

## ANALISIS MULTIDIMENSIONAL DE DATOS

Raul Gouet \*

## COMPONENTES PRINCIPALES REGRESION LINEAL.-

El propósito de este cursillo es presentar dos métodos clásicos del análisis estadístico de datos multivariados, intentando un equilibrio razonable entre teoría y aplicaciones. La teoría, que espero despierte interés en aquellos con vocación matemática, usa herramientas clásicas de Algebra Lineal y Probabilidades, contenidas normalmente en los programas de carreras científico-tecnológicas.

Completando la teoría, veremos también elementos de interpretación de resultados que luego aplicaremos a algunos conjuntos de datos en el microcomputador. Todos están invitados a cocinar las recetas teóricas, incluso los matemáticos.

\* Facultad de Ciencias Físicas y Matemáticas, Universidad de Chile.

## 1.- EL ANALISIS DE DATOS.-

En estudios de diversa naturaleza y origen, se presenta ge  
neralmente una etapa de mediciones empíricas que arrojan volúmen  
es importantes de datos. De ellos esperamos una mejor comprens  
ión del fenómeno en estudio.

El AD se ocupa del análisis estadístico de datos que sur-  
gen en muestras multivariadas. Su propósito primario es reducir  
grandes masas de números, produciendo una síntesis significativ  
iva y razonablemente completa de toda información que contenga  
la muestra. En este proceso de resumen de los datos, aceptamos  
perder información para ganar significación.

En general, la práctica del AD es descriptiva, en el sentid  
o de rechazar la imposición de un modelo probabilístico para  
los datos, por ejemplo normalidad. Esta posición, de aparente  
objetividad, impide estudios de estabilidad de resultados, a  
través de procedimientos clásicos de inferencia estadística.

Hay una gran variedad de técnicas y recetas que suelen llam  
arse de AD. Por una parte están los métodos lineales, que tien  
en su origen en la escuela anglo sajona durante la primera mid  
ad de este siglo, y que adquieren verdadera popularidad con el  
uso masivo de computadores.

La escuela francesa, en los años 60 y 70, remoza y formaliz  
a elegantemente estas técnicas, llamadas también factoriales,  
revelando por ejemplo la potencia del Análisis de Correspondenc  
ias.

La otra cara de la medalla reú  
ne una extensa fauna de métod  
os exóticos, ni lineales ni convencionales, sobre los cuales

hay pocos resultados teóricos pero abundante controversia.

En este curso abordaremos, con la perspectiva de la escuela francesa, los métodos de Análisis en Componentes Principales y Regresión Lineal.

#### 1.1.- POBLACION Y MUESTRA.-

Antes de detallar el formalismo matemático, vamos a precisar el sentido de algunos términos de uso corriente en estadística: Se entenderá que una población es una colección de objetos que pertenecen a una clase bien definida e identificable. Una muestra será cualquier subconjunto de la población.

La distinción entre muestra y población cobra interés, cuando buscamos hacer inferencias que generalicen las conclusiones válidas en la muestra hacia la población completa. Ello obliga en principio a diseñar cuidadosos esquemas de muestreo, pero la realidad suele apartarse de los buenos deseos del estadístico, y de la muestra a menudo no se sabe de donde viene ni como llegó. En ciertas disciplinas, como la Arqueología, el investigador debe conformarse con lo escasamente disponible.

Aunque habitualmente el objetivo final de un estudio es el conocimiento de la población, nuestro análisis estará centrado en la muestra, la que exploraremos desde diversos ángulos.

### 1.2.- VARIABLES E INDIVIDUOS.-

A los términos variable e individuo asignaremos un significado preciso. Un individuo es cualquier objeto o entidad que forme parte de la muestra. Una variable (también llamada carácter) es un atributo cuantitativo que poseen los individuos. Podemos imaginarla como una regla o aplicación que asocia a cada individuo un número real.

En un estudio se medirán típicamente muchas variables, sobre uno o más conjuntos de individuos. Estos números, dispuestos en arreglos rectangulares (matrices), llamadas tablas de datos, constituyen el punto de partida para el AD.

### 1.3.- COMPONENTES PRINCIPALES (ACP).-

Desarrollado por Hotelling en los años 30, el ACP es un típico representante de los métodos factoriales, que tiene actualmente gran popularidad, aunque se reconoce que no es ni el único ni el mejor método de análisis factorial disponible.

Desde el punto de vista de la representación de los individuos en un espacio vectorial, el objetivo del ACP es simplificar, reduciendo la dimensión del espacio, con la menor pérdida de información posible. Desde el punto de vista de las variables, el ACP pretende concentrar la información en una base de variables no correlacionadas, que tengan máxima importancia.

#### 1.4.- REGRESION LINEAL (RL).-

Los métodos de regresión son una panacea estadística. Todo el mundo tiene a su alcance, ya sea en calculadora o computador un programa que hace regresiones. Se usa la regresión con múltiples propósitos y frecuentemente se abusa.

Nadie quiere dejar de ajustar un modelo si tiene la oportunidad: para predecir, extrapolar, aproximar, explicar, descubrir una ley, verificar una hipótesis, etc.

Es la opinión de muchos que la RL es probablemente una de las técnicas más potentes para analizar datos, y se abusa de ella porque parece decirnos mucho más de lo que los datos son capaces de dar.

Los métodos de regresión se ocupan del problema general de describir relaciones entre variables (estadísticas), "explicando" una de ellas en términos de las restantes.

La RL se practica con variados objetivos: remover efectos de variables; descubrir una ley empírica; predecir comportamiento; controlar un sistema etc.

El tratamiento sistemático de los problemas relacionados con la regresión necesita dos o más cursos de un semestre. Aquí estaremos satisfechos mirando algunas "componentes principales" que sirvan como estímulo en la investigación personal posterior.

## 2.- DEFINICIONES BASICAS.-

Supongamos que la muestra consiste en  $n$  individuos  $i = 1, 2, \dots, n$  sobre los cuales se han medido  $p$  variables cuantitativas  $v = 1, 2, \dots, p$ .

Cada variable  $v$  en  $v$  puede verse como una aplicación

$$X^v : \mathcal{I} \longrightarrow \mathbb{R}$$

$$i \longmapsto X^v(i) = X_i^v$$

Siendo  $X_i^v$  el valor que toma la variable  $v$  sobre el individuo  $i$ .

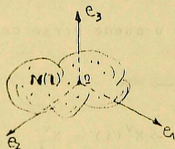
El arreglo rectangular de  $p$  filas y  $n$  columnas conteniendo la información relativa a las  $p$  variables sobre todos los individuos, lo llamaremos tabla de datos

$$X = \begin{array}{c} \longleftarrow i \longrightarrow \\ \left[ \begin{array}{c} X_i^1 \dots \\ X_i^2 \dots \\ \vdots \\ X_i^p \dots \end{array} \right] \begin{array}{l} \uparrow \\ v \\ \downarrow \end{array} \end{array}$$

Un individuo  $i$  en  $\mathcal{I}$  queda (para efectos del análisis) completamente caracterizado por los valores  $X_i^1, X_i^2, \dots, X_i^p$ . Esto permite representarlo en el espacio  $E = \mathbb{R}^p$  como el vector  $X^i$  cuyas coordenadas en la base canónica son los números  $X_i^v$ .

Es natural que el espacio  $E$  sea llamado Espacio de individuos.

La colección de vectores de  $E$  que representa a los individuos de la muestra se llama Nube de individuos:  $N(\mathcal{I}) = \{X_i; i \in \mathcal{I}\}$



Tenemos una primera representación de los datos en un espacio vectorial, desde la perspectiva de los individuos.

En ciertos problemas es útil (y necesario) asignar a cada individuo un número no negativo que refleje su importancia relativa. Esto se logra con pesos  $P_i$ , normalizados de manera que  $0 \leq P_i \leq 1$  y  $\sum_i P_i = 1$ . Cuando no hay forma de justificar un trato discriminatorio se escogen valores  $P_i = 1/n$ .

Consideremos ahora una variable  $v \in \mathcal{V}$ . Todo lo que sabemos de ella está en la colección de valores que toma en la muestra  $\mathcal{I}$ . Es entonces razonable asociarle el vector  $X^v$  en  $F = \mathbb{R}^n$ , cuyas coordenadas en la base canónica  $\{f_j; j \in \mathcal{I}\}$  son los datos  $X^v_1, X^v_2, \dots, X^v_n$ .

El espacio  $F$  se conoce como Espacio de variables y la colección  $N(\mathcal{V}) = \{X^v; v \in \mathcal{V}\}$  es la nube de variables.

## 2.1.- REPRESENTACION EN ESPACIOS DUALES.-

Sean  $E^*$ ,  $F^*$  los espacios duales de  $E$  y  $F$  respectivamente. Para designar la imagen  $u(x)$  de un vector  $X$  en  $E$  (Resp.  $F$ ) por la forma lineal  $u$  en  $E^*$  (resp.  $F^*$ ) escribiremos  $\langle u, x \rangle$  o  $\langle x, u \rangle$ , identificando  $E$  con el doble dual  $E^{**}$ .

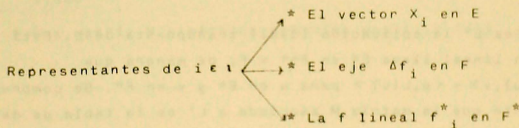
Supongamos que  $E^*$  y  $F^*$  están dotados de las bases duales  $\{e_v^*; v \in v\}$  y  $\{f_i^*; i \in i\}$ . Estas bases se caracterizan porque a asignan a un vector sus coordenadas en las bases canónicas. Por ejemplo  $\langle e_v^*, X_i \rangle = X_i^v$  y  $\langle f_j^*, X_i^v \rangle = X_i^v$ .

Lo anterior revela que  $e_v^*$  representa efectivamente a la variable  $v$  puesto que asocia al individuo  $i$ , a través de su vector  $X_i$ , el valor  $X_i^v = X^v(i)$ . Con idéntico argumento podemos convencernos que  $f_i^*$  representa al individuo  $i$ .

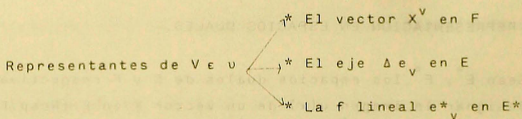
Finalmente, con un poco de imaginación aceptaremos que los ejes de  $E$  representan variables y los de  $F$  individuos.

Para designar al eje (subespacio vectorial (sev) de dimensión 1) engendrado por el vector  $u$ , escribimos  $\Delta u$ .

A manera de resumen tenemos:







## 2.2.- LAS APLICACIONES $X$ y $X'$ .-

Veremos como se relacionan las distintas representaciones de individuos o variables y el papel que tiene la tabla de datos en dichas relaciones.

Pongamos en correspondencia cada  $f^*$ , en  $F^*$  con  $X_i$  en  $E$ . Esto define la aplicación lineal  $L : F^* \rightarrow E$  caracterizada por  $L(f^*_i) = X_i$ .

Esta transformación tiene asociada la matriz  $M$  con respecto a las bases respectivas de  $F^*$  y  $E$ , cuyas columnas son las coordenadas en la base de  $E$ , de las transformadas de la base de  $F^*$ :  $L(f^*_i) = X_i = \sum_V X^V_i e_V$ .

Luego, la  $i$ -ésima columna de  $M$  está formada por los números  $X^1_i, X^2_i, \dots, X^D_i$ , de la  $i$ -ésima columna de la tabla de datos  $X$ . Por lo tanto,  $X$  es la matriz asociada a  $L$ .

Sea  $L'$  la aplicación (dual) transpuesta de  $L$ . Esta transformación lineal lleva  $E^*$  en  $F^{**} = F$ , de manera que  $\langle L'(u), v \rangle = \langle u, L(v) \rangle$  para  $u$  en  $E^*$  y  $v$  en  $F^*$ . Se comprueba fácilmente que la matriz  $M$  asociada a  $L'$  es la tabla de datos "transpuesta"  $X'$ .

## 2.3.- METRICAS.-

Las acciones de clasificar y comparar, requieren una noción de proximidad entre los objetos estudiados. Por esta razón debemos metrizar los espacios  $E$  y  $F$ .

Para disponer del concepto de ángulo trabajaremos con métricas euclidianas, es decir, que se originen en un producto interno.

A una forma bilineal simétrica definida positiva

$M : E \times E \rightarrow \mathbb{R}, (x, y) \rightarrow M(x, y)$  podemos asociarle:

\* Una forma cuadrática  $M : E \rightarrow \mathbb{R}, x \rightarrow M(x) = M(x, x)$

\* Un isomorfismo  $M : E \rightarrow E^*, x \rightarrow M(x)$  en  $E^*$   
 $\langle M(x), y \rangle = M(x, y)$

\* Una norma euclidiana  $\|x\|_M = \sqrt{M(x)}$

\* Una distancia euclidiana  $d(x, y) = \|x - y\|_M$

\* Una noción de  $M$ -ortogonalidad  $x$  y  $y$  si  $M(x, y) = 0$

\* Una matriz  $M$  simétrica definida positiva (dp)

## 2.4.- METRICAS EUCLIDIANAS EN E, E\*, F, F\*.-

Sea M una métrica sobre E. Para cada par de puntos de la nube  $N(t)$  calculamos su distancia como  $\|x_i - x_k\|_M$ . Dado que los individuos también se representan en  $F^*$ , es natural pensar que la transformación X induce una métrica "heredada de M" que llamaremos W. Para el par de individuos i, k pedimos que se verifique que  $\|x_i - x_k\|_M = \|f_i^* - f_k^*\|_W$  o más generalmente  $\|a - b\|_W = \|X(a) - X(b)\|_M$  para a, b en  $F^*$  cualesquiera.

Puesto de manera equivalente tendremos:

$$\|a\|_W = \|X(a) - X(b)\|_M \text{ para cualquier } a \text{ en } F^*, \text{ o bien}$$

$$W(a, b) = M(X(a), X(b)) \quad a, b \text{ en } F^*$$

$$\text{Pero } M(X(a), X(b)) = \langle M'X(a), X(b) \rangle$$

$$= \langle X' \cdot M'X(a), b \rangle$$

$$\text{luego } \langle W(a), b \rangle = \langle X' \cdot M'X(a), b \rangle \text{ para todo } a, b \text{ en } F^*$$

$$\text{por lo tanto } W = X' \cdot M'X$$

En el espacio de variables procedentes de manera análoga. Dada una métrica euclidiana N sobre F construimos la métrica V sobre  $E^*$  a través de la aplicación  $X'$ . Para ello imponemos la condición  $\|a - b\|_V = \|X'(a) - X'(b)\|_N$  para todos a, b en  $E^*$ , es decir,

$$V(a, b) = N(X'(a), X'(b))$$

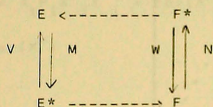
$$= \langle N'X'(a), X'(b) \rangle$$

$= \langle X'N'X'(a), b \rangle$  para todos  $a, b$  en  $E^*$ , luego

$$V = X'N'X'$$

## 2.5.- ESQUEMA DE DUALIDAD.-

Los espacios y transformaciones con los que hemos complicado progresivamente el panorama, son presentados en el esquema siguiente:



Donde podemos apreciar como las métricas sobre  $E$  y  $F$  inducen de manera natural las estructuras euclidianas sobre los duales.

### NOTA:

Las formas bilineales  $V$  y  $W$  serán dp o semi dp según si  $X$  es inyectiva o no. En efecto, si  $a \in \ker(X)$  entonces  $W(a, a) = M(X(a), X(a)) = 0$ .

El rango de  $W$  es el mismo que el de  $X$ . Como generalmente tendremos más individuos que variables ( $n > p$ ), y  $n = \dim(\ker(X)) + \text{rango}(x)$ , resulta  $\dim(\ker(X)) \geq n - p > 0$  luego  $W$  no es dp.

Razonando de manera idéntica se verifica que  $V$  es dp si  $\text{rango}(x') = p$  lo que equivale a pedir que no haya variables

colineales.

## 2.6.- METRICA DE LOS PESOS EN F.-

Hay una métrica para las variables que permite interpretar geoméricamente medias, varianzas y covarianzas. Para calcular la distancia entre dos variables se consideran las diferencias en cada individuo ponderado por su respectivo peso.

Esta métrica tiene asociada la matriz diagonal  $D_p$  dada por:

$$D_p = \begin{bmatrix} P_1 & & & & \\ & \cdot & & & \\ & & \cdot & & \\ & & & \cdot & \\ & & & & P_n \end{bmatrix}$$

La media aritmética de la variable  $v$  en  $U$  se calcula como

$$me(x^v) = \sum_i P_i X_i^v \quad v \text{ en } U$$

Sea  $\mathbf{1}$  el vector de  $F$  de componentes  $1,1,1,\dots,1$  en la base canónica. Entonces  $D_p(X^v, \mathbf{1}) = \sum_i P_i X_i^v$ ,  $v$  en  $U$ .

Por otra parte,  $\|\mathbf{1}\|_{D_p} = \sum P_i = 1$ , luego  $\mathbf{1}$  es vector unitario y  $me(X^v)$  puede interpretarse como la norma de la proyección  $D_p$ -ortogonal de  $X^v$  sobre el eje  $\Delta \mathbf{1}$ . Puesto que la proyección ortogonal es el vector a mínima distancia,  $me(x^v)\mathbf{1}$  en  $F$  es la variable constante que mejor aproxima a  $X^v$  en el sentido de la métrica  $D_p$ .

La variable  $X^v - m_e(X^v)_1$ , que es la proyección sobre el ortogonal  $\Delta^1_1$ , se llama centrada porque su media es nula.

Para mayor comodidad vamos a suponer en adelante que todas las variables de la tabla de datos son centradas, es decir  $Dp(X^v, 1) = 0$  para todo  $v$  en  $u$ . Veremos luego que esto equivale a fijar el origen del sistema de coordenadas de  $E$  en el centro de gravedad de la nube  $N(1)$ .

Sean  $X^v$  y  $X^s$  dos variables (centradas). Entonces

$$Dp(X^v, X^s) = \sum_i P_i X^v_i X^s_i = \text{cov}(X^v, X^s) \text{ y } \|X^v\|_{Dp}^2 = Dp(X^v, X^v) =$$

$$\sum_i P_i (X^v_i)^2 = \text{var}(X^v)$$

Si miramos el esquema de dualidad con  $N = Dp$  encontramos que la métrica  $V$  sobre el dual  $E^*$  tiene asociada la matriz de covarianzas  $X'DpX'$ .

## 2.7.- CENTRO DE GRAVEDAD E INERCIA.-

Como si se tratara de un sistema de partículas, el centro de gravedad de la nube de individuos  $N(1)$  se calcula como

$$g = \sum_i P_i X_i$$

Y veremos que las coordenadas de  $g$  son las medias:

$g^v = \sum_i P_i X^v_i = m_e(x^j)$ . Luego, las variables que describen la muestra son centradas sii  $g = 0$ , es decir, el origen de  $E$  esta en  $g$ .

Continuando con la analogía física definamos el momento de inercia de  $N(\mathcal{V})$  con respecto a un punto  $a$  en  $E$  como:

$$I_a = \sum_i p_i \|x_i - a\|_M^2$$

Este valor lo interpretamos como una medida de la deformación que debemos introducir para proyectar la nube  $N(\mathcal{V})$  sobre el punto  $a$ .

$N(\mathcal{V})$  se concentra en  $a$  si  $I_a = 0$ .

Un resultado importante sobre la inercia se obtiene del siguiente cálculo.

$$\begin{aligned} \|x_i - a\|_M^2 &= \|x_i - g - (g - a)\|_M^2 \\ &= \|x_i - g\|_M^2 + \|g - a\|_M^2 - 2M(x_i - g, g - a) \end{aligned}$$

Sumando sobre  $\mathcal{V}$  tenemos  $\sum_i p_i \|x_i - a\|_M^2 = \sum_i p_i \|x_i - g\|_M^2 + \|g - a\|_M^2 - 2M(\sum_i p_i (x_i - g), g - a)$ , para llegar a :

$$I_a = I_g + \|g - a\|_M^2$$

Lo que significa que  $g$  es el punto con respecto al cual la inercia de  $N(\mathcal{V})$  es mínima. Podemos imaginar que  $g$  es el subespacio (afín si las variables no fueran centradas) de dimensión cero que mejor aproxima a la nube.

### 3.0.- COMPONENTES PRINCIPALES.-

El ACP es un método que busca reducir la complejidad de la representación de los datos. Desde el punto de vista de los individuos, nos interesamos en proyectar la nube  $N(1)$  sobre un SEV (afín) de baja dimensión, de tal manera que se preserve toda o gran parte de la estructura de ángulos y proximidades de la representación original.

¿Por qué reducir la dimensión?. Simplemente porque somos incapaces de ver en dimensión mayor que 3.

Desde la perspectiva de las variables, el ACP intenta sustituirlas por una colección reducida de descriptores no correlacionados, que capturen la mayor parte de la variabilidad de la muestra (máxima varianza).

Nos ubicamos en el espacio  $E$  de los individuos para plantear una sucesión de problemas que corresponden a reducciones óptimas de los datos, en orden creciente de complejidad.

P0. Determinar el punto más cercano a  $N(1)$ . Es decir, el sev (afín) de dimensión cero, óptimo para proyectar la nube. Averiguar si  $N(1)$  se concentra en dicho punto.

P1. Determinar la recta más cercana a  $N(1)$ . Es decir, el sev (afín) de dimensión uno que provoque la menor distorsión al proyectar la nube. Averiguar si  $N(1)$  se concentra en dicha recta.



.  
.  
.

Pk. Determinar el mejor sev (afín) de dimensión  $k$  para proyectar la nube, etc.

Usando el criterio de la inercia, la solución de PO ya la encontramos en el centro de gravedad  $g$ . Necesitamos extender este criterio a los sev para resolver P1, P2,...

Sea  $W \subset E$  un sev (no afín).  $E = W \oplus W'$ . Un vector  $x_i$  de la nube  $N(\nu)$  se descompone de manera única como  $x_i = a_i + b_i$  donde  $a_i \in W$  y  $b_i \in W'$ .

Si proyectamos  $x_i$   $M$ -ortogonalmente sobre  $W$  hay una pérdida equivalente a la norma de la parte de  $x_i$  ortogonal a  $W$ , es decir,  $\|b_i\|_M^2$ . Sumando las pérdidas individuales adecuadamente ponderadas, llegamos a la inercia de  $N(\nu)$  con respecto a  $W$ .

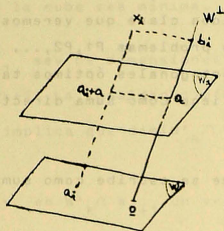
$$I_W = \sum_i p_i \|b_i\|_M^2$$

Es evidente que  $I_W = 0$  si  $N(\nu)$  está contenida en  $W$ .

Del teorema de Pitágoras ( $\|x_i\|^2 = \|a_i\|^2 + \|b_i\|^2$ ) recuperamos una identidad bastante útil:

$$I_g = I_W + I_{W'} \quad (g=0)$$

La inercia con respecto a un subespacio afín  $W_1$  (no contiene al cero) se define de manera muy natural. Sea  $W$  un sev de  $E$ ,  $a \in W'$  y  $W_1 = W + a$ . Entonces  $W_1 \cap W' = \{a\}$  y



$x_i = a_i + b_i = (a_i + a) + (b_i - a)$   $a_i \in W, b_i \in W'$ , luego  
 $I_{W_1} = \sum_i p_i \|b_i - a\|_M^2$ . Usando la identidad  
 $\|b_i - a\|^2 = \|b_i\|^2 + \|a\|^2 - 2M(b_i, a)$  llegamos a  $I_{W_1} = I_W +$   
 $+ \|a\|^2 - 2M(\sum_i p_i b_i, a)$  pero  $g = 0 = \sum_i p_i x_i = \sum_i p_i a_i + \sum_i p_i b_i$ .  
 Finalmente

$$I_{W_1} = I_W + \|a\|^2 \geq I_W$$

Esto significa que  $W$  es mejor ( en términos de menor distorsión ) que cualquier ser afín paralelo. Por lo tanto, los problemas  $P_1, P_2, \dots, P_k$  tendrán por solución sólo sev (conteniendo al cero =  $g$ ).

## 3.1.- DESCOMPOSICION DE LA INERCIA.-

Un par de resultados clave que veremos a continuación, permitirán reemplazar los problemas  $P_1, P_2, \dots$  por la búsqueda de una sucesión de ejes ortogonales óptimos tales que el sev  $W_k$ , solución de  $P_k$ , se obtiene como suma directa de los primeros  $k$  ejes.

Sea  $W$  sev de  $E$  que se escribe como suma directa de dos sev  $M$ -ortogonales  $w_1$  y  $w_2$ :

$$W = w_1 \oplus w_2 \quad w_1' \quad w_2' \quad y$$

$E = W \oplus W' = w_1 \oplus w_2 \oplus W'$ , luego  $x_i = a_i + b_i = c_i + d_i + b_i$  con  $a_i \in W$ ,  $b_i \in W'$ ,  $c_i \in w_1$ ,  $d_i \in w_2$ .

Entonces

$$Iw' = \sum_i p_i \|a_i\|^2 = \sum_i p_i \|c_i\|^2 + \sum_i p_i \|d_i\|^2$$

$$Iw' = Iw_1' + Iw_2'$$

Si preferimos una expresión en términos de  $Iw$ , podemos escribir  $Iw = Ig - Iw_1' - Iw_2'$ .

El subespacio óptimo de dimensión  $k + 1$  contiene un subespacio óptimo de dimensión  $k$ .

Esto significa que el sev  $W_{k+1}$ , solución de  $P_{k+1}$ , puede sustituirse a partir de  $W_k$  y un eje adicional, con respecto al cual la inercia de la nube sea mínima.

Sean  $W_k$  y  $W_{k+1}$  sev de dimensiones respectivas  $k$  y  $k+1$ , soluciones de  $P_k$  y  $P_{k+1}$ . La igualdad  $\dim(W'_k) + \dim(W_{k+1}) = (p-k) + (k+1) = p+1$  implica que  $\dim(W'_k \cap W_{k+1}) \geq 1$ .

Sea entonces  $v$  en  $W'_k \cap W_{k+1}$  un vector no nulo y descompongamos  $W_{k+1}$  en suma directa como  $W_{k+1} = \Delta v \oplus R$ , donde  $R$  (de dim  $k$ ) es el suplementario  $M$ -Ortogonal de  $\Delta v$  en  $W_{k+1}$ .

Definamos también el sev  $u$  (de dim  $k+1$  porque  $v \in W'_k$ ) como  $u = \Delta v \oplus W_k$  y comparemoslo con  $W_{k+1}$  en términos de la inercia.

Aplicando las fórmulas recientemente deducidas para la descomposición de la inercia llegamos a:

$$I_{W_{k+1}} = I_g - I \Delta'v - I r' \quad y$$

$$I_u = I_g - I \Delta'v - I w'_k$$

Puesto que  $W_k$  es el sev óptimo de dimensión  $k$ , tenemos  $I w'_k \geq I r'$  luego, de las ecuaciones anteriores deducimos que

$$I_{W_{k+1}} \geq I_u$$

Pero  $W_{k+1}$  es óptimo y por lo tanto debe verificarse la igualdad  $I_{W_{k+1}} = I_u$ , que implica finalmente  $I_R = I_{W_k}$ .

Hemos comprobado que  $W_{k+1}$  contiene un sev óptimo de dimensión  $k$  y que además, en virtud de las relaciones

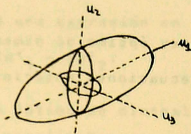
$$u = \Delta v \otimes W_{k+1} \quad \text{y} \quad I_u = I_g - I \Delta' v - I_{W'}$$

es posible construir  $W_{k+1}$  "pegando" a  $W_k$  un eje  $\Delta v$  ortogonal a  $W_k$ , de mínima inercia. Resulta claro ahora que los problemas  $P_1, P_2, \dots$  son equivalentes a la siguiente secuencia de tareas:

$T_1$  : Determinar  $\Delta u_1$  más cercano a  $N(1)$  (mínima inercia)

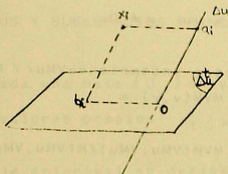
$T_2$  : Determinar  $\Delta u_2$  de mínima inercia,  $M$ -ortogonal a  $\Delta u_1$ .

$TK$ . Determinar  $\Delta u_k$  de mínima inercia,  $M$ -ortogonal a los ejes  $\Delta u_1, \Delta u_2, \dots, \Delta u_{k-1}$  (es decir, al sev engendrado por  $u_1, u_2, \dots, u_k$ )



## 3.2.- CALCULOS EXPLICITOS.-

Ahora comenzamos a resolver concretamente las tareas  $T_1, T_2, \dots$



$$x_i = a_i + b_i \quad I \Delta^t u = \sum_i p_i \|a_i\|^2$$

Asociemos a  $u$  una forma lineal  $v = Mu$  en  $E^*$  y una variable  $c = X'v = X'Mu$  en  $F$  (ver esquema de dualidad). Entonces

$$a_i = M(x_i, u)u \quad y \quad M(x_i, u) = \langle x_i, Mu \rangle = \langle Xf^*_i, Mu \rangle = \langle f^*_i, X'Mu \rangle = \langle f^*_i, c \rangle.$$

$$\begin{aligned} \text{Luego } c &= \sum_i M(x_i, u)f_i \quad y \quad \text{además } I \Delta^t u = \sum_i p_i M(x_i, u)^2 = \\ \|c\|_{Dp}^2 &= \langle DpC, C \rangle = \langle DpX'Mu, X'Mu \rangle = \langle X DpX'Mu, Mu \rangle \\ &= \langle VMu, Mu \rangle = M(u, VMu) \\ &= \langle MVMu, u \rangle = MVM(u, u). \end{aligned}$$

Llegamos así a la fórmula

$$I \Delta^t u = M(u, VMu) = M(u, VMu) = MVM(u, u)$$

que nos indica que  $N(t) C \Delta^t u$  sii  $VMu = 0$ , es decir, si  $u$  es vector propio de  $VM$  asociado al valor propio cero.

Supongamos que  $N(1)$  no está contenida en  $\Delta'u$  (hay parte de la nube que se estira a lo largo de  $\Delta u$ ). Aplicando la desigualdad de Schwarz veremos que el eje engendrado por  $VMu$  es mejor que  $\Delta u$ , en el sentido que  $I\Delta VMu \leq I\Delta u$ . Luego, un eje inmejorable tendrá que ser invariante frente a  $VM$ . En otras palabras  $I\Delta VMu = I\Delta u$  sii  $VMu = \mu u$   $\mu \neq 0$ .

Consideremos el vector unitario  $v = VMu / \sqrt{M(VMu, VMu)}$ . Es claro que  $I\Delta'VMu = I\Delta'v = MVM(v, v) =$

$$= MVM(VMu, VMu) / M(VMu, VMu)$$

$$= MVM(VMu, VMu) / MVM(u, VMu)$$

Pero

$$I\Delta'u = M(u, VMu) \leq M(VMu, VMu) / M(u, VMu) =$$

$$= MVM(u, VMu) / MVM(u, u) \leq$$

$$\leq MVM(VMu, VMu)MVM(u, u) / MVM(u, u)MVM(u, VMu) = I\Delta'v$$

lo que equivale a  $I\Delta VMu \leq I\Delta u$ .

Si  $\Delta u$  es un eje inmejorable, es decir, si  $I\Delta VMu = I\Delta u$  entonces  $M(u, VMu)^2 = M(VMu, VMu)$ , igualdad que se cumple sii  $VMu = \mu u$   $\mu \neq 0$ . Además  $I\Delta'u = M(u, VMu) = \mu M(u, u) = \mu$

La inercia  $I\Delta u$  es mínima si  $u$  pertenece al sev propio asociado al mayor valor propio de  $VM$ .

Hemos ligado el problema de los ejes de mínima inercia con el estudio de valores y vectores propios de  $VM$ . Es útil mencio-

nar ahora que VM tiene valores propios reales no negativos  $\mu_1 \geq \mu_2 \geq \dots \mu_p \geq 0$ . Además podemos construir una base M-ortonormal  $\{u_j ; j = 1, 2, \dots, p\}$  de vectores propios de VM.

### 3.3.- EJES, PLANOS Y SUBESPACIOS PRINCIPALES.-

Suponemos dada una base  $\{U_v\}$  de vectores propios de VM asociados a los valores propios  $\mu_1 \geq \mu_2 \geq \dots \geq \mu_p \geq 0$ .

El primer eje principal se define como el sev de dimensión 1  $\Delta u$  tal que  $I \Delta u$  es mínima.

Escribiendo  $u$  en la base de vectores propios, tendremos  $u = \sum_v a_v U_v$  y  $\sum_v (a_v)^2 = 1$ . Además  $I \Delta' u = M(u, VMu) = \sum_v \mu_v (a_v)^2$ .

El problema de cálculo del primer eje principal se reduce a

$$\begin{aligned} \text{Minimizar} \quad & \sum_v \mu_v (a_v)^2 \\ \text{s.a} \quad & \sum_v (a_v)^2 = 1 \end{aligned}$$

Cuya solución es  $a_1 = 1$ ,  $a_v = 0 \forall v \neq 1$ . Esto muestra que el primer eje principal es  $\Delta u_1$  y  $I \Delta' u_1 = \mu_1$ .

El segundo eje principal corresponde al eje de mínima inercia, ortogonal a  $\Delta u_1$ . Se puede probar, con un cálculo muy similar al anterior, que dicho eje es  $\Delta u_2$  y que  $I \Delta' u_2 = \mu_2$ .



Los primeros dos ejes principales determinan el plano principal  $W_2 = \Delta u_1 \otimes \Delta u_2$ , sev de mínima inercia:  $Iw'_2 = \mu_1 + \mu_2$ .

Continuando con el proceso descrito llegamos los  $k$  primeros ejes principales  $\Delta u_1, \Delta u_2, \dots, \Delta u_k$ , los que engendran el sev principal  $W_k = \Delta u_1 \otimes \dots \otimes \Delta u_k$  y  $Iw'_k = \mu_1 + \dots + \mu_k$ .

### 3.4.- CALIDAD DE LA REPRESENTACION.-

Lo primero es saber que entendemos por representación. Recordemos que estamos tras una presentación simple de los datos y para ello vamos a proyectar la nube  $N(i)$  en su sev de dimensión baja, digamos  $k$ , perdiendo un mínimo de información.

Todo el desarrollo teórico anterior nos lleva a escoger  $W_k$  para proyectar. La pregunta es ¿qué tan buena es la representación? y la respuesta dependerá del valor de un índice que definiremos más abajo.

La información que contienen los datos podemos imaginarla asociada a la variabilidad de los mismos.

Si proyectamos  $N(i)$  al centro de gravedad perdemos toda variabilidad, luego es natural interpretar el valor  $Ig$  como una medida de la información total de la muestra. Por otra parte, puesto que  $Iw_k$  mide la distorsión debida a la proyección sobre  $W_k$ , debemos interpretarla como pérdida de información al proyectar. Alternativamente,  $Iw'_k = Ig - Iw_k$  puede verse como la parte de información que explica el sev  $W_k$ .

El índice porcentaje de inercia explicada por el sev principal  $W_k$  se define como

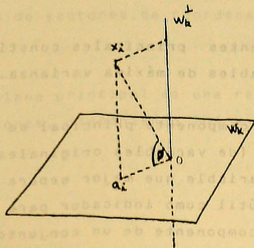
$$i_k = Iw'_k / I g$$

pero  $Iw'_k = \mu_1 + \dots + \mu_k$ ,  $I g = I E' = I \Delta' u_1 + \dots + I \Delta' u_p = \mu_1 + \dots + \mu_p = \text{traza}(VM)$ . Luego

$$i_k = \frac{\sum_{i=1}^k \mu_i}{\text{traza}(VM)}$$

Si  $i_k = 1$  la representación en  $W_k$  es perfecta y  $\mu_v = 0$  para  $v = k+1, k+2, \dots, p$ . En tal caso  $\text{rango}(VM) = k = \text{rango}(X)$  lo que significa (suponiendo  $n > p$ ) que hay multicolinealidad (redundancia lineal) en las variables.

El porcentaje de inercia es un índice global de calidad. Es posible y a veces recomendable recurrir a índices que muestran si uno o varios individuos están bien representados por su proyección en  $W_k$ .



$$\text{Por ejemplo } \cos^2(\theta) = (M(x_i, u_1))^2 + \dots + (M(x_i, u_k))^2 / \|x_i\|^2$$

## 3.5.- LAS COMPONENTES PRINCIPALES.-

Asociemos a los ejes principales  $\Delta U_v$  en E, variables  $C^j$  en F, llamadas componentes principales, a través de las aplicaciones M y  $X'$  del esquema de dualidad:

$$C^v = X' MU_v \quad v = 1, 2, \dots, p$$

Notemos sus propiedades más importantes:

1) Son centradas, es decir de media nula  $me(C^v) = Dp(C^v, 1) =$

$$\langle Dp1, X' MU_v \rangle = \langle XDp1, MU_v \rangle = \langle 0, MU_v \rangle = 0.$$

2) Son no correlacionadas:  $COV(C^v, C^s) = Dp(C^v, C^s) =$

$$\langle DpX' MU_v, X' MU_s \rangle = \langle XDpX' MU_v, MU_s \rangle =$$

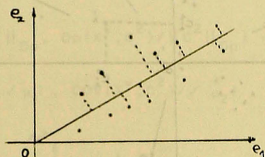
$$\langle VMu_v, MU_s \rangle = \mu_v M(U_s, U_v) = 0$$

3)  $var(C^v) = Dp(C^v) = \langle VMu_v, MU_v \rangle = \mu_v M(U_v, U_v) = \mu_v$ .

Las componentes principales constituyen un sistema Dp-ortogonal de variables de máxima varianza. Una base para F.

La primera componente principal se interpreta como la combinación lineal (de variables originales) que tiene varianza máxima. Es la variable que mejor separa a los individuos y en ese sentido es útil como indicador para ordenarlos. Por ejemplo la primera componente de un conjunto de variables económicas, culturales, demográficas, etc., serviría para definir un

indicador de desarrollo de países o regiones.



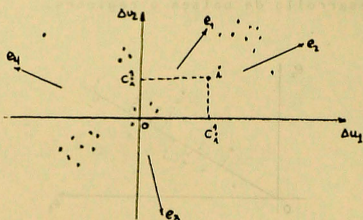
### 3.6.- SALIDAS GRAFICAS.-

Gran parte del interés del ACP está en los gráficos de los resultados, que permiten visualizar relaciones entre variables así como la organización de los individuos.

En aplicaciones prácticas del método, se seleccionan los primeros ejes tales que el porcentaje de inercia explicado por ellos sea al menos 90%, pero habitualmente nunca se analizan más de 4 ó 5.

La proyección de la nube  $N^{(1)}$  sobre el plano principal produce una colección de vectores de coordenadas  $M(x_i, U_1)$ ,  $M(x_i, U_2)$   $i = 1, 2, \dots, n$ .

El gráfico del plano principal es una representación cartesiana de las proyecciones, disponiendo los ejes  $\Delta U_1$  y  $\Delta U_2$  perpendicularmente.



El gráfico puede enriquecerse incorporando por ejemplo identificadores de individuos y los porcentajes de inercia de los ejes.

Al clasificar visualmente a los individuos es interesante dar sentido a los grupos en términos de las variables antiguas. Si proyectamos los vectores de la base canónica de  $E\{e_1, e_2, \dots, e_p\}$  (que representan a las viejas variables), orientaremos el plano y podremos explicar las clasificaciones y proponer interpretación de las componentes principales.

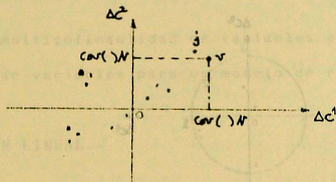
Para completar el análisis puede ser interesante mirar otros planos como  $\Delta U_1 \otimes \Delta U_3$ ,  $\Delta U_2 \otimes \Delta U_3$  etc.

### 3.7.- REPRESENTACION DE VARIABLES.-

Usando la base de  $F$ , de las componentes principales normalizadas  $\{C^1/\sqrt{\mu_1}, \dots, C^p/\sqrt{\mu_p}\}$  representamos las variables en un gráfico cartesiano, cuya calidad depende del índice de inercia. Si éste es alto, los ángulos y distancias entre las proyecciones reflejarán los ángulos y distancias entre las variables originales.

Escogiendo las dos primeras componentes, la variable  $X^V$  es tará representada por el punto de coordenadas

$$\left( D_p(X^V, C^1) / \|C^1\|_{D_p}, D_p(X^V, C^2) / \|C^2\|_{D_p} \right) = \\ \left( \text{COV}(X^V, C^1) / \sqrt{\mu_1}, \text{COV}(X^V, C^2) / \sqrt{\mu_2} \right)$$



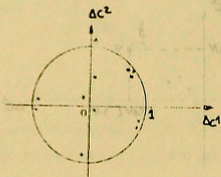
### 3.8.- LA ELECCION DE LAS METRICAS.-

Dos son las métricas  $M$  que se usan comúnmente en el espacio  $E$ . La primera de ellas es la identidad o métrica euclidiana usual. Según esta métrica todas las direcciones tienen igual importancia en el cálculo de la distancia y esto hace que los resultados del ACP dependan de las unidades de medición de las variables.

La otra alternativa es una distancia que otorga pesos a las direcciones, en razón inversa a la dispersión de las variables que representa. Esta métrica se apoya en el principio según el cual se otorga menor peso a las medidas hechas con instrumentos de poca precisión. Su matriz con respecto a la base canónica es diagonal con los recíprocos de las varianzas de  $x^1, x^2, \dots, x^p$ .

Se puede demostrar que la elección de esta métrica es equivalente a usar la identidad sobre los datos estandarizados (variables centradas y de varianza 1).

En estas condiciones la matriz  $V = XDpX'$  de covarianzas se convierte en matriz de correlaciones y el gráfico de las variables se transforma en el círculo de correlaciones,



Justamente porque ahora las coordenadas de  $X^V$  son  $\text{cor}(X^V, C^1)$  y  $\text{cor}(X^V, C^2)$ . Debido a que  $\|X^V\| = 1 = \text{var}(X^V)$  resulta

$$1 = \left\| \sum_S (Dp(X^V, C^S) / \sqrt{\mu_S}) C^S / \sqrt{\mu_S} \right\|_{Dp}^2 = \sum_S (Dp(X^V, C^S) / \sqrt{\mu_S})^2,$$

luego los puntos del gráfico de variables caen en el interior (o el borde) del círculo unitario. La cercanía al borde indica mejor calidad de la representación.

### 3.9.- APLICACIONES DEL ACP.-

El ACP puede ser empleado provechosamente para responder a problemas como los siguientes:

- 1) Clasificar individuos descritos por variables cuantitativas.
- 2) Detectar observaciones aberrantes o anómalas.

- 3) Detectar factores o variables "escondidas"
- 4) Economizar espacio de almacenamiento de datos.
- 5) Construir indicadores (desarrollo, inteligencia, etc.)
- 6) Representar gráficamente una estructura de correlaciones.
- 7) Hacer una representación euclidiana de una tabla de distancias.
- 8) Detectar multicolinealidad de variables en un modelo lineal.
- 9) Seleccionar variables para un modelo de regresión lineal.

#### 4.- REGRESION LINEAL.-

Entra en escena una nueva variable que vamos a representar como el vector  $y$  en el espacio  $F$ .

El problema que abordaremos en esta sección del curso consiste en explicar el comportamiento de  $y$  en términos de las variables  $x^1, x^2, \dots, x^p$ .

El enunciado del problema es particularmente vago porque un estudio de regresión puede tener muy variados objetivos:

- 1) Describir la dependencia entre  $x^1, x^2, \dots, x^p$  e  $y$ .
- 2) Estudiar el comportamiento de  $y$  retirando los efectos de con fusión de otras variables.
- 3) Determinar una ley empírica que relacione  $y$  con las otras va riables.
- 4) Ajustar un modelo para suavisar datos.



- 5) Contribuir con evidencia empírica a un análisis causal.
- 6) Predecir el comportamiento de  $y$  cuando se conocen los valores de  $X^1, \dots, X^P$  para un nuevo individuo.
- 7) Controlar un sistema complejo.

La perspectiva de la regresión para responder a esta serie de tareas distintas es la construcción de un modelo que ligue  $y$  con  $X^1, \dots, X^P$ , a través de una función  $f$  sencilla (interpretable si es posible) tal que para cada  $i$  en  $i$  el valor  $y_i$  se determine aproximadamente como  $f(X_i^1, \dots, X_i^P)$ . La diferencia se atribuye a un error aditivo  $e_i$  desconocido.

Tenemos entonces:

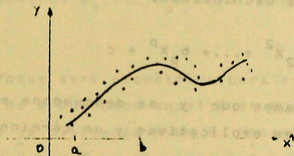
$$y_i = f(X_i^1, X_i^2, \dots, X_i^P) + e_i \quad i = 1, 2, \dots, n$$

Para escoger adecuadamente la función  $f$  es necesario definir un criterio, una familia de candidatos (paramétrica o no) y dejar que los datos decidan.

En el marco de la regresión lineal la función  $f$  se toma lineal en los argumentos, es decir, los candidatos para modelo tienen la forma  $f(X_1, \dots, X_p) = \sum_1 b_1 x_1$  donde  $b_1, b_2, \dots, b_p$  son números reales llamados parámetros del modelo

$$y_i = b_1 X_i^1 + b_2 X_i^2 + \dots + b_p X_i^P \quad i = 1, 2, \dots, n \quad (1)$$

La razón para trabajar con modelos aparentemente tan simples como los lineales se debe a que por una parte la aproximación lineal puede ser suficiente en cierto rango de las variables  $X^1, \dots, X^p$ , lo que no justifica una función más compleja. Por otra parte, el cálculo numérico y el análisis teórico de los modelos lineales es harto más simple que las alternativas no lineales.



En la figura anterior se presenta una forma de dependencia entre  $y$  y  $x^1$  que sería perfectamente aproximada por una recta en el rango  $[a, b]$ .

La variable  $y$  y las restantes  $X^1, \dots, X^p$ , tienen roles asimétricos y suelen llamarse respectivamente variable a explicar y variables explicativas.

Tradicionalmente el estudio del modelo de RL va precedido de una serie hipótesis probabilísticas sobre las variables y los errores. En una primera etapa vamos a prescindir de ellas para enfrentar el problema con espíritu descriptivo.

## 4.1.- NOTACION MATRICIAL-VECTORIAL.-

Presentaremos el modelo de RL usando los mismos espacios que intervienen en el ACP (ver esquema de dualidad).

Definamos el vector de errores  $e$  en  $F$  cuyas coordenadas en la base canónica son los  $e_i$ . Entonces las ecuaciones en (1) la fundimos en una escribiendo:

$$y = b_1 X^1 + b_2 X^2 + \dots + b_p X^p + c \quad (2)$$

Donde apreciamos que  $y$  se descompone en una combinación lineal de variables explicativas y un término de error  $e$ .

Sea  $w$  el sev de  $F$  engendrado por las variables explicativas  $X^1, \dots, X^p$ . La ecuación (2) nos dice que  $y$  se explica mediante una componente a lo largo de  $w$  más un error.

En la ecuación (1) podemos interpretar al vector de coeficientes  $B$  como una forma lineal en  $E^*$  porque asocia al individuo  $X_i$  en  $E$  el valor  $\sum_v b_v X_i^v = B(X_i) = \langle B, X_i \rangle$ .

En este sentido  $B$  se comporta como el representante en  $E^*$  de una variable. Para recuperar la auténtica variable de  $F$  que  $B$  representa recurrimos a la aplicación  $X'$  y obtenemos  $X'B$ . Pero si  $B = \sum_v b_v e_v^*$  tenemos  $X'B = \sum_v b_v X'e_v^* = \sum_v b_v X^v$ .

De manera que la ecuación (2) es equivalente a la expresión matricial

$$y = X'B + e$$

Es bueno notar que  $W = X'(E^*)$  y que  $\dim(W) = \text{rango}(X') =$  número de variables linealmente independientes.

**PROBLEMA CENTRAL: ENCONTRAR B EN  $E^*$ , OPTIMO.**

El problema planteado se conoce como problema de estimación de los parámetros  $b_1, \dots, b_p$ . Para resolverlo necesitamos un criterio.

Nuestro enfoque será geométrico para explicar o aproximar y con un vector de  $W$  buscamos aquel  $\hat{y}$  en  $W$  más cercano de  $a$  y en el sentido de la distancia  $N$  de  $F$ . Necesariamente  $\hat{y}$  se describe como  $\hat{y} = X \cdot \hat{B}$  con  $\hat{B}$  en  $E^*$ , aunque  $\hat{B}$  puede no ser único. Consideramos a  $\hat{B}$  como solución del problema central.

La unicidad de  $\hat{B}$  se consigue si  $X'$  es inyectiva en efecto si  $X'B_1 = X'B_2$  tendríamos  $\ker(X') \neq \{0\}$

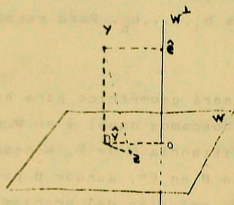
Si  $X'$  no es inyectiva se dice que el parámetro no es identificable, problema que está ligado a la estimabilidad (insesgada) de combinaciones lineales de los  $b_v$ .

Desde el punto de vista de las variables, cuando  $X'$  no es inyectiva se habla de multicolinealidad.

El enfoque geométrico que hemos adoptado corresponde al famoso criterio de mínimos cuadrados. Si  $N = 1$  se habla de mínimos cuadrados ordinarios. Si  $N$  es cualquier métrica euclidiana se le llama mínimos cuadrados generalizados.

## 4.2.- SOLUCION DEL PROBLEMA CENTRAL.-

Nos basamos en el siguiente resultado muy conocido, en el punto  $\hat{y}$  de  $W$  más cercano de  $y$  en el sentido de  $N$  es la proyección  $N$ -ortogonal de  $y$  sobre  $W$ . Para convencernos miremos la figura y los cálculos a continuación



Designemos por  $A$  el proyector  $N$ -ortogonal sobre  $W$ . Entonces  $y = \hat{y} + \hat{e} = A(y) + (1-A)(y)$ .

Si  $Z$  en  $W$  es cualquiera tenemos:  $d(y, x)^2 = \|y - Z\|^2 = \|y - \hat{y} + \hat{y} - Z\|^2 = \|y - \hat{y}\|^2 + \|\hat{y} - Z\|^2 \geq \|y - \hat{y}\|^2 = d(y, \hat{y})^2$ .

Busquemos ahora la forma explícita para  $\hat{y}$  (o para el proyector  $A$ ), partiendo del simple hecho que  $\hat{e}$  (vector de residuos) es ortogonal a  $W = X'(E^*)$ .

$$\begin{aligned} \text{Para cualquier } e^*_v, \hat{e} &= \langle NX'e^*_v, \hat{e} \rangle = \langle e^*_v, XN\hat{e} \rangle \\ &= \langle e^*_v, XN(y - \hat{y}) \rangle = 0 \end{aligned}$$

de donde salen las famosas ecuaciones normales:

$$XNy = XN\hat{y} = XNX'\hat{B} = V\hat{B}$$

Si  $X'$  es inyectiva  $V$  es invertible y podemos despejar la solución  $\hat{B}$  única:

$$\hat{B} = V^{-1}XNy = (XNX')^{-1}XNy.$$

Cuando  $X'$  no es inyectiva se puede recurrir a inversas generalizadas para seleccionar alguna solución de las ecuaciones normales.

La expresión para el proyector en el modelo de rango completo ( $X'$  inyectiva) es

$$A = X'(XNX')^{-1}XN.$$

#### 4.3.- CALIDAD DE LA REGRESION.-

La regresión será buena en la medida que la parte de  $y$  que no se logra explicar sea pequeña. En otras palabras, mejor es la calidad mientras más chico sea  $\hat{\epsilon}$ .

El índice  $r^2$  de correlación lineal múltiple mide la calidad global de la regresión en los términos expresados más arriba:

$$r^2 = \|\hat{y}\|^2 / \|y\|^2$$

y no es otra cosa que el  $\cos^2$  del ángulo entre  $y$  y el sev

w.

En general la calidad de la regresión depende del comportamiento global de los residuos  $\hat{\epsilon}$ , cuyo análisis detallado sirve para comprender y corregir el modelo.

#### 4.4.- HIPOTESIS PROBABILISTICAS.-

El modelo lineal  $y = X'B + e$  suele adornarse con muchas hipótesis probabilísticas que permiten, entre otras cosas, demostrar la optimalidad del criterio de mínimos cuadrados y efectuar test de hipótesis estadísticas sobre los parámetros.

Estas hipótesis facilitan el tratamiento teórico elegante del modelo de RL pero la validez de las mismas es difícilmente verificable a partir de los datos.

Consideremos algunas de ellas:

- H1)  $y$  es un vector aleatorio.
- H2) Las variables explicativas no son aleatorias
- H3) El vector de errores tiene esperanza nula.  $E(e_i) = 0$ .
- H4) Los errores son no correlacionados.  $E(e_i e_k) = 0$ .
- H5) Los errores son independientes.
- H6) Las varianzas de los  $e_i$  son iguales.
- H7)  $e$  tiene distribución normal multivariada.

#### 4.5.- TEOREMA DE GAUSS MARKOV.-

Con H1, H2 y H3 se puede demostrar un interesante teorema sobre la calidad del estimador de mínimos cuadrados  $\hat{B}$  y además indicar la métrica  $N$  que se debe escoger.

En la fórmula  $\hat{B} = (XNX')^{-1}XNy$  observamos que  $\hat{B}$  es función lineal de  $y$ .

El teorema de Gauss Markov afirma que entre todos los estimadores  $\bar{B}$  lineales en  $y$ ,  $\hat{B}$  es mejor, en el sentido que la diferencia de las matrices de varianzas-covarianza  $V(\bar{B}) - V(\hat{B})$  es semi dp, siempre que la métrica  $N$  se tome como la inversa de la matriz de varianzas covarianzas del vector de errores  $E(ee') = \Sigma$ .

Si agregamos la hipótesis  $H7$  se puede reforzar la conclusión del teorema de Gauss Markov, diciendo que  $\hat{B}$  es el mejor entre todos los estimadores insesgados.

Finalmente corresponde notar que la hipótesis de normalidad permite, al menor en principio conocer la distribución de  $\hat{y}$ ,  $\hat{B}$ ,  $\hat{e}$  y de otros estadísticos de interés, brindando al usuario la teoría de la inferencia estadística clásica; estimación y test de hipótesis.

#### 4.6.- INTERPRETACION Y ESTABILIDAD DE B.

El modelo lineal provoca en el analista la legítima tentación de interpretar los coeficientes. Imaginarlos por ejemplo como elasticidad precio-demanda, índice de reactividad, coeficiente de aislación, etc. Esta práctica es interesante pero peligrosa puesto que el significado de un coeficiente depende no sólo de la variable que multiplica sino también de las restantes (recordar que en general las variables explicativas están correlacionadas).

Por otra parte el coeficiente que multiplica a una variable dada cambia si el conjunto de variables "acompañantes" es alterado.



Otro problema delicado es la estabilidad de  $\hat{B}$ , que podemos formular como una pregunta: ¿qué confianza se puede tener en el valor de  $\hat{B}$  calculado. La respuesta depende, por un lado de la precisión de los algoritmos de cálculo, (factor no despreciable) y por otro de las relaciones de colinealidad entre las variables explicativas.

Para enfrentar situaciones concretas se procede a estimar la varianza del estimador del coeficiente de interés y luego se usa ese valor como indicador de calidad.

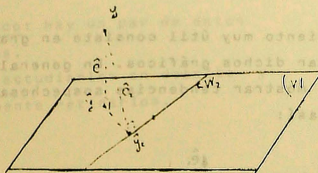
Hay un último aspecto importante que mencionar la significación de  $\hat{B}$ . Esto se refiere a decidir si uno o varios betas son significativamente no nulos, o si por el contrario, son distintos de cero solo por azar.

En términos del modelo, la significación de los betas habla de la relevancia que las variables correspondientes tienen en el modelo. Si demostramos (o decidimos) que  $b_v$  debería ser cero, es tamos diciendo que  $X^v$  debe salir dle modelo, porque en nada contribuye.

Este tipo de problema cae en el ámbito del análisis de varianza y los test F que no se tratan en este apunte. Sin embargo hay una interesante aproximación geométrica al problema de la significación que vale la pena considerar:

Supongamos que  $B$  se particiona en dos pedazos  $B_1$  y  $B_2$ . Sos pechamos que las variables asociadas a  $B_1$  no son relevantes, es decir, creemos que  $B_1$  debería ser cero.

Para comprobar esto procedemos como sigue: estimar el mo delo completo, luego estimar el modelo botando las variables sospechosas y finalmente comparar. ¿Cómo?, decidiendo en base a los tamaños de los respectivos vectores de residuos.



$W$  = Espacio engendrado por todas las variables explicativas.

$W_2$  = Espacio de las variables no sospechosas

$\hat{y}$  = Proyección N-ortogonal en  $W$ .

$\hat{y}_2$  = Proyección N-ortogonal en  $W_2$

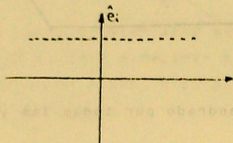
$\hat{e}$  =  $y - \hat{y}$ ,  $\hat{e}_2 = y - \hat{y}_2$

La regla de decisión es botar las variables sospechosas si  $\|\hat{e}_2\|$  no es significativamente más grande que  $\|\hat{e}\|$

## 4.7.- LOS RESIDUOS.-

El vector de los residuos  $\hat{e}$  contiene valiosa información sobre datos aberrantes, errores en la especificación del modelo, validez de los supuestos probabilísticos, etc.

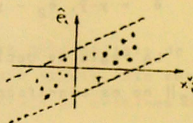
Un procedimiento muy útil consiste en graficar los residuos e interpretar dichos gráficos. En general un gráfico de residuos no debe mostrar tendencias sospechosas. Debe verse aproximadamente así:



Situaciones que se apartan de este caso ideal, admiten algunas explicaciones como veremos en las figuras a continuación:

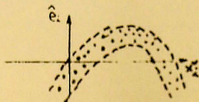
Diagnóstico: falta término lineal en la variable  $X^V$ .

Remedio: agregarlo



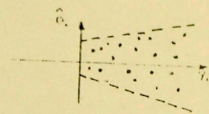
Diagnóstico: falta término cuadrático en  $X^V$ .

Remedio: agregarlo



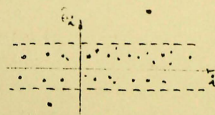
Diagnóstico: la igualdad de varianzas de los  $y_i$  es falsa.

Remedio: Cambiar métrica N o transformar la variable  $y$ .



Diagnóstico: hay un par de datos aberrantes.

Remedio: estudiarlos en detalle y eventualmente retirarlos.



Apenas hemos tocado la superficie de la regresión lineal, en teoría o práctica. A quienes deseen profundizar en el tema recomiendo el excelente clásico "Applied Regression Analysis" de N. Draper y H. Smith (2° edición).

## 5.- REFERENCIAS.-

- [1] Draper N. Smith H., Applied Regression Analysis, Sec. ed. Wiley. 1979.
- [2] Mosteller F. Tuckey J., Data Analysis and Regression, Ad. Wesley, 1977.
- [3] Pages J.P. Cailliez F., Introducción a L'analyse des Donnees, Smash, 1976.

## 6.- ABREVIATURAS.-

ACP = Análisis en Componentes Principales  
 dp = definida positiva  
 etc = etcétera  
 RL = Regresión lineal  
 sev = subespacio vectorial  
 sii = si y sólo si