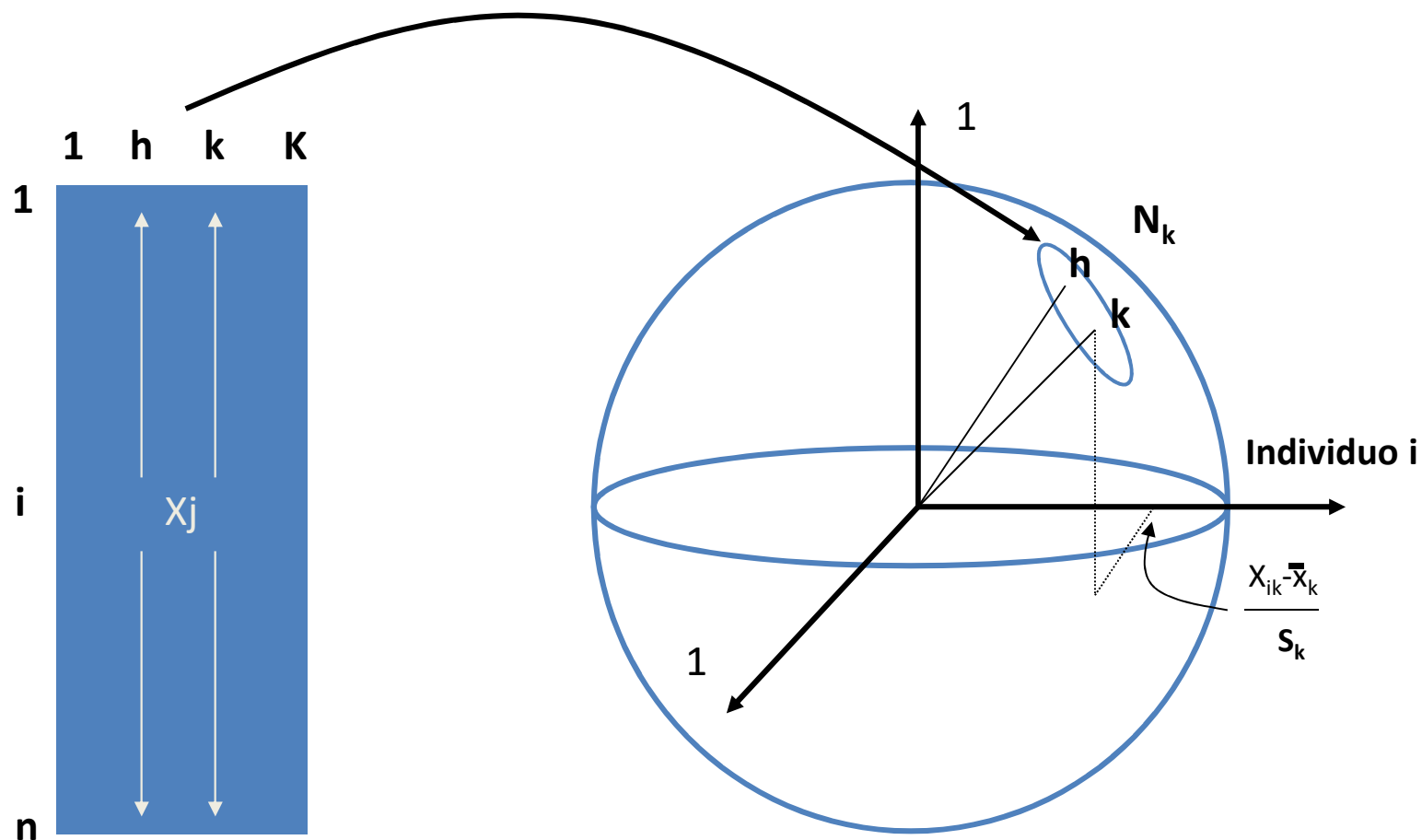
Generalización del Análisis Geométrico

La Nube de Individuos



Individuos como yuxtaposición de “n”
filas en R^k

La Nube de Variables



Norma de la Variable $X_k=1$
 (Hiperesfera)

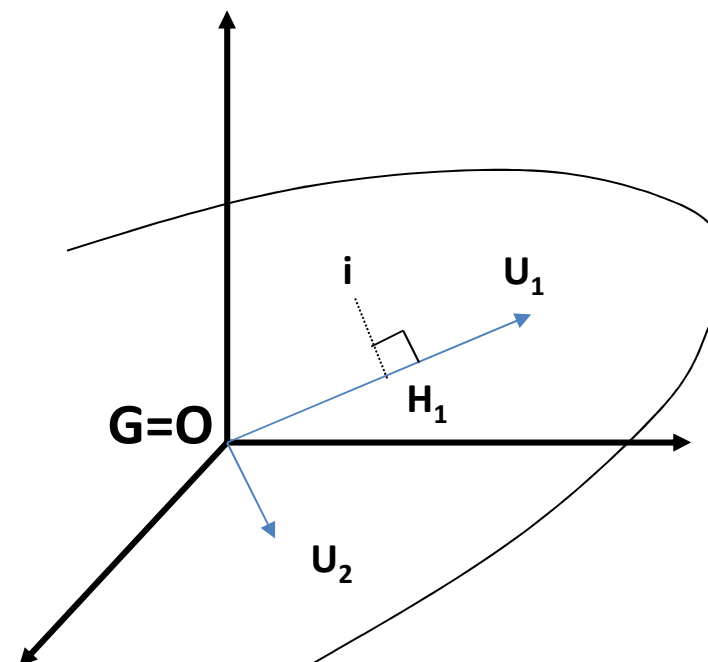
El Ajuste de la Nube de Individuos

Criterio

Máxima Inercia
respecto al Centro de
Gravedad G



Ejes Factoriales de
Máximo
Alargamiento de la
Nube



El individuo “i” se proyecta sobre u_1 en H_i

Se busca en primer lugar u_i que maximiza $\sum OH_i^2$

Se busca u_2 , ortogonal a u_1 , que satisface el mismo criterio

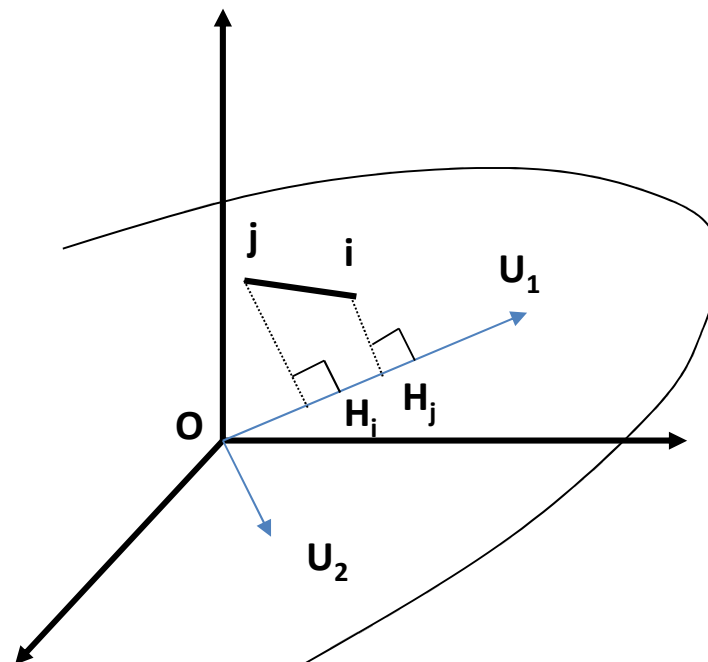
Si los individuos tienen pesos diferentes el criterio
Consiste en maximizar $\sum_i p_i OH_i^2$

Distancias entre Individuos

El eje u_1 maximiza $\sum_i \sum_j (OH_i - OH_j)^2$ o sea

$\sum_i \sum_j d^2(H_i, H_j)$ se acerca lo más posible a

$\sum_i \sum_j d^2(i, j)$



El Ajuste de la Nube de Variables

Criterio

Máxima Inercia
Proyectada

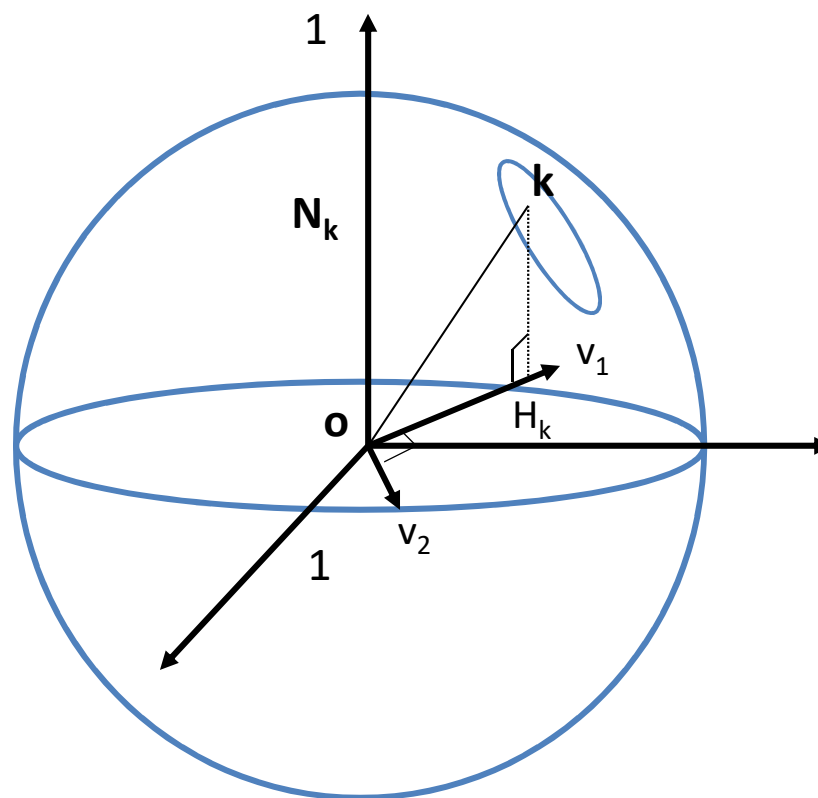


Ejes Factoriales
Maximizan la Suma de
Cosenos al Cuadrado

H_k : proyección sobre v_1 del punto que representa la variable k .

Se busca v_1 que maximice $\sum_k OH_k^2$

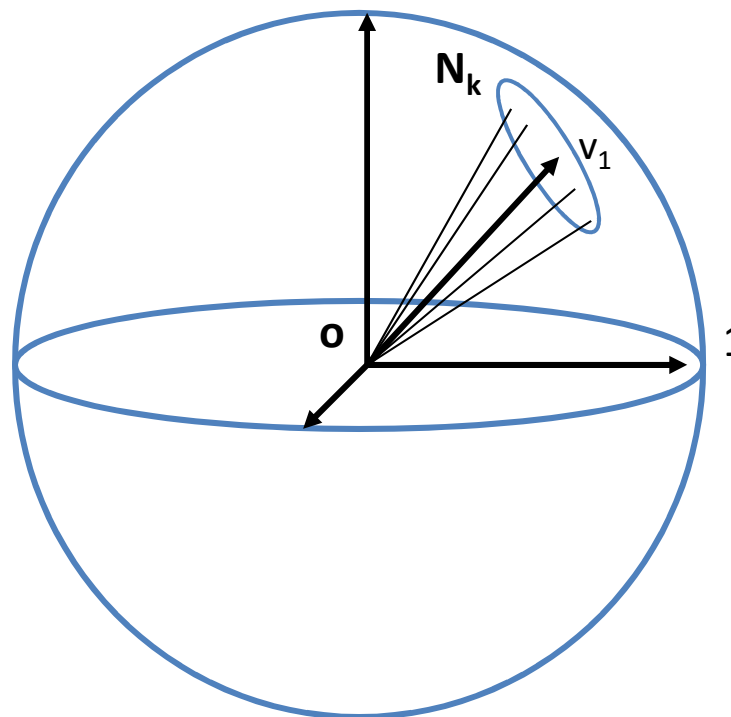
Después se busca v_2 ortogonal a v_1 que satisface el mismo criterio



El Efecto Talla en R^n

La nube N_k se concentra en un pequeño sector de la esfera.

La proyección de las variables sobre el primer eje factorial V_1 informa de la posición de N_k en relación a O .



Los Datos

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

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

Diagram illustrating the structure of a multivariate data matrix $X_{n \times p}$:

- The matrix is organized by rows and columns.
- Columns represent variables: x_1, x_2, \dots, x_p . An arrow labeled "Variables" points to these column labels.
- Rows represent observations: the first row is labeled "1ra observación" and the last row is labeled "n-ésima observación". Arrows point from these labels to the corresponding rows in the matrix.
- Individual elements are labeled as x_{ij} , where i is the observation index and j is the variable index.



Los Individuos

Lo que se evalúa es su **semejanza**. Dos individuos se asemejan más cuanto más próximos sean sus valores en el conjunto de las variables.

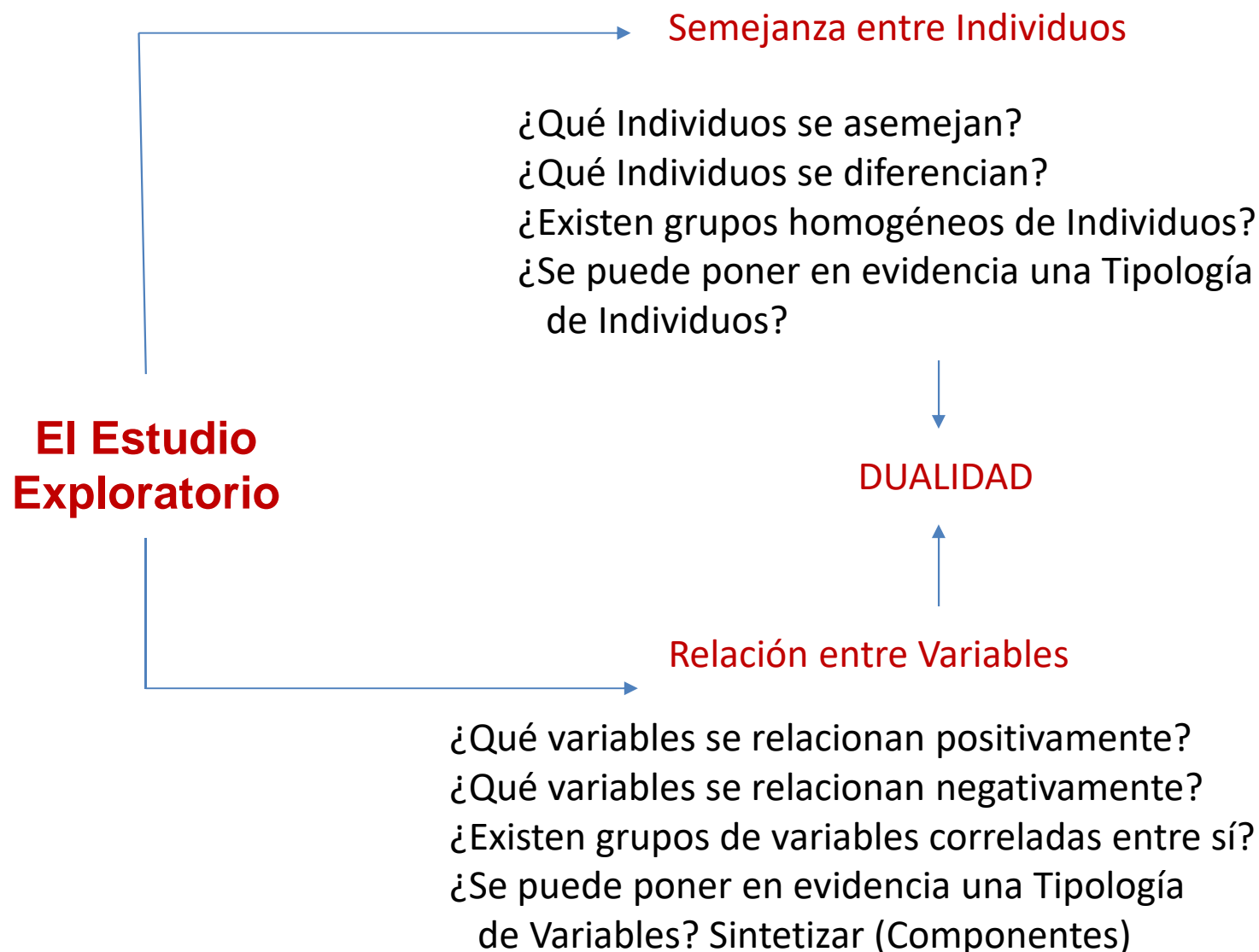
$$d^2(i,j) = \sum (x_{ik} - x_{jk})^2$$



Las Variables

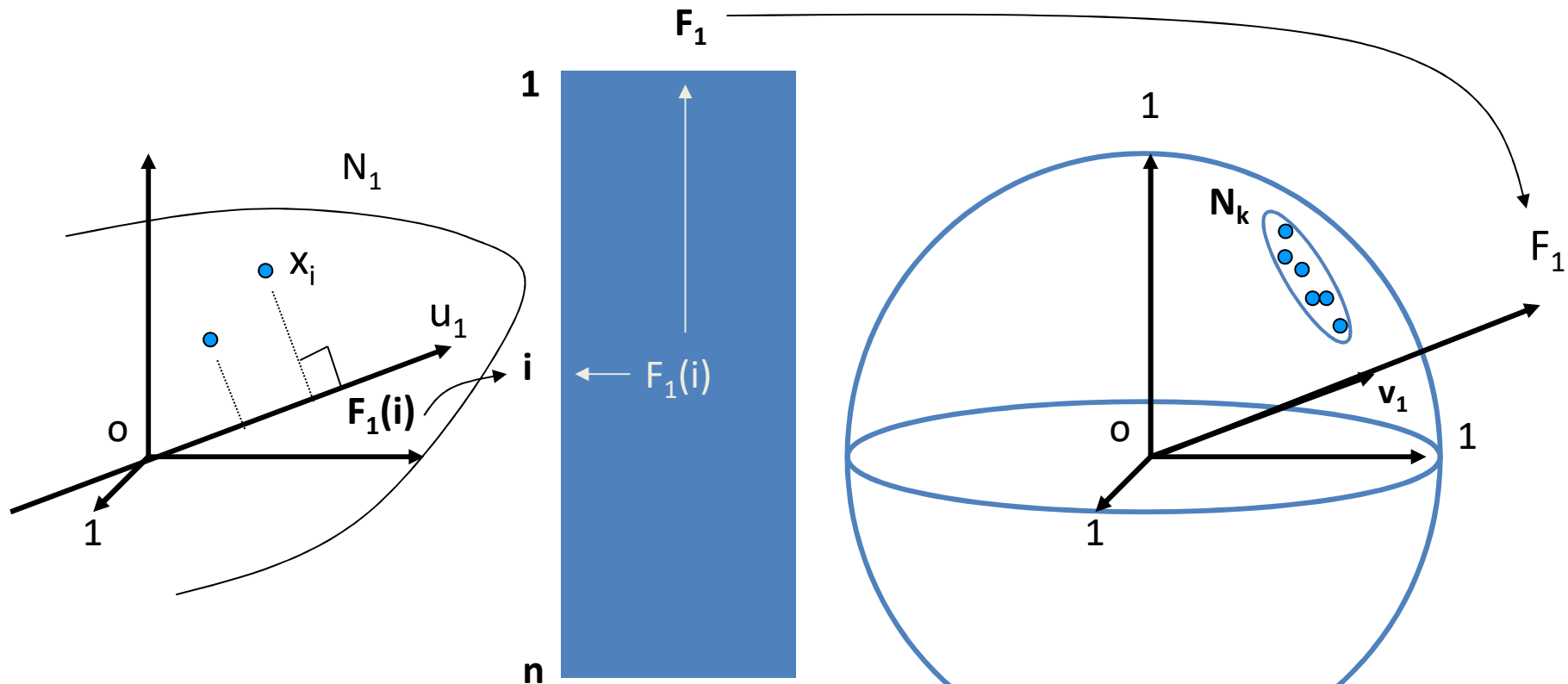
Lo que se evalúa es su **relación**. La relación entre dos variables se mide por el coeficiente de correlación lineal, en otros casos, se utiliza la covarianza.

$$r(k,h) = \frac{\text{Cov}(k,h)}{\sqrt{\text{var}(k) \text{var}(h)}}$$



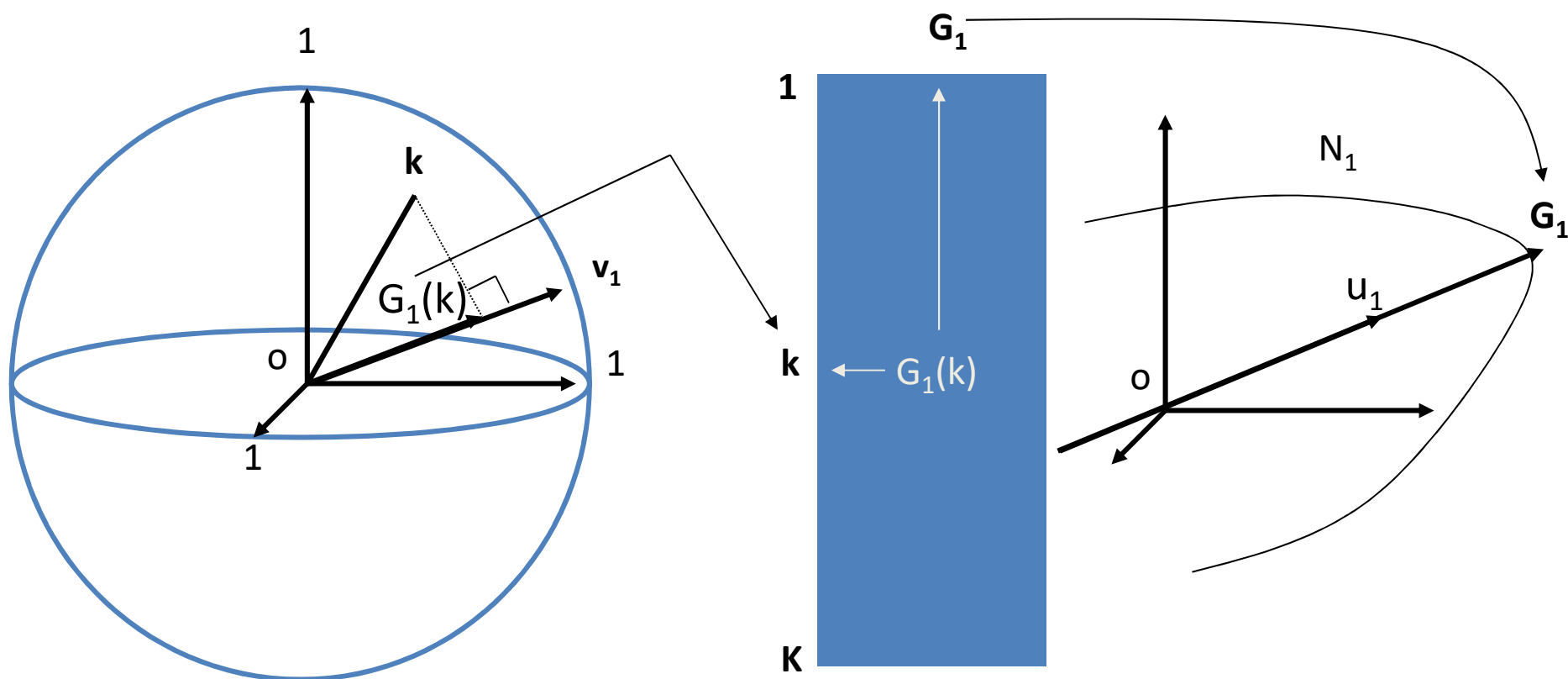


Primera Forma de Dualidad



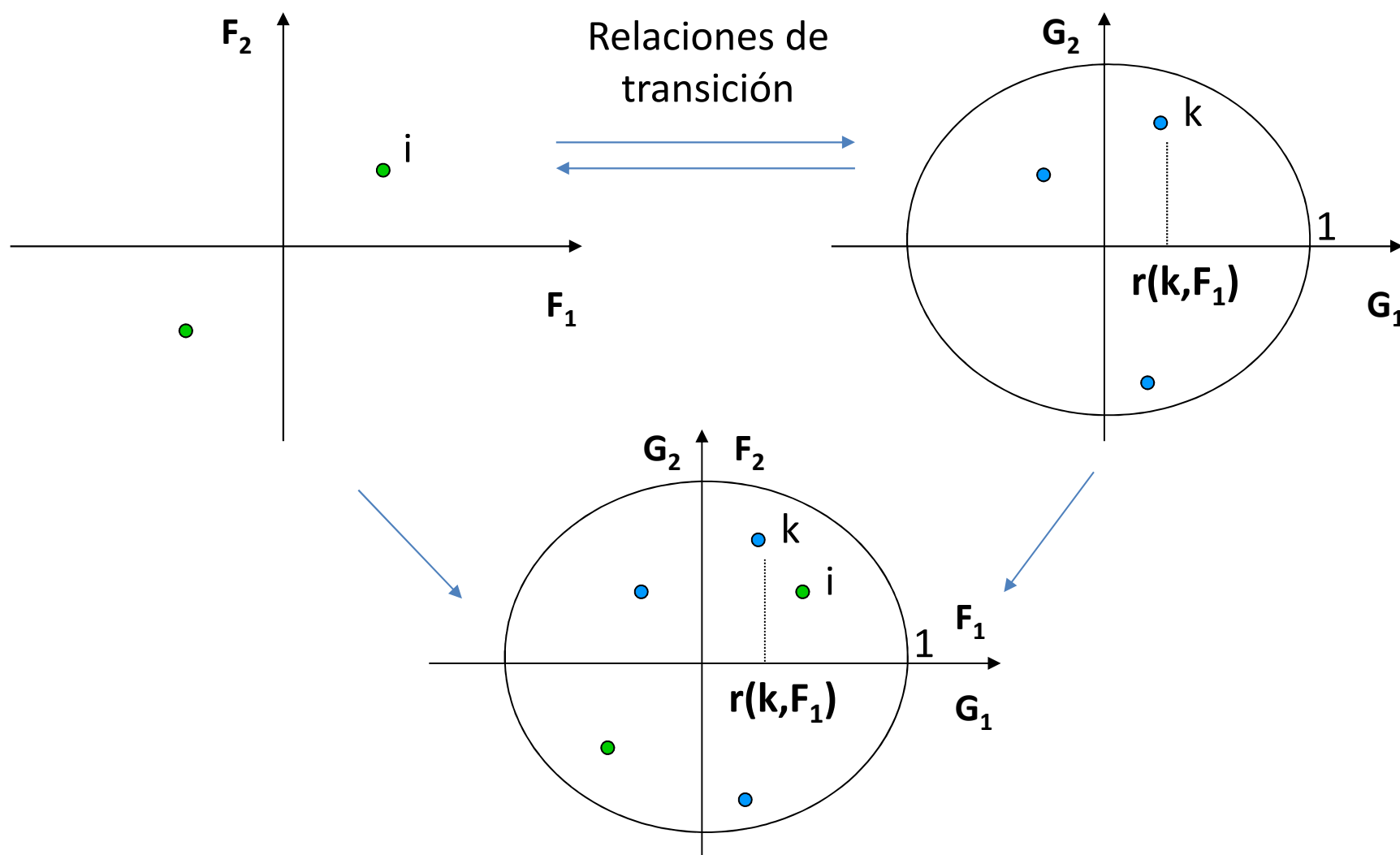
$$F_s(i) = \frac{1}{\sqrt{\lambda_s}} \sum_s \frac{x_{ik} - \bar{x}_k}{s_k} G_s(k)$$

Segunda Forma de Dualidad



$$G_s(i) = \frac{1}{\sqrt{\lambda_s}} \sum_i \frac{X_{ik} - \bar{x}_k}{S_k} F_s(k)$$

Análisis Dual





UNIVERSIDAD NACIONAL
DE INGENIERÍA

Escuela Profesional de Ingeniería Estadística – FIEECS

ESTADÍSTICA MULTIVARIADA – Análisis de Componentes Principales

Prof. Luis Huamanchumo de la Cuba

Cálculo de las Componentes Principales

Definición Operativa

La Componente Principal es una combinación lineal de ‘p’ variables aleatorias X_1, X_2, \dots, X_p , representa un nuevo eje de coordenadas (rotación) cuya dirección es de máxima variabilidad.

Dado el vector aleatorio $X = [X_1, X_2, \dots, X_p]$ y su covarianza Σ con valores característicos $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$

$$y_1 = L_1^t X = l_{11}X_1 + l_{21}X_2 + \dots + l_{p1}X_p$$

$$y_2 = L_2^t X = l_{12}X_1 + l_{22}X_2 + \dots + l_{p2}X_p$$

$$\begin{array}{c} \cdot \\ \cdot \\ \cdot \end{array} \quad \begin{array}{c} \cdot \\ \cdot \\ \cdot \end{array}$$

$$y_p = L_p^t X = l_{1p}X_1 + l_{2p}X_2 + \dots + l_{pp}X_p$$

$$\text{Var}(Y_i) = L_i^t \Sigma L_i \quad i = 1, 2, \dots, p$$

$$\text{Cov}(Y_i, Y_k) = L_i^t \Sigma L_k \quad i, k = 1, 2, \dots, p$$

Proceso Iterativo de Maximización

Formulación Matemática

K componentes principales

Primera iteración

$$\begin{aligned} \text{Max} \quad & L_1^t S L_1 \\ \text{s.a.} \quad & L_1^t L_1 = 1 \end{aligned}$$

Segunda iteración

$$\begin{aligned} \text{Max} \quad & L_2^t S L_2 \\ \text{s.a.} \quad & L_2^t L_2 = 1 \\ & L_2^t L_1 = 0 \end{aligned}$$

K-ésima iteración

$$\begin{aligned} \text{Max} \quad & L_K^t S L_K \\ \text{s.a.} \quad & L_k^t L_k = 1 \\ & L_k^t L_i = 0 \quad \forall i \end{aligned}$$



UNIVERSIDAD NACIONAL
DE INGENIERÍA

Escuela Profesional de Ingeniería Estadística – FIEECs

ESTADÍSTICA MULTIVARIADA – Análisis de Componentes Principales

Prof. Luis Huamanchumo de la Cuba

CRITERIO DE SELECCIÓN DE COMPONENTES PRINCIPALES



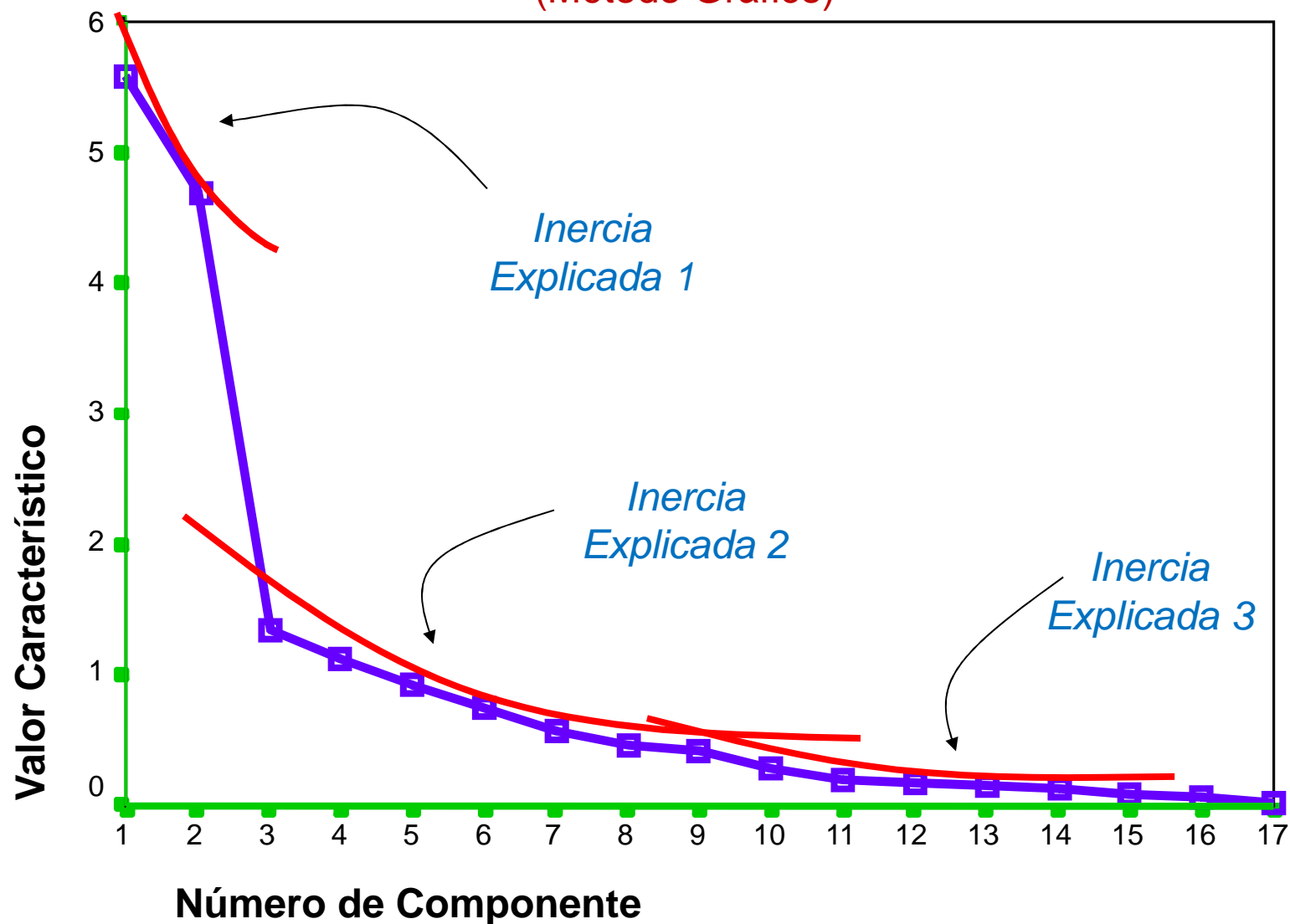
MÉTODO DE INERCIA TOTAL (Coeficiente de Inercia)

$$r_k = \frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{\text{Tr}(\Lambda)} \quad k \leq p$$

Se escogen las “p” componentes principales que explican el 100r% de la variabilidad total



MÉTODO DE ARCOS (Método Gráfico)



MÉTODO DE INERCIA PROMEDIO

$$\begin{array}{l} \text{Componente} \\ \text{Inercial "i-ésimo"} \geq 1 \\ i=1,\dots,p \end{array}$$

Se selecciona la Componente Principal cuyo valor característico es mayor o igual a 1

$$\text{Tr}(\Lambda) = p = \sum_i \lambda_i$$

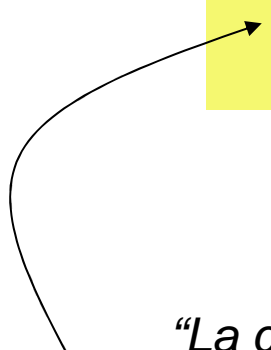
“p” número de componentes principales totales

Formalización de las Formas de Dualidad

Relación entre los espacios \mathcal{R}^p y \mathcal{R}^n

Definición

Los valores propios $\mu_1, \mu_2, \dots, \mu_p$ asociados a los vectores propios e_1, e_2, \dots, e_p de XX^t son iguales respectivamente a los valores propios $\lambda_1, \lambda_2, \dots, \lambda_p$ de X^tX


$$\mu_1 = \lambda_1, \mu_2 = \lambda_2, \dots, \mu_p = \lambda_p$$

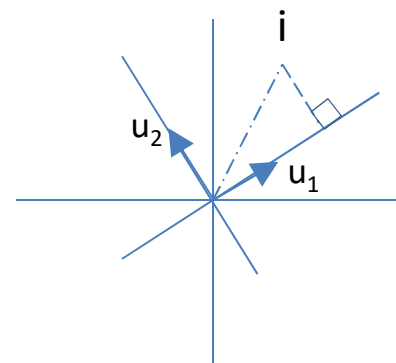
“La cantidad de información recogida por los ejes respectivos en ambos espacios es la misma”

ENFOQUE GEOMÉTRICO

Análisis en R^p :

X_i (observación i)
 $p \times 1$

$$\text{Comp}_{\bar{u}_1} x_i = F_1(i) = x_i' \bar{u}_1$$



$$F_{\alpha}() = X \bar{u}_{\alpha}$$

$(n \times p)(p \times 1)$

$$\begin{aligned} \max \quad & \bar{u}_{\alpha}' X' X \bar{u}_{\alpha} \\ \text{s.a.} \quad & u_{\alpha}' u_{\alpha} = 1 \end{aligned}$$

Sabiendo que:

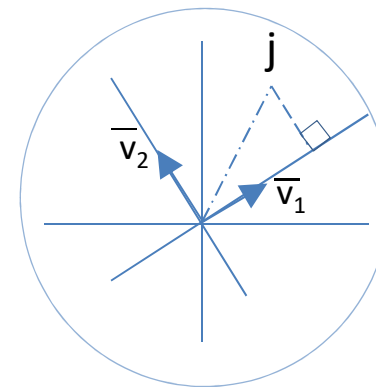
$$X'X u_{\alpha} = \lambda_{\alpha} u_{\alpha} \quad \text{--- (a)}$$

$$u_{\alpha}' X' X u_{\alpha} = \lambda_{\alpha} \quad (\text{máximo})$$

Análisis en R^n :

X_j (variable j)
 $n \times 1$

$$\text{Comp}_{\bar{v}_1} x_j = G_1(j) = x_j' \bar{v}_1$$



$$G_\alpha() = X' \bar{v}_\alpha$$

(n x p)(p x 1)

$$\begin{aligned} \max \quad & \bar{v}_\alpha' X X' \bar{v}_\alpha \\ \text{s.a.} \quad & \bar{v}_\alpha' \bar{v}_\alpha = 1 \end{aligned}$$

Sabiendo que:

$$X X' \bar{v}_\alpha = \mu_\alpha \bar{v}_\alpha \quad \text{--- (b)}$$

$$\bar{v}_\alpha' X X' \bar{v}_\alpha = \mu_\alpha \quad (\text{máximo})$$

Análisis entre R^p y R^n :

De (a): $X'X \bar{u}_\alpha = \lambda_\alpha \bar{u}_\alpha$ premultiplicando por X

$$(XX')X \bar{u}_\alpha = \lambda_\alpha X \bar{u}_\alpha \Rightarrow \bar{v}_\alpha \propto X \bar{u}_\alpha$$

Como μ_α es máximo y asociado a $\bar{v}_\alpha \Rightarrow \mu_\alpha \geq \lambda_\alpha$ (v)

De (b): $XX' \bar{v}_\alpha = \mu_\alpha \bar{v}_\alpha$ premultiplicando por X'

$$(X'X)X' \bar{v}_\alpha = \mu_\alpha X' \bar{v}_\alpha \Rightarrow \bar{u}_\alpha \propto X' \bar{v}_\alpha \quad \text{de (b)}$$

Como λ_α es máximo y asociado a $\bar{u}_\alpha \Rightarrow \lambda_\alpha \geq \mu_\alpha$ (vi)

De (v) y (vi):

$$\lambda_\alpha = \mu_\alpha$$

Sabemos que:

$$\bar{u}_\alpha \propto X' \bar{v}_\alpha$$

$$\bar{u}_\alpha = k X' \bar{v}_\alpha$$



$$\bar{u}_\alpha' \bar{u}_\alpha = k^2 \bar{v}_\alpha' X X' \bar{v}_\alpha$$

$$\bar{u}_\alpha' \bar{u}_\alpha = k^2 \mu_\alpha$$



$$K = 1 / \sqrt{\mu_\alpha}$$

Luego,

$$\sqrt{\mu_\alpha} \bar{u}_\alpha = X' \bar{v}_\alpha$$

$$\sqrt{\mu_\alpha} \bar{u}_\alpha \bar{v}_\alpha' = X' \bar{v}_\alpha \bar{v}_\alpha'$$

$$\sum_\alpha \sqrt{\mu_\alpha} \bar{u}_\alpha \bar{v}_\alpha' = \sum_\alpha X' \bar{v}_\alpha \bar{v}_\alpha'$$

$$X' \cong \sum_\alpha \sqrt{\mu_\alpha} \bar{u}_\alpha \bar{v}_\alpha'$$

----- comparar con (iv)



Definición

Conocido los vectores propios de un sub espacio se pueden obtener los del otro sin necesidad de una nueva factorización. Así tenemos,

$$V_{\alpha} = (1 / \lambda_{\alpha}) X^t u_a$$

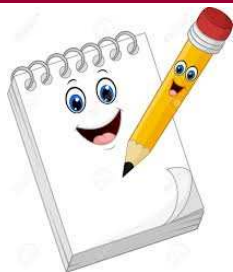
$$u_{\alpha} = (1 / \lambda_{\alpha}) X V_a$$

Definición

Existe una proporcionalidad entre las coordenadas de los puntos individuos sobre el eje factorial α en \mathcal{R}^p , \mathbf{xv}_a y las componentes del vector unitario director del eje α en el otro espacio \mathbf{V}_α

$$\mathbf{G}_\alpha = \mathbf{X}^t \mathbf{U}_\alpha = \sqrt{\lambda_\alpha} \mathbf{V}_\alpha$$

$$\mathbf{G}_\alpha(j) = \sum_i \mathbf{x}_{ij} u_{\alpha i} = \sqrt{\lambda_\alpha} \mathbf{v}_{\alpha i}$$



Análisis sobre las variables y sobre los individuos.

Ejemplo

Matriz de dispersión sobre las variables



$$X^T X$$

(p x p)

$$\frac{X^T X}{n - 1}$$

Matriz de dispersión sobre los individuos



$$X X^T$$

(n x n)

$$\frac{X X^T}{p - 1}$$

variables

$$X = \begin{bmatrix} -2 & -2 & -1 \\ 0 & 1 & 0 \\ 2 & 2 & 2 \\ 0 & -1 & -1 \end{bmatrix}$$

observaciones

variables

$$\mathbf{X}'\mathbf{X} = \begin{bmatrix} 8 & 8 & 6 \\ 8 & 10 & 7 \\ 6 & 7 & 6 \end{bmatrix}$$

$$\text{Tr}(\mathbf{X}'\mathbf{X})=24$$

Raíces características de $\mathbf{X}'\mathbf{X}$:

$$l_1 = 22.2819$$

$$l_2 = 1.0000$$

$$l_3 = .7181$$

 $\mathbf{X}'\mathbf{X}$: es definida positiva de rango 3

observaciones

$$\mathbf{X}\mathbf{X}' = \begin{bmatrix} 9 & -2 & -10 & 3 \\ -2 & 1 & 2 & -1 \\ -10 & 2 & 12 & -4 \\ 3 & -1 & -4 & 2 \end{bmatrix}$$

$$\text{Tr}(\mathbf{X}\mathbf{X}')=24$$

Raíces características de $\mathbf{X}\mathbf{X}'$:

$$l_1 = 22.2819$$

$$l_2 = 1.0000$$

$$l_3 = .7181$$

$$l_4 = 0$$

 $\mathbf{X}\mathbf{X}'$: es semidefinida positiva de rango 3

Vectores característicos de $X'X$:

$$U = \begin{bmatrix} -.574 & .816 & -.066 \\ -.654 & -.408 & .636 \\ -.493 & -.408 & -.768 \end{bmatrix}$$

$$U'U=I$$

Vectores característicos de XX' :

$$U^* = \begin{bmatrix} .625 & -.408 & -.439 \\ -.139 & -.408 & .751 \\ -.729 & 0 & -.467 \\ .243 & .816 & .156 \end{bmatrix}$$

$$U^{*'}U^*=I$$

Valores característicos L de $X'X$ y XX' :

$$L = \begin{bmatrix} 22.2819 & 0 & 0 \\ 0 & 1.0000 & 0 \\ 0 & 0 & .7181 \end{bmatrix}$$

$$L^{1/2} = \begin{bmatrix} 4.7204 & 0 & 0 \\ 0 & 1.0000 & 0 \\ 0 & 0 & .8474 \end{bmatrix}$$

$$\mathbf{V} = \mathbf{U}\mathbf{L}^{1/2} = \begin{bmatrix} -2.707 & .816 & .056 \\ -3.089 & -.408 & -.539 \\ -2.326 & -.408 & .651 \end{bmatrix} \quad \mathbf{V}^* = \mathbf{U}^*\mathbf{L}^{1/2} = \begin{bmatrix} 2.9549 & -.408 & -.372 \\ -.656 & -.408 & .636 \\ -3.441 & 0 & -.396 \\ 1.147 & .816 & .132 \end{bmatrix}$$

Escalamiento de los vectores característicos por $\mathbf{L}^{1/2}$

Se verifica que:

$$\mathbf{V}\mathbf{V}' = \mathbf{X}'\mathbf{X} \text{ y } \mathbf{V}^*\mathbf{V}^{*'} = \mathbf{X}\mathbf{X}' \quad \text{-----} \quad (\text{i})$$

Debido a la descomposición espectral:

$$\mathbf{X}'\mathbf{X} = \mathbf{U}\mathbf{L}\mathbf{U}' \text{ y } \mathbf{X}\mathbf{X}' = \mathbf{U}^*\mathbf{L}\mathbf{U}^{*'} \quad \text{-----} \quad (\text{ii})$$




$$\begin{matrix} \mathbf{W} = \mathbf{U}\mathbf{L}^{-1/2} = \begin{bmatrix} -.122 & .816 & -.078 \\ -.139 & -.408 & .751 \\ -.104 & -.408 & -.906 \end{bmatrix} \\ \text{(p} \times \text{p)} \end{matrix} \quad \begin{matrix} \mathbf{W}^* = \mathbf{U}^*\mathbf{L}^{-1/2} = \begin{bmatrix} .132 & -.408 & -.518 \\ -.029 & -.408 & .886 \\ -.154 & 0 & -.552 \\ .051 & .816 & .184 \end{bmatrix} \\ \text{(n} \times \text{p)} \end{matrix}$$

Escalamiento de los vectores característicos por $\mathbf{L}^{-1/2}$

Se verifica (i) y (ii) para \mathbf{W} y \mathbf{W}^*


Scores de las observaciones (CP):


$$\mathbf{Y} = \mathbf{XW} = \mathbf{U}^*$$

(n x p)

“Los scores de las n observaciones (CP) son iguales a los vectores característicos obtenidos de la matriz suma de cuadrados de observaciones – ver (ii)”

Scores de las variables (CP):


$$\mathbf{Y}^* = \mathbf{X}'\mathbf{W}^* = \mathbf{U}$$

(p x p)

“Los scores de las p variables (CP) son iguales a los vectores característicos obtenidos de la matriz suma de cuadrados de variables – ver (ii)”

DESCOMPOSICIÓN POR VALOR SINGULAR (DVS)

Los datos originales se pueden obtener de las CP:

$$\mathbf{X} = \mathbf{YV}' \quad (\text{iii})$$

O también en términos de las transformaciones:

*Descomposición por
valor singular (DVS)*


$$\mathbf{X} = \mathbf{U} * \mathbf{L}^{1/2} * \mathbf{U}'$$

(iv)

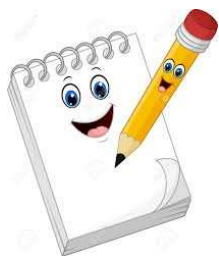
En DVS una matriz de datos \mathbf{X} se descompone como el producto de vectores característicos de $\mathbf{X}'\mathbf{X}$, los vectores característicos de \mathbf{XX}' y una función de sus raíces características \mathbf{L} .

$$\mathbf{X} = \mathbf{Y}\mathbf{L}^{1/2}\mathbf{U}' = \mathbf{U} * \mathbf{L}^{1/2} * \mathbf{Y}'$$

RESULTADOS IMPORTANTES

- 1 Ninguna combinación lineal estandarizada de $x_{(p \times 1)}$ tiene varianza mayor que λ_1 , la varianza del primer componente principal
- 2 Si la matriz de covarianzas de \mathbf{x} tiene rango “ r ” ($r < p$), entonces, la variación de \mathbf{x} puede ser totalmente explicada por las primeras “ r ” componentes principales.
- 3 Si $Y_1 = e_1 X$, $Y_2 = e_2 X$, ..., $Y_p = e_p X$ son los componentes principales obtenidos a partir de Σ , entonces,

$$\rho_{Y_i, X_k} = \frac{e_{ki} \sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}}$$



Ejemplo

Se realiza un estudio cuyo objetivo es poner de relieve los factores que diferencian al máximo las marcas entre si y determinar aquellas marcas que el conjunto de encuestados considera semejantes.

Marcas	Características		
	Elegancia	Comodidad	Deportivo
A	2	3	6
B	3	2	4
C	4	5	4
D	5	5	4
E	8	9	6
F	9	7	7

Paso 1.- Tipificado de datos

$$X_{ij} = \frac{z_{ij} - \bar{z}_j}{s_j \sqrt{n}}$$
$$i = 1, \dots, 6$$
$$j = 1, \dots, 3$$
$$n = 6$$

Así obtenemos,

$$X = \begin{pmatrix} -0.51 & -0.38 & 0.28 \\ -0.35 & -0.55 & -0.39 \\ -0.19 & -0.03 & -0.39 \\ -0.03 & -0.03 & -0.39 \\ -0.46 & 0.67 & 0.28 \\ 0.62 & 0.03 & 0.63 \end{pmatrix}$$

Paso 2.- Cálculo de la Matriz de Correlación

$$\mathbf{C} = \mathbf{X}^t \mathbf{X} = \begin{pmatrix} 1 & 0.89 & 0.58 \\ 0.89 & 1 & 0.52 \\ 0.58 & 0.52 & 1 \end{pmatrix}$$

Paso 3.- Cálculo de la Inercia Explicada

$$\begin{vmatrix} 1 - \lambda & 0.89 & 0.58 \\ 0.89 & 1 - \lambda & 0.52 \\ 0.58 & 0.52 & 1 - \lambda \end{vmatrix} = 0 \quad \Rightarrow \quad \begin{aligned} \lambda_1 &= 2.334 \\ \lambda_2 &= 0.560 \\ \lambda_3 &= 0.106 \end{aligned}$$

Se verifica que:

$$\text{Tr}(\mathbf{C}) = \lambda_1 + \lambda_2 + \lambda_3 = 3$$



Componentes Principales e Inercia

Componente i-ésimo	Inercia Total (λ_i)	Inercia % (λ_i / p)*100%	r_p
1	2.334	78.16%	78.16
2	0.560	13.36%	91.52
3	0.106	8.48%	100.00

Paso 4.- Obtención de las Representaciones Gráficas

Dirección del Primer Eje

Sea \mathbf{V}_α el vector que sigue la dirección del primer eje, entonces,

$$\mathbf{V}_\alpha = \begin{pmatrix} v_{11} \\ v_{12} \\ v_{13} \end{pmatrix} \quad \text{donde} \quad v_{11}^2 + v_{12}^2 + v_{13}^2 = 1 \quad \text{..... (i)}$$

y cumple que: $(\mathbf{C} - \lambda_1 \mathbf{I}) \mathbf{V}_\alpha = 0 \quad \text{..... (ii)}$

De (i) y (ii) se obtiene que: $\mathbf{V}_\alpha = \begin{pmatrix} 0.622 \\ 0.603 \\ 0.505 \end{pmatrix}$

Paso 5.- Proyección de las Marcas sobre el Primer Eje:

$$F_1 = \sqrt{\frac{n}{p}} \times \mathbf{V}_1 \quad \text{donde} \quad \sqrt{\frac{n}{p}} = \sqrt{\frac{6}{3}}$$

$$F_1 = \sqrt{\frac{6}{3}} \begin{pmatrix} -0.51 & -0.38 & 0.28 \\ -0.35 & -0.55 & -0.39 \\ -0.19 & -0.03 & -0.39 \\ -0.03 & 0.03 & 0.39 \\ 0.46 & 0.67 & 0.28 \\ 0.62 & 0.03 & 0.63 \end{pmatrix} \begin{pmatrix} 0.622 \\ 0.603 \\ 0.505 \end{pmatrix} = \begin{pmatrix} -0.57 \\ -0.06 \\ -0.47 \\ -0.33 \\ 1.17 \\ 1.25 \end{pmatrix}$$

$$F_2 = \sqrt{\frac{n}{p}} \times \mathbf{V}_2 = \begin{pmatrix} 0.77 \\ -0.01 \\ -0.38 \\ -0.45 \\ -0.24 \\ 0.31 \end{pmatrix}$$



Para obtener las coordenadas de los puntos variables (características), haremos uso de la relación de ambos espacios:

$$\mathbf{G}_1 = \mathbf{X}^t \mathbf{e}_1 = \sqrt{\lambda_1} \mathbf{v}_1$$

De donde,

$$\mathbf{G}_1 = \begin{bmatrix} G_1(\text{elegancia}) \\ G_1(\text{comodidad}) \\ G_1(\text{deportivo}) \end{bmatrix} = \sqrt{2.347} \begin{bmatrix} 0.622 \\ 0.603 \\ 0.505 \end{bmatrix} = \begin{bmatrix} 0.95 \\ 0.92 \\ 0.77 \end{bmatrix}$$

$$\mathbf{G}_2 = \begin{bmatrix} G_2(\text{elegancia}) \\ G_2(\text{comodidad}) \\ G_2(\text{deportivo}) \end{bmatrix} = \begin{bmatrix} -0.22 \\ -0.31 \\ 0.64 \end{bmatrix}$$



CASO: CONTAMINACIÓN AMBIENTAL

Datos: Airpollution.txt

Tomado de: Multivariate Analysis of Variance. Johnson&Wichern

Variables

TMR: Tasa de mortalidad total

SMIN: lectura de sulfato quincenal más pequeña

SMEAN: lectura de sulfato quincenal promedio

SMAX: lectura de sulfato quincenal máximo

PMIN: lecturas de partículas suspendidas cada dos semanas - mínimas.

PMEAN: lecturas de partículas suspendidas cada dos semanas - promedio.

PMAX: lecturas de partículas suspendidas cada dos semanas - máximas.

PM2: Densidad de población por milla cuadrada x 0.1

GE65: Porcentaje de POBLACIÓN al menos 65 x 10

PERCWH: Porcentaje de blancos en población

NONPOOR: Porcentaje de familias con ingresos por encima del nivel de pobreza

LPOP: Logaritmo (base 10) de la población x 10.