

FRANCISCO MONTES SUAY

ESTADÍSTICA MULTIVARIANTE (EXPLICADA)

UNIVERSITAT DE VALÈNCIA

Índice general

1	Álgebra matricial	5
2	Muestras de una normal multivariante	15
3	Coeeficientes de correlación muestrales	45
4	El estadístico T^2 generalizado	71
5	Clasificación de observaciones	89
6	La distribución de la matriz de covarianzas muestral y de la varianza muestral generalizada	109

Capítulo 1

Álgebra matricial

Se pretende en este primer capítulo recordar algunos conceptos referentes al cálculo matricial y materias con él relacionadas, por ser este útil matemático de gran importancia para facilitar los cálculos que el manejo de muestras multivariantes implican.

Evidentemente vamos a insistir en algunos conceptos de mucha importancia, a saber, definiciones referentes a los distintos tipos de matrices, operaciones elementales con vectores y matrices, como suma, producto, producto por un escalar, etc. así como la correspondencia biunívoca existente entre las matrices y las formas lineales, el producto interno vectorial, la transposición de matrices, la definición y obtención del determinante de una matriz cuadrada, etc.

Ocupémonos ahora de otros conceptos, videtur de algunos de ellos han sido debidamente propiedades de interés.

1. INVERSA DE UNA MATRIZ

Sólo existe en para las matrices cuadradas ~~no~~ regulares. Viene definida por $A \cdot A^{-1} = A^{-1} \cdot A = I$.

Se impone que su existencia está ligada al valor del determinante de A , por cuanto si este es 0, aquella no existe. Propiedades de interés:

- la inversa de una matriz simétrica también lo es
- la inversa de la traspuesta es la traspuesta de la inversa.
- la inversa del producto invierte

$$(ABC)^{-1} = C^{-1} \cdot B^{-1} \cdot A^{-1}$$

- la inversa de una matriz diagonal es una matriz diagonal cuyos elementos son los inversos de los elementos primitivos.

2. RANGO DE UNA MATRIZ

Sea A una matriz $K \times P$, con $K \leq P$. El rango de A es el número de rectas fila linealmente independientes en la matriz. Si $P \leq K$, la definición es análoga substituyendo fila por columna. Se demuestra, que en cualquier caso el rango de una matriz es único, independientemente de si obtenidos considerando las rectas fila o columna.

Atendiendo al valor del determinante de los menores complementarios en una matriz A , podemos afirmar que:

la matriz A es de rango r si contiene al menos un menor $r \times r$ distinto de cero, siendo todos los menores en dimensiones mayores que r .

Algunas propiedades de interés:

- El rango de A' y el de A son iguales.
- El rango de $A'A$ es igual al de A e igual al de AA' .
- El rango de A no varía si se pre o post-multiplica A por una matriz no singular.

Como consecuencia de c) podemos afirmar que el rango de una matriz A no cambia si

- Se intercambian las filas (columnas) cualesquiera
- multiplicando cada fila (columna) por un escalar
- añadiendo a una fila (columna) los elementos de otra únicamente multiplicados por un escalar.

y todo ello por que las operaciones dadas en estas tres propiedades pueden llevarse a cabo mediante premultiplicación (caso de las filas) o postmultiplicación (caso de columnas) de la matriz A por una matriz no singular.

$$A = \begin{bmatrix} 5 & 1 \\ -2 & 3 \\ 3 & 2 \end{bmatrix}$$

$$E = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$EA = \begin{bmatrix} -2 & 3 \\ 5 & 1 \\ 3 & 2 \end{bmatrix}$$

Operaciones sobre las filas y columnas nos permiten transformar una matriz A de rango r , en una matriz canónica F , que tiene la forma

$$F = \begin{bmatrix} 1 & 0 & \dots & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & \dots & 0 \end{bmatrix}$$

es decir, 1's en las r primeras posiciones de la diagonal principal y ceros en el resto. Si designamos por $R_s \rightarrow R_p$ y $C_s \rightarrow C_q$ estos cambios en filas y columnas podemos escribir

$$R_p \rightarrow R_s A C_s \rightarrow C_q = P A Q = F$$

Si A es cuadrada y de rango n (máximo), entonces

$$R_q \rightarrow R_s A = I$$

y de aquí

$$A^{-1} = R_p \rightarrow R_s I$$

Este hecho es aprovechado para invertir matrices de forma rápida y sencilla, útil en los casos de matrices de elevada dimensión.

3. INVERSA GENERALIZADA

La definición dada anteriormente es sólo aplicable a matrices cuadradas no singulares. Se trata de introducir un concepto de inversa que generalice el anterior.

La G -inversa de una matriz A de cualquier dimensión viene dada por:

$$(1) \quad A G A = A$$

Si la dimensión de A es $p \times q$ y su rango r , la dimensión de G es $q \times p$ y su rango $s \leq r$.

Existen otras definiciones, por ejemplo la debida a Penrose que exige además

$$G A G = G \quad (G A)' = G A \quad (A G)' = A G$$

La obtención de la G definida en (1) puede hacerse a partir de la reducción canónica de A , como sigue:

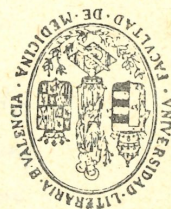
$$P A Q = F = \begin{bmatrix} D & 0 \\ 0 & 0 \end{bmatrix}$$

donde D es una matriz diagonal $r \times r$ (no necesariamente I) y el resto matrices nulas de dimensión adecuada. P y Q las anteriormente definidas a partir de operaciones en filas y columnas. Definimos ahora

$$F^{-1} = \begin{bmatrix} D^{-1} & 0 \\ 0 & 0 \end{bmatrix}$$

entonces $G = Q F^{-1} P$, es una inversa generalizada de A , que desde luego no es única. Si A es cuadrada y no singular, entonces $G = A^{-1}$ y es única.

EXAMENES



4. SISTEMAS DE ECUACIONES LINEALES

Recordemos que un sistema de la forma

$$a_{11}x_1 + \dots + a_{1n}x_n = c_1$$

$$a_{m1}x_1 + \dots + a_{mn}x_n = c_m$$

admite una representación matricial mediante

$$AX = C$$

con

$$A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix}$$

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$$C = \begin{bmatrix} c_1 \\ \vdots \\ c_m \end{bmatrix}$$

si $C = [0]$, el sistema se llama homogéneo

(1) El sistema posee una solución si y solo si la matriz ampliada $[A|C]$ de dimensión $m \times (n+1)$ tiene el mismo rango que la matriz A . (Se dice entonces que el sistema es consistente)

Estudiaremos los tipos de sistemas y sus soluciones en función del rango de A y $[A|C]$.

a) Sistemas no homogéneos: A cuadrada y no singular

la solución es única y viene dada por

$$X = A^{-1}C$$

b) Sistemas no homogéneos: A $m \times n$ y de rango r

Si se satisface (1) existe una solución que puede obtenerse hallando r ecuaciones linealmente independientes, resolviendo para r incógnitas en función de las constantes y de las restantes $n-r$ incógnitas.

c) Sistemas homogéneos

Entonces el rango de $[A|0]$ es siempre el mismo que el de A por tanto siempre tienen solución.

c.1) Si $\text{rango}(A) = n =$ número de incógnitas, solución única dada por $X = [0]$

c.2) En cualquier otro caso solución del tipo expresada en b).

En cualquier caso la solución puede expresarse en función de la misma generalización mediante la siguiente expresión debida a Rao (1962), siendo el sistema consistente.

$$X^* = GC + (GA - I)Z$$

donde Z es un vector $n \times 1$ de constantes arbitrarias, obsérvese que si existe A^{-1} , entonces $X^* = A^{-1}C$, como ya sabemos.

5. MATRICES ORTOGONALES

Un vector se dice que es ortogonal respecto de otro si su producto interior es nulo.

Una matriz T se dice ortogonal si sus columnas (o filas) son vectores ortogonales, siendo cada una de ellas una vector de norma unidad. De esta definición se desprende que:

$$TT' = T'T = I$$

algunas propiedades de interés:

- Las columnas de una matriz ortogonal son ortogonales.
- El determinante de una matriz ortogonal es siempre 1 o -1.
En efecto $|TT'| = |I| = 1$ y $|T| = |T'| = \pm 1$
- El producto de matrices ortogonales de la misma dimensión es también una matriz ortogonal

5. FORMAS CUADRÁTICAS

Una forma cuadrática en las variables x_1, \dots, x_n es una expresión del tipo

$$f(x_1, \dots, x_n) = a_{11}x_1^2 + \dots + a_{nn}x_n^2 + 2a_{12}x_1x_2 + \dots + 2a_{1n}x_1x_n + \dots + 2a_{n-1,n}x_{n-1}x_n =$$

$$= \sum_i \sum_j a_{ij}x_i x_j \quad \text{con } a_{ij} = a_{ji} \quad \text{pudiendo ser uno algún } a_{ij}$$

De inmediato se observa que una forma cuadrática admite notación matricial, como sigue:

$$f(x_1, \dots, x_n) = x' A x$$

donde A es una matriz simétrica $n \times n$. Las formas cuadráticas juegan un importante papel en estadística, tanto univariante como multivariante. Por ejemplo, la suma de los cuadrados de las desviaciones respecto de la media muestral

$$\sum_{i=1}^N (x_i - \bar{x})^2 = \sum_{i=1}^N x_i^2 - \frac{1}{N} \left(\sum_{i=1}^N x_i \right)^2$$

es una forma cuadrática cuya matriz de coeficientes viene dada por:

$$A = \begin{bmatrix} \frac{N-1}{N} & -\frac{1}{N} & \dots & -\frac{1}{N} \\ -\frac{1}{N} & \frac{N-1}{N} & \dots & -\frac{1}{N} \\ \vdots & \vdots & \ddots & \vdots \\ -\frac{1}{N} & -\frac{1}{N} & \dots & \frac{N-1}{N} \end{bmatrix}$$

Carácter de una forma cuadrática

- Definida positiva, semidefinida positiva.

Una matriz simétrica A y su forma cuadrática asociada se dice que es definida positiva si $x' A x > 0$, $\forall x \neq 0$.

Si $x' A x \geq 0$, $\forall x \neq 0$ entonces se la denomina semidefinida positiva.

- Definida negativa, semidefinida negativa.

Id. con los signos cambiados.

- Indefinida

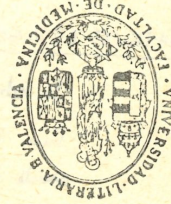
Se puede tomar cualquier valor positivo, negativo o nulo.

Algunas propiedades de las formas cuadráticas definidas y semidefinidas positivas

- Las formas cuadráticas definidas positivas tienen matrices de rango completo. Es posible mediante complementarios de cuadrados menores reducir dichas formas a expresiones del tipo

$$d_1 y_1^2 + \dots + d_n y_n^2 \quad d_i > 0, \forall i$$

EXAMENES



Algo similar ocurre con las semidefinidas positivas, pero ahora la expresión es

1 (3)

$$dx_1^2 + \dots + dx_r^2 \quad dx_i > 0 \quad r \leq n \text{ es el rango de la ~~forma cuadrática~~ matriz.}$$

2) No obstante para determinar el carácter de una forma basta a dar ahora una condición necesaria y suficiente. Formaremos en primer lugar la cuestión de los determinantes de los menores principales, a saber

$$p_0 = 1 \quad p_1 = a_{11} \quad p_2 = \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} \quad \dots \quad p_i = \begin{vmatrix} a_{11} & \dots & a_{1i} \\ \vdots & & \vdots \\ a_{ii} & \dots & a_{ii} \end{vmatrix} \quad \dots \quad p_n = |A|$$

Si el rango de A es r, deduce que A es regular si $p_r \neq 0$ y no hay dos p_i consecutivos en la sucesión iguales a cero. Siempre es posible, intercambiando filas y columnas, colocar una matriz simétrica en su forma regular. Entonces, si A es una matriz dispuesta regularmente, tenemos:

2.a) Una condición ^{he.} suficiente para la definición positiva es que $p_i > 0, i = 1, \dots, n$

2.b) Id. para la semidefinición positiva es que $p_1 > 0, \dots, p_r > 0$ y los restantes $n-r$ p_i sean nulos, donde $r \leq n$.

Se puede demostrar aplicando esto que la forma cuadrática de la suma de los cuadrados de las derivadas respecto de la media es semidefinida positiva con rango $n-1$ para un matriz. Los distintos p_i tienen un valor

$$p_i = \frac{1}{n} (n-1-i+1).$$

6. RAICES Y VECTORES CARACTERÍSTICOS DE UNA MATRIZ

Las raíces características de una matriz A, $p \times p$, son las soluciones de la ecuación determinante

$$|A - \lambda I| = 0$$

El determinante es un polinomio de grado p en λ y tiene por tanto p raíces. Aplicando el desarrollo de Laplace ~~de~~ determinante característico, podemos escribir

$$|A - \lambda I| = (-\lambda)^p + S_1 (-\lambda)^{p-1} + \dots + S_{p-1} (-\lambda) + |A|$$

donde S_i es ~~la~~ la suma de todos los determinantes de los menores principales de dimensión i . Si sea en particular la suma de los elementos de la diagonal principal de A, o sea $\text{tr}(A)$. De la teoría de las soluciones de una ecuación polinomial, tenemos de inmediato que:

- 1) El producto de las raíces características de A es igual a $|A|$
- 2) La suma de las raíces características de A es igual a $\text{tr}(A)$

Algunas propiedades de interés y de posterior uso de las raíces características son:

- a) Las raíces características de una matriz simétrica con elementos reales son todas reales.
- b) Las raíces características de una matriz definida positiva son todas positivas.
- c) Para una matriz $n \times n$ semidefinida positiva de rango $r \leq n$, hay exactamente r raíces características positivas y $n-r$ nulas.
- d) Las raíces características no nulas del producto AB son iguales a las raíces no nulas de BA. Por tanto $\text{tr}(AB) = \text{tr}(BA)$
- e) Las raíces características de una matriz diagonal son los elementos de la diagonal.

A cada raíz característica podemos asociar un vector característico, cuyos elementos satisfacen la ecuación homogénea

$$[A - \lambda_i I] x_i = 0$$

Como el determinante $|A - \lambda_i I| = 0$ existe siempre una solución no trivial para x_i que quedará determinado a excepción de un factor de escala. Dada la importancia de los vectores y raíces características, obtenidos a partir de matrices simétricas, damos a continuación algunas propiedades interesantes.

- a) Si $\lambda_i \neq \lambda_j$ son raíces características de una matriz simétrica A , entonces x_i y x_j son vectores asociados son ortogonales. En efecto:

$$Ax_i = \lambda_i x_i \quad Ax_j = \lambda_j x_j$$

premultiplicando por x_j' y x_i' respectivamente

$$x_j' Ax_i = \lambda_i x_j' x_i \quad y \quad x_i' Ax_j = \lambda_j x_i' x_j$$

y dada la simetría de A , tenemos

$$\lambda_i x_j' x_i = \lambda_j x_i' x_j$$

y siendo $\lambda_i \neq \lambda_j$ sólo puede ocurrir que $x_j' x_i = 0$.

- b) Para cualquier matriz simétrica real A existe una matriz ortogonal P tal que

$$P'AP = D$$

donde D es una matriz diagonal cuyos elementos son las raíces características de A . Las columnas de P ^{pueden} ser los vectores característicos normalizados de A .

Estas propiedades para las matrices simétricas tienen una aplicación inmediata en las formas cuadráticas. En efecto, si aplicamos ~~esta~~ la transformación ortogonal

$$x = Py$$

a las variables de la forma cuadrática $x'Ax$, tendremos:

$$x'Ax = y'P'APy = y'Dy = \lambda_1 y_1^2 + \dots + \lambda_r y_r^2$$

7. MATRICES PARTICIONADAS

Se trata de presentar una matriz ~~como~~ función de las submatrices que la componen y que vienen determinadas por agrupaciones de determinadas filas y columnas requiriendo únicamente de homogeneidad, etc... Puede entonces una matriz aparecer como

$$A = \begin{bmatrix} A_{11} & \dots & A_{1n} \\ \vdots & & \vdots \\ A_{m1} & \dots & A_{mn} \end{bmatrix}$$

donde A_{ij} contiene r_i filas y c_j columnas y donde todas las submatrices de una misma fila tienen el mismo número de filas y todas las de una misma columna el mismo número de columnas. Con estas matrices pueden llevarse a cabo las mismas operaciones que en el caso no agrupado, teniendo en cuenta, en todo caso, la naturaleza no scalar de los elementos. Así por ejemplo si A y B tienen submatrices de igual dimensión podemos escribir

$$A+B = \begin{bmatrix} A_{11}+B_{11} & \dots & A_{1n}+B_{1n} \\ \vdots & & \vdots \\ A_{m1}+B_{m1} & \dots & A_{mn}+B_{mn} \end{bmatrix} \quad \text{o bien}$$

$$AB = \begin{bmatrix} \sum_{i=1}^n A_{1i} B_{i1} & \dots & \sum_{i=1}^n A_{1i} B_{ip} \\ \vdots & & \vdots \\ \sum_{i=1}^n A_{mi} B_{i1} & \dots & \sum_{i=1}^n A_{mi} B_{ip} \end{bmatrix}$$

Existen expresiones para la matriz inversa de una matriz particionada en función de la submatriz que la conforman, ¹ (4) así como para el determinante caso de tratarse de una matriz cuadrada no singular. Por ejemplo, para una matriz de la forma

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \quad \text{con } A_{11}, A_{22} \text{ cuadradas y } \text{no singulares} \text{ dada su naturaleza de matriz principal}$$

$$A^{-1} = \begin{bmatrix} (A_{11} - A_{12}A_{22}^{-1}A_{21})^{-1} & -(A_{11} - A_{12}A_{22}^{-1}A_{21})^{-1}A_{12}A_{22}^{-1} \\ -A_{22}^{-1}A_{21}(A_{11} - A_{12}A_{22}^{-1}A_{21})^{-1} & A_{22}^{-1} + A_{22}^{-1}A_{21}(A_{11} - A_{12}A_{22}^{-1}A_{21})^{-1}A_{12}A_{22}^{-1} \end{bmatrix}$$

para el determinante

$$|A| = |A_{11}| \cdot |A_{22} - A_{21}A_{11}^{-1}A_{12}| \quad \text{si } A_{11} \text{ es no singular, o bien}$$

$$|A| = |A_{22}| \cdot |A_{11} - A_{12}A_{22}^{-1}A_{21}| \quad \text{si } A_{22} \text{ es no singular.}$$

8. DIFERENCIACIÓN CON VECTORES Y MATRICES.

Sea $f(x)$ una función continua de los elementos del vector $x' = [x_1, \dots, x_p]$ cuyas primeras y segundas derivadas parciales existen para todos los puntos x de una determinada región del espacio euclídeo p -dimensional. El vector operador derivada parcial lo definimos como

$$\frac{\partial}{\partial x} = \begin{bmatrix} \frac{\partial}{\partial x_1} \\ \vdots \\ \frac{\partial}{\partial x_p} \end{bmatrix}$$

aplicándolo a $f(x)$ nos lleva a $\frac{\partial f(x)}{\partial x} = \left[\frac{\partial f(x)}{\partial x_1} \quad \dots \quad \frac{\partial f(x)}{\partial x_p} \right]$.

Algunos casos de funciones y sus derivadas especialmente importantes son:

1) la función constante

2) $f(x) = a'x$, $\frac{\partial f(x)}{\partial x} = a'$

3) $f(x) = x'Ax$ una forma cuadrática. Ecribamos la forma de la manera siguiente:

$$x'Ax = \sum_i \sum_j a_{ij} x_i x_j = \sum_i a_{ii} x_i^2 + 2 \sum_{i < j} a_{ij} x_i x_j + \sum_{\substack{i \neq j \\ j \neq i}} a_{ij} x_i x_j$$

derivando parcialmente respecto de x_i

$$\frac{\partial f(x)}{\partial x_i} = 2a_{ii}x_i + 2 \sum_{\substack{j=1 \\ j \neq i}}^p a_{ij}x_j = 2 \sum_{j=1}^p a_{ij}x_j = 2a_i'x$$

o por a_i' designamos la i -ésima fila de la matriz simétrica A . Entremos el vector de las derivadas parciales vendrá dado por

$$\frac{\partial f(x)}{\partial x} = 2Ax$$

4) Si se trata de una función cuadrática más general $f(x) = (a - Cx)'K(a - Cx)$ con K matriz $N \times N$, $C = [C_1, \dots, C_p]$ matriz de constantes $N \times p$, a un vector $N \times 1$.

haciendo $u = a - Cx$, las derivadas de $h(x)$ pueden calcularse aplicando la regla de la cadena:

$$\frac{\partial h(x)}{\partial x_i} = \sum_{j=1}^N \frac{\partial (u'ku)}{\partial u_j} \cdot \frac{\partial u_j}{\partial x_i} = -2 \sum_{j=1}^N k_j u \cdot c_{ji} = -2 u' k C_i \quad i=1, \dots, p$$

y de aquí

$$\frac{\partial h(x)}{\partial x} = -2 C' k (a - Cx).$$

5) La matriz de las derivadas de segundo orden (parciales) se la denomina Hessiano. En el caso de la forma cuadrática $f(x) = x'Ax$ dicha matriz es $2A$. El Hessiano es necesariamente simétrico si nuestras condiciones originales de continuidad y existencia de las derivadas parciales de primer y segundo orden son satisfechas por $f(x)$.

Determinación de Máximos y mínimos

Ya conocemos las condiciones necesarias para la existencia de máximos o mínimos. En cualquier lo más interesante para nosotros es las aplicaciones. Si queremos obtener los extremos estacionarios de $h(x)$ del apartado 4) anterior, tendremos que el valor x_0 donde ocurren el extremo debe satisfacer:

$$C' k C x = C' k a$$

y si $C' k C$ es de rango completo p , tendremos una solución única dada por $x = (C' k C)^{-1} C' k a$. Si k es definida positiva y C tiene rango p , el Hessiano $2C' k C$ es ~~positivo~~ positivo y representará de un mínimo.

Máximos y mínimos restringidos

Se utiliza ahora el método de los multiplicadores de Lagrange. Conste sencillamente en lo siguiente

$$h(x, \lambda) = f(x) - \lambda [g(x) - c]$$

donde $f(x)$ es la función a maximizar o minimizar y $g(x) = c$ es la restricción. Para un valor estacionario restringido tendremos

$$\frac{\partial h(x, \lambda)}{\partial x} = \frac{\partial f(x)}{\partial x} - \lambda \frac{\partial g(x)}{\partial x} = 0$$

$$\frac{\partial h(x, \lambda)}{\partial \lambda} = -g(x) + c = 0$$

En la práctica se obtiene x después de haber eliminado λ en el sistema de ecuaciones.

Veamos un ejemplo para el caso de una forma cuadrática $f(x) = x'Ax$, con A definida positiva, sujeta a la restricción $x'x = 1$. Entonces

$$h(x, \lambda) = x'Ax - \lambda (x'x - 1)$$

entonces

$$\frac{\partial h(x, \lambda)}{\partial x} = 2Ax - 2\lambda x = 0 \rightarrow [A - \lambda I]x = 0$$

que es la ecuación característica de A . Remultiplicando por x' y haciendo de la restricción tenemos

$$x'[A - \lambda I]x = 0 \rightarrow x'Ax - \lambda = 0 \rightarrow \lambda = x'Ax$$

entonces si $f(x)$ debe ser un máximo, λ deberá ser la mayor de todas las raíces características, y x el vector característico asociado. La restricción aún de que dicho vector debe tener módulo 1. Análogamente para el mínimo de $f(x)$.

La derivada de los determinante de A respecto de un elemento A_{ij} podemos encontrarla a partir del desarrollo de $|A|$ en cofactores de la i fila y la j columna. En efecto:

$$\frac{\partial |A|}{\partial A_{ij}} = \frac{\partial}{\partial A_{ij}} (a_{i1} A_{11} + \dots + a_{ij} A_{ij} + \dots + a_{in} A_{in}) = A_{ij}$$

Para derivar una matriz respecto de los elementos que la constituyen tendremos. Sea X una matriz $m \times n$ un elemento genérico x_{ij} y derivada respecto a este elemento será

$$\frac{\partial X}{\partial x_{ij}} = J_{ij}$$

donde J_{ij} es una matriz $m \times n$ con un 1 en la posición ij -ésima y ceros en el resto. Si X es simétrica

$$\frac{\partial X}{\partial x_{ij}} = J_{ij} + J_{ji} \quad i \neq j$$

La regla para derivar un producto de matrices es similar a la utilizada para escalar. Supongamos entonces que $X \in Y$ son dos matrices conformes y que ambas son función vectorial de z a través de $x_{ij}(z)$ e $y_{ij}(z)$ entonces

$$\frac{\partial XY}{\partial z} = \frac{\partial X}{\partial z} Y + X \frac{\partial Y}{\partial z}$$

Esta fórmula nos va a permitir obtener la derivada de la inversa de una matriz cuadrada no singular.

Veamos

$$I = XX^{-1}$$

$$\frac{\partial I}{\partial x_{ij}} = \frac{\partial X}{\partial x_{ij}} \cdot X^{-1} + X \cdot \frac{\partial X^{-1}}{\partial x_{ij}} = 0$$

$$J_{ij} X^{-1} + X \frac{\partial X^{-1}}{\partial x_{ij}} = 0 \rightarrow \frac{\partial X^{-1}}{\partial x_{ij}} = -X^{-1} J_{ij} X^{-1}$$

si X es simétrica

$$\frac{\partial X^{-1}}{\partial x_{ij}} = \begin{cases} -X^{-1} J_{ii} X^{-1} & i=j \\ -X^{-1} (J_{ij} + J_{ji}) X^{-1} & i \neq j \end{cases}$$

Capítulo 2

Muestras de una normal multivariante

1. VARIABLES ALEATORIAS MULTIDIMENSIONALES

Un vector aleatorio multidimensional \mathbf{X} , es un vector del tipo

$$\mathbf{X}' = [x_1, \dots, x_p]$$

cuyos elementos son variables aleatorias continuas cuyas densidades vienen dadas por $f_i(x_i)$, y cuyas funciones de distribución son $F_i(x_i)$ $i=1, \dots, p$. Recordemos que la función de distribución conjunta viene dada por

$$F(x_1, \dots, x_p) = \Pr(\mathbf{X}_1 \leq x_1, \dots, \mathbf{X}_p \leq x_p)$$

Si la función es absolutamente continua, existe $f(x_1, \dots, x_p)$ tal que

$$F(x_1, \dots, x_p) = \int_{-\infty}^{x_1} \dots \int_{-\infty}^{x_p} f(x_1, \dots, x_p) dx_1, \dots, dx_p$$

Si las variables $\mathbf{X}_1, \dots, \mathbf{X}_p$ son independientes la obtención de $f(x_1, \dots, x_p)$ y de $F(x_1, \dots, x_p)$ se simplificará enormemente aplicando el teorema de factorización, obteniendo

$$f(x_1, \dots, x_p) = f_1(x_1) \dots f_p(x_p)$$

$$F(x_1, \dots, x_p) = F_1(x_1) \dots F_p(x_p)$$

No siempre puede llevarse a cabo este supuesto de independencia, depende de las condiciones del problema. En cualquier caso, en análisis multivariante no siempre se verificará la independencia entre las variables observadas. La densidad conjunta de cualquier subconjunto de variables del conjunto original, se puede obtener mediante:

$$g(x_1, \dots, x_p) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f(x_1, \dots, x_{p+q}) dx_{p+1}, \dots, dx_{p+q}$$

para la función de distribución

$$G(x_1, \dots, x_p) = F(x_1, \dots, x_p, \infty, \dots, \infty)$$

Es importante que los límites de integración estén bien definidos, aunque formalmente no importan ser definidos entre $-\infty$ y $+\infty$. Por ejemplo, para X_1, X_2 variables aleatorias i.i.d. en f.d.p. $g(x)$ y f.d. $G(x)$ si queremos obtener las f.d.p. y f.d. de las nuevas variables

$$X_2 = \max(X_1, X_2)$$

$$X_1 = \min(X_1, X_2)$$

se obtiene fácilmente

$$f(x_1, x_2) = 2g(x_1)g(x_2) \quad -\infty < x_1 \leq x_2 < +\infty, \quad 0 \text{ en el resto}$$

las marginales vienen dadas por

$$f_1(x_1) = \int_{-\infty}^{+\infty} 2g(x_1)g(x_2)dx_2 = \int_{-\infty}^{x_1} 0 + \int_{x_1}^{+\infty} 2g(x_1)g(x_2)dx_2 = 2g(x_1)[1 - G(x_1)]$$

analogamente

$$f_2(x_2) = \int_{-\infty}^{+\infty} 2g(x_1)g(x_2)dx_1 = \int_{-\infty}^{x_2} 2g(x_1)g(x_2)dx_1 + \int_{x_2}^{+\infty} 0 = 2G(x_2)g(x_2)$$

Observa además que $f(x_1, x_2) \neq f_1(x_1)f_2(x_2)$ no hay pues independencia.

Condicionales

Se trata de obtener la densidad de X_1, \dots, X_p dado que las restantes variables han tomado un valor determinado $X_i = x_i \quad i = p+1, \dots, p+q$. Resulta que dicha densidad viene dada por

$$h(x_1, \dots, x_p / x_{p+1}, \dots, x_{p+q}) = \frac{f(x_1, \dots, x_{p+q})}{g(x_{p+1}, \dots, x_{p+q})}$$

donde $f(x_1, \dots, x_{p+q})$ es la densidad conjunta y $g(x_{p+1}, \dots, x_{p+q})$ la marginal de las variables que condicionan. En el caso de independencia la condicional se transforma simplemente en la marginal correspondiente. La función de distribución se obtiene mediante integración adecuada. Por ejemplo, continuando con las variables X_1, X_2 antes definidas.

$$f(x_1/x_2) = \frac{g(x_1)}{g(x_2)} \quad -\infty < x_1 \leq x_2 < \infty$$

o en el caso

la distribución condicional es

$$F(x_1/x_2) = \begin{cases} \frac{g(x_1)}{g(x_2)} & -\infty < x_1 \leq x_2 < \infty \\ 1 & x_2 \leq x_1 < \infty \end{cases}$$

Momentos

El vector media viene dado por

$$E(X') = [E(X_1), \dots, E(X_p)]$$

La extensión del concepto de varianzas y covarianzas al caso multidimensional viene dado por

$$E\{[X - E(X)][X - E(X)]'\} = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1p} \\ \sigma_{12} & \sigma_{22} & \dots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1p} & \sigma_{2p} & \dots & \sigma_{pp} \end{bmatrix} = \Sigma$$

Matriz que se denomina matriz de varianzas de X

Las expresiones matriciales para las varianzas y covarianzas de determinados cambios de variables vienen dadas por:

a) $Y = a_1 X_1 + \dots + a_p X_p = a'X \quad \text{var}(Y) = \text{var}(a'X) = \sum \sum a_i a_j \sigma_{ij} = a' \Sigma a$

b) $Y = a'X, Z = b'X \quad \text{cov}(Y, Z) = \sum a_i b_j \sigma_{ij} = a' \Sigma b$

c) Expresando de forma más general el cambio; A matriz $r \times p$ B matriz $s \times p$

$$Y = AX \quad Z = BX \quad \text{then} \quad \text{cov}(Y, Y) = A \Sigma A'$$

$$\text{cov}(Z, Z) = B \Sigma B'$$

$$\text{cov}(Y, Z) = A \Sigma B'$$

Finalmente el coeficiente de correlación viene definido mediante

$$\rho_{ij} = \frac{\text{cov}(X_i, X_j)}{\sqrt{\text{var}(X_i) \cdot \text{var}(X_j)}}$$

La matriz de correlaciones viene dada por

$$P = \begin{bmatrix} 1 & \rho_{12} & \dots & \rho_{1p} \\ \rho_{21} & 1 & \dots & \rho_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{p1} & \rho_{p2} & \dots & 1 \end{bmatrix}$$

La relación existente entre las matrices de covarianza y de correlaciones viene dada por:

$$P = D\left(\frac{1}{\sigma_i}\right) \Sigma D\left(\frac{1}{\sigma_i}\right)$$

$$\Sigma = D(\sigma_i) P D(\sigma_i)$$

donde $D(\sigma_i)$ es la matriz diagonal de las desviaciones típicas de las variables.

2. LA DISTRIBUCIÓN NORMAL MULTIVARIANTE

En adelante las muestras a las que hagamos referencia procederán de poblaciones normales multivariantes lo que pueden considerarse como una restricción puede justificarse en base a los siguientes motivos:

- a) Cualquier vector aleatorio que surja como la suma de un gran número de vectores aleatorios distribuidos idéntica e independientemente, tiene una distribución normal multivariante a medida que el número de estos vectores crece sin límite. Esto no es más que aplicar el teorema Central del límite en versión multivariante. Este tipo de modelos suma parece ser bastante realista a la hora de explicar gran parte de los fenómenos que surgen en biología y en ciencias del comportamiento.
- b) Diferentes modelos para los vectores considerados podrían conducir a muy diferentes distribuciones conjuntas de los elementos ~~considerados~~ cuya complejidad matemática impediría el desarrollo de las distribuciones muestrales de los ~~test~~ estadísticos y de los ~~estadísticos~~ usuales. Tales distribuciones deberían ser conocidas para cada una de las poblaciones fundamentales del modelo. Sin embargo, parece probable que con la excepción de algunos casos patológicos, el teorema central del límite, en su versión multivariante garantiza que las distribuciones para grandes muestras de los ~~test~~ estadísticos conducen a conclusiones similares acerca del estado de la naturaleza.

Hacia este, ~~sigamos~~, justificación del posterior uso, abando de la normal multivariante vamos en ella.

Recordemos que la distribución normal para el caso univariante tiene la expresión, en forma de función de densidad,

$$\phi(x) = \frac{1}{\sqrt{2\pi} \sigma} \exp \left[-\frac{1}{2} \left(\frac{x - \mu}{\sigma} \right)^2 \right] \quad -\infty < x < \infty$$

Para el caso de n variables, x_1, \dots, x_n , todas normales e independientes tenemos:

$$\phi(x_1, \dots, x_n) = \frac{1}{(\sqrt{2\pi})^n \sigma_1 \dots \sigma_n} \exp \left[-\frac{1}{2} \sum_{i=1}^n \left(\frac{x_i - \mu_i}{\sigma_i} \right)^2 \right]$$

y si escribimos

$$x' = [x_1, \dots, x_n] \quad \mu' = [\mu_1, \dots, \mu_n] \quad \Sigma = \begin{bmatrix} \sigma_1^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_n^2 \end{bmatrix}$$

la densidad conjunta anterior adquiere la expresión

$$\phi(x) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} z' \Sigma^{-1} z \right] \quad \text{con } z' = (x - \mu)' \quad , \quad z = (x - \mu)$$

~~Sea esta expresión~~

Esta expresión es una expresión particular de la densidad de una normal multivariante, porque, en efecto, normalmente el vector x considerado está constituido por variables aleatorias independientes. La generalización de la anterior expresión se obtiene haciendo que Σ sea una matriz simétrica definida positiva, entonces la función $\phi(x)$ es positiva para cualquier valor de x y se impone además que

$$\int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} \phi(x_1, \dots, x_n) dx_1 \dots dx_n = 1$$

es decir, $\phi(x)$ es una función de densidad. En esta notación tenemos que:

- el i -ésimo elemento del vector μ es la esperanza de x_i
- el i -ésimo elemento de la diagonal principal de Σ es la varianza de x_i
- el ij -ésimo elemento, σ_{ij} , de la matriz Σ es la covarianza de x_i y x_j . Evidentemente si las $P(1-1)/2$ covarianzas son nulas, entonces las x_i son independientes.

Consideremos con algún detalle el caso $n=2$, de gran importancia en estadística técnica.

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \quad \Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{bmatrix} = \begin{bmatrix} \sigma_{11} & \rho \sigma_1 \sigma_2 \\ \rho \sigma_1 \sigma_2 & \sigma_{22} \end{bmatrix}$$

la densidad conjunta será

$$\phi(x_1, x_2) = \frac{1}{2\pi \sigma_1 \sigma_2 \sqrt{1-\rho^2}} \exp \left\{ -\frac{1}{2} \frac{1}{1-\rho^2} \left[\left(\frac{x_1 - \mu_1}{\sigma_1} \right)^2 - 2\rho \frac{x_1 - \mu_1}{\sigma_1} \cdot \frac{x_2 - \mu_2}{\sigma_2} + \left(\frac{x_2 - \mu_2}{\sigma_2} \right)^2 \right] \right\}$$

si tipificamos las variables mediante la transformación

$$z_i = \frac{x_i - \mu_i}{\sigma_i} \quad \text{tenemos}$$

$$\phi(z_1, z_2) = \frac{1}{2\pi \sqrt{1-\rho^2}} \exp \left\{ -\frac{1}{2} \frac{1}{1-\rho^2} \left[z_1^2 - 2\rho z_1 z_2 + z_2^2 \right] \right\}$$

Posteriormente veremos cómo esta densidad depende del papel de ρ .

Ejes principales en la densidad multivariante

El exponente

$$(x - \mu)' \Sigma^{-1} (x - \mu)$$

es una forma cuadrática y por tanto especifica la ecuación de un elipsoide en el espacio n -dimensional cuando lo igualamos a alguna constante positiva. La familia de elipsoides que se genera cuando igualamos a distintos valores, c , de esta constante se caracterizan porque todos ellos están centrados en μ . El primer eje principal del elipsoide es aquella línea que pasa a través de su mayor dimensión. Si representamos a cualquier uno de los elipsoides que pasan por μ , mediante su punto de corte en la superficie del elipsoide, el eje principal tendrá una longitud que maximizará la longitud del segmento que μ y x determinan. El cuadrado de dicha longitud viene dado por

$$(x - \mu)' (x - \mu)$$

El cálculo de este primer eje vendrá pues dado mediante la solución de un problema de máximo condicionado. A saber, maximizar

$$(x - \mu)' (x - \mu)$$

Sujeto a la restricción

$$(x - \mu)' \Sigma^{-1} (x - \mu) = c$$

Aplicando en los multiplicados de lazo, tendremos

$$f(x) = (x-\mu)'(x-\mu) - \lambda [(x-\mu)' \Sigma^{-1} (x-\mu) - c]$$

y derivando

$$\frac{\partial f(x)}{\partial x} = 2(x-\mu) - 2\lambda \Sigma^{-1} (x-\mu) = 2[I - \lambda \Sigma^{-1}] (x-\mu) = 0$$

e igualando a cero

$$[I - \lambda \Sigma^{-1}] (x-\mu) = 0$$

Es decir, que las coordenadas del eje mayor deben satisfacer la anterior ecuación, que como Σ^{-1} es no singular, admite también la forma:

$$[\Sigma - \lambda I] (x-\mu) = 0$$

lo que da a entender que las coordenadas del eje principal son proporcionales a los ^{1º} los ^{1º} vectores característicos de Σ . Pero vamos a ver. Si premultiplicamos por $4(x-\mu)'$, tendremos:

$$4(x-\mu)' [x-\mu] - 4\lambda (x-\mu)' \Sigma^{-1} (x-\mu) = 0$$

$$4(x-\mu)' (x-\mu) = 4\lambda c$$

Por tanto para c fijo la longitud del eje principal se maximiza tomando λ igual a la mayor de las raíces características de Σ . Por tanto podemos afirmar:

"La posición del primer eje principal del elipsoide de concentración viene especificada por los nuevos directores del vector característico normalizado α_1 asociados a la mayor raíz característica λ_1 de Σ ".

Obsérvese que la longitud del eje viene entonces dada por $l = 2\sqrt{\lambda_1 c}$

El segundo eje principal vendrá dado por λ_2 , cuando dicho valor es la segunda mayor raíz característica de Σ . El resto se sigue hasta obtener los n ejes principales del elipsoide. Si las raíces características son distintas todas ellas,

$$\lambda_1 > \lambda_2 > \dots > \lambda_{n-1} > \lambda_n$$

entonces, recordando una propiedad que asegura que los vectores asociados son ortogonales

$$\alpha_i' \alpha_j = 0 \quad , \quad i \neq j$$

tendremos que los n ejes están unívocamente determinados y además, son todos ellos perpendiculares entre sí. Si hubieran dos raíces consecutivas iguales, quisiéramos decir, que el elipsoide tiene una sección circular sobre el plano que las direcciones de los vectores asociados independientes determinan. En este, aunque pueden encontrarse ejes perpendiculares entre sí, su posición no es única. En general si λ_i es una raíz característica, cuyo orden de multiplicidad es r_i , cada uno de los r_i ejes a determinar puede ser elegido perpendicular con los $r_i - 1$ restantes y con los $n - r_i$, aunque dichos ejes pueden ocupar una infinidad de direcciones principales. En tal caso, el elipsoide tiene forma lupusferica en el subespacio que determinan las r_i raíces características y se dice que tiene variación isotópica en este subespacio. Posteriormente nos ocuparemos más ampliamente de este tema cuando lo relacionemos con una determinada técnica estadística.

Consideremos ahora la nueva variable $Y' = [Y_1 \dots Y_n]$ cuyos elementos tienen valores en los ejes principales del elipsoide de concentración. Esta variable está relacionada con el vector original X , mediante la transformación

$$Y = A'(X - \mu)$$

donde la i -ésima columna de la matriz A es el vector característico normalizado α_i . La ortogonalidad de A implica que la transformación consiste en una rotación rígida de los ejes originales hasta coincidir con los ejes principales del elipsoide regido de una traslación del antiguo centro al centro μ del elipsoide. La matriz de covarianzas del nuevo vector Y vendrá dada por

$$\Sigma_Y = A' \Sigma A$$

que como sabemos es una matriz diagonal cuyos valores (en la diagonal principal son) λ_i . Por tanto la varianza de la variable que se mide en el i -ésimo eje principal vendrá dada por

$$\text{var}(Y_i) = \alpha_i' \Sigma \alpha_i = \lambda_i$$

por tanto

$$\text{cov}(Y_i, Y_j) = \alpha_i' \Sigma \alpha_j = 0 \quad i \neq j$$

a condición de que todas las λ_i sean distintas o bien, en caso contrario, los α_i se hayan elegido ortogonales.

Así pues la transformación de los ejes principales da lugar a variables lineales cuyas varianzas son proporcionales a la longitud de los ejes (a ser medido más exactamente).

Apliquemos este método a la normal bivariate con los componentes tipificados por aquellos de la uniplicidad.

~~Definamos la familia de elipsoides (elipses, ahora) mediante~~

$$h = \phi(z_1, z_2)$$

Su ecuación vendrá dada por

$$(1-\rho^2)c = z_1^2 - 2\rho z_1 z_2 + z_2^2$$

con $c = -2 \log(2\pi h \sqrt{1-\rho^2})$. Las raíces características de $\Sigma = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$ son $\lambda_1 = 1+\rho$, $\lambda_2 = 1-\rho$ y los vectores asociados $\alpha_1' = [\frac{1}{2}\sqrt{2}, \frac{1}{2}\sqrt{2}]$, $\alpha_2' = [\frac{1}{2}\sqrt{2}, -\frac{1}{2}\sqrt{2}]$.

Si ρ es positivo el eje principal mayor es la línea $z_1 = z_2$ y el menor la línea $z_1 = -z_2$. En caso contrario se invierten los ejes. Cuando $\rho = 0$, la elipse es un círculo y existen una infinidad de ejes principales, coincidiendo a los anteriores y a los de las variables originales. Señalamos por último que en este caso particular en que las varianzas son iguales (e iguales a 1) los ejes tienen la misma posición, independientemente del valor que tome ρ .

3. DISTRIBUCIONES MARGINALES Y CONDICIONALES DE LAS VARIABLES MULTINORMALES

Comenzaremos demostrando una proposición de utilidad posterior para determinar la distribución marginal de cualquier subconjunto de variables del vector $X = (X_1 \dots X_n)$.

PROPOSICIÓN.- Se X un vector aleatorio con n componentes con una distribución normal multivariante. Entonces

$$Y = CX$$

se distribuye como $N(C\mu, C\Sigma C')$ para C no singular.

Demostremos: Revenemos que la densidad de Y se obtiene a partir de la densidad de X , cambiando 2 (4)

X por

$$X = C^{-1}Y$$

y multiplicando por el jacobiano de la transformación. Sabríamos que en el caso de una transformación lineal viene dado por $\text{mod } |C^{-1}|$. Tendremos que

$$\text{mod } |C^{-1}| = \frac{1}{\text{mod } |C|} = \sqrt{\frac{1}{|C|^2}} = \sqrt{\frac{|\Sigma|}{|C| \cdot |\Sigma| \cdot |C|}} = \frac{|\Sigma|^{1/2}}{|C \Sigma C'|^{1/2}}$$

La forma cuadrática del exponente $Q = (X - \mu)' \Sigma^{-1} (X - \mu)$ se transformará en

$$\begin{aligned} Q &= (C^{-1}Y - \mu)' \Sigma^{-1} (C^{-1}Y - \mu) = (C^{-1}Y - C^{-1}C\mu)' \Sigma^{-1} (C^{-1}Y - C^{-1}C\mu) = \\ &= [C^{-1}(Y - C\mu)]' \Sigma^{-1} [C^{-1}(Y - C\mu)] = (Y - C\mu)' (C^{-1})' \Sigma^{-1} C^{-1} (Y - C\mu) = \\ &= (Y - C\mu)' (C \Sigma C')^{-1} (Y - C\mu) \end{aligned}$$

La densidad de Y vendrá entonces dada por

$$\begin{aligned} f(y) &= \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{1/2}} \cdot \frac{|\Sigma|^{1/2}}{|C \Sigma C'|^{1/2}} \cdot \exp \left\{ -\frac{1}{2} (Y - C\mu)' (C \Sigma C')^{-1} (Y - C\mu) \right\} = \\ &= \frac{1}{(2\pi)^{\frac{n}{2}} |C \Sigma C'|^{1/2}} \cdot \exp \left\{ -\frac{1}{2} (Y - C\mu)' (C \Sigma C')^{-1} (Y - C\mu) \right\}. \end{aligned}$$

Lo que demuestra la proposición.

Vamos ahora de estudiar las marginales. Vamos para ello a dividir el conjunto de variable en dos subconjuntos X_1, \dots, X_q , y X_{q+1}, \dots, X_n y formaremos los vectores

$$X^{(1)} = \begin{bmatrix} X_1 \\ \vdots \\ X_q \end{bmatrix}, \quad X^{(2)} = \begin{bmatrix} X_{q+1} \\ \vdots \\ X_n \end{bmatrix}$$

de manera que

$$(1) \quad X = \begin{bmatrix} X^{(1)} \\ X^{(2)} \end{bmatrix} = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix}$$

supongamos además que las variables que representan el vector X tienen una distribución conjunta normal multivariante con medias

$$E(X^{(1)}) = \mu^{(1)}, \quad E(X^{(2)}) = \mu^{(2)}$$

y covarianzas

$$\begin{aligned} E((X^{(1)} - \mu^{(1)})(X^{(1)} - \mu^{(1)})') &= \Sigma_{11} \\ E((X^{(2)} - \mu^{(2)})(X^{(2)} - \mu^{(2)})') &= \Sigma_{22} \\ E((X^{(1)} - \mu^{(1)})(X^{(2)} - \mu^{(2)})') &= \Sigma_{12} = \Sigma_{21}' = 0. \end{aligned}$$

Decimos entonces que el vector X ha sido particionado en (1) en dos subvectores, que el vector medio ha sido particionado análogamente en subvectores y que la matriz de covarianzas

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} = \begin{bmatrix} \Sigma_{11} & 0 \\ 0 & \Sigma_{22} \end{bmatrix}$$

he sido fraccional, análogamente, en matrices.

hemos demostrado ahora, que en las condiciones anteriormente impuestas a las matrices de Σ , las variables $X^{(1)}$ y $X^{(2)}$ se distribuyen independientemente y normalmente.

la inversa de Σ vendrá dada por

$$\Sigma^{-1} = \begin{bmatrix} \Sigma_{11}^{-1} & 0 \\ 0 & \Sigma_{22}^{-1} \end{bmatrix}$$

en cuanto a la forma cuadrática del exponente, tenemos

$$Q = (x - \mu)' \Sigma^{-1} (x - \mu) = [(x^{(1)} - \mu^{(1)})', (x^{(2)} - \mu^{(2)})'] \begin{bmatrix} \Sigma_{11}^{-1} & 0 \\ 0 & \Sigma_{22}^{-1} \end{bmatrix} \begin{bmatrix} (x^{(1)} - \mu^{(1)})' \\ (x^{(2)} - \mu^{(2)})' \end{bmatrix} =$$

$$= [(x^{(1)} - \mu^{(1)})' \Sigma_{11}^{-1}, (x^{(2)} - \mu^{(2)})' \Sigma_{22}^{-1}] \begin{bmatrix} x^{(1)} - \mu^{(1)} \\ x^{(2)} - \mu^{(2)} \end{bmatrix} =$$

$$= (x^{(1)} - \mu^{(1)})' \Sigma_{11}^{-1} (x^{(1)} - \mu^{(1)}) + (x^{(2)} - \mu^{(2)})' \Sigma_{22}^{-1} (x^{(2)} - \mu^{(2)}) = Q_1 + Q_2$$

por otra parte $|\Sigma| = |\Sigma_{11}| \cdot |\Sigma_{22}|$. Entonces la densidad de X podemos escribir como

$$f(x) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \cdot e^{-\frac{1}{2} Q} = \frac{1}{(2\pi)^{n_1/2} |\Sigma_{11}|^{1/2}} \cdot e^{-\frac{1}{2} Q_1} \cdot \frac{1}{(2\pi)^{n_2/2} |\Sigma_{22}|^{1/2}} \cdot e^{-\frac{1}{2} Q_2} = f(x^{(1)}) \cdot f(x^{(2)})$$

que es el producto de dos normales multivariantes, a saber, $N(\mu^{(1)}, \Sigma_{11})$, $N(\mu^{(2)}, \Sigma_{22})$.

la marginal de $X^{(1)}$ vendrá dada por la integral

$$\int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} f(x) dx_{q+1} \dots dx_n = f(x^{(1)}) \cdot \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} f(x^{(2)}) dx_{q+1} \dots dx_n = f(x^{(1)})$$

es decir, la marginal de $X^{(1)}$ es una $N(\mu^{(1)}, \Sigma_{11})$. Análogamente obteníamos $N(\mu^{(2)}, \Sigma_{22})$ para la marginal de $X^{(2)}$ y como además $N(\mu, \Sigma)$ puede factorizarse como producto de ambas marginales podemos afirmar que ambas variables son independientes, podemos enunciar el siguiente teorema:

TEOREMA :- Si X_1, \dots, X_n tienen una densidad conjunta normal multivariante, una condición necesaria y suficiente para que un subconjunto de estas variables y el subconjunto complementario sean independientes es que cada covarianza de una variable de un conjunto y una variable del otro sea nula.

la necesidad resulta de inmediato teniendo en cuenta que

$$\sigma_{ij} = E[(X_i - \mu_i)(X_j - \mu_j)] = \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} (x_i - \mu_i)(x_j - \mu_j) f(x^{(1)}) f(x^{(2)}) dx^{(1)} \cdot dx^{(2)} =$$

$$= \left\{ \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} (x_i - \mu_i) f(x^{(1)}) dx^{(1)} \right\} \left\{ \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} (x_j - \mu_j) f(x^{(2)}) dx^{(2)} \right\} = 0.$$

por tanto $\sigma_{ij} = 0$ y vemos que independencia e independencia son equivalentes en el caso de normalidad.

Consideremos el caso especial en que n vector de una normal ~~multivariante~~ bivariente. Es decir

2 (5)

$$\mathbf{X}^{(1)} = \mathbf{X}_1, \quad \mathbf{X}^{(2)} = \mathbf{X}_2, \quad \mu^{(1)} = \mu_1, \quad \mu^{(2)} = \mu_2, \quad \Sigma_{11} = \sigma_{11} = \sigma_1^2, \quad \Sigma_{22} = \sigma_{22} = \sigma_2^2, \quad \Sigma_{12} = \Sigma_{21} = \sigma_{12} = \sigma_1 \sigma_2 \rho_{12}.$$

Así, si $\mathbf{X}_1, \mathbf{X}_2$ tienen distribución normal bivariente, son independientes si y sólo si $\rho_{12} = 0$. En este caso la marginal de \mathbf{X}_i es normal con media μ_i y varianza σ_i^2 . Podemos concluir el siguiente

COROLARIO.- Si \mathbf{X} es un vector normal multivariante $N(\mu, \Sigma)$ y si un conjunto de variables de \mathbf{X} es independiente con otro conjunto (el complementario), las distribuciones marginales del conjunto es una normal multivariante con medias, varianzas y covarianzas obtenidas tomando las correspondientes componentes de μ y Σ respectivamente.

Pero vamos a ir más allá en nuestro resultado. Vamos a suponer que este resultado es válido aun cuando los subvectores obtenidos no sean independientes. Vamos para ello la siguiente transformación lineal no singular de los subvectores

$$\mathbf{Y}^{(1)} = \mathbf{X}^{(1)} + \mathbf{M} \mathbf{X}^{(2)}$$

$$\mathbf{Y}^{(2)} = \mathbf{X}^{(2)}$$

eligiendo \mathbf{M} de tal manera que los componentes de $\mathbf{Y}^{(1)}$ sean independientes con los componentes de $\mathbf{Y}^{(2)} = \mathbf{X}^{(2)}$. La matriz \mathbf{M} debe satisfacer, para ello, la siguiente ecuación

$$\begin{aligned} 0 &= E(\mathbf{Y}^{(1)} - E(\mathbf{Y}^{(1)}))(\mathbf{Y}^{(2)} - E(\mathbf{Y}^{(2)}))' = E((\mathbf{X}^{(1)} + \mathbf{M} \mathbf{X}^{(2)} - E(\mathbf{X}^{(1)}) - \mathbf{M} E(\mathbf{X}^{(2)}))(\mathbf{X}^{(2)} - E(\mathbf{X}^{(2)}))' = \\ &= E((\mathbf{X}^{(1)} - E(\mathbf{X}^{(1)})) + (\mathbf{M} \mathbf{X}^{(2)} - \mathbf{M} E(\mathbf{X}^{(2)})))(\mathbf{X}^{(2)} - E(\mathbf{X}^{(2)}))' = \Sigma_{12} + \mathbf{M} \Sigma_{22} \end{aligned}$$

y de aquí

$$\mathbf{M} = -\Sigma_{12} \Sigma_{22}^{-1}$$

y por tanto

$$\mathbf{Y}^{(1)} = \mathbf{X}^{(1)} - \Sigma_{12} \Sigma_{22}^{-1} \mathbf{X}^{(2)}$$

Es decir

$$\mathbf{Y} = \begin{bmatrix} \mathbf{Y}^{(1)} \\ \mathbf{Y}^{(2)} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & -\Sigma_{12} \Sigma_{22}^{-1} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{X}^{(1)} \\ \mathbf{X}^{(2)} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & -\Sigma_{12} \Sigma_{22}^{-1} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \mathbf{X}$$

Si una transformación no singular de \mathbf{X} , también podemos, de acuerdo con la primera proposición, una distribución normal dada por

$$E \begin{bmatrix} \mathbf{Y}^{(1)} \\ \mathbf{Y}^{(2)} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & -\Sigma_{12} \Sigma_{22}^{-1} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mu^{(1)} \\ \mu^{(2)} \end{bmatrix} = \begin{bmatrix} \mu^{(1)} - \Sigma_{12} \Sigma_{22}^{-1} \mu^{(2)} \\ \mu^{(2)} \end{bmatrix} = \begin{bmatrix} \mu^{(1)} \\ \mu^{(2)} \end{bmatrix} = \mu$$

y matriz de covarianzas

$$\mathbf{C}(\mathbf{Y}, \mathbf{Y}') = \begin{bmatrix} \mathbf{I} & -\Sigma_{12} \Sigma_{22}^{-1} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\Sigma_{21} \Sigma_{22}^{-1} & \mathbf{I} \end{bmatrix} = \begin{bmatrix} \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} & \mathbf{0} \\ \mathbf{0} & \Sigma_{22} \end{bmatrix}$$

Queda entonces que de acuerdo con el teorema anterior $\mathbf{Y}^{(1)}$ y $\mathbf{Y}^{(2)}$ son independientes y aplicando el corolario $\mathbf{X}^{(2)} = \mathbf{Y}^{(2)}$ tiene como marginal una normal $N(\mu^{(2)}, \Sigma_{22})$. Podemos entonces enunciar el siguiente teorema

TEOREMA.- Si \mathbf{X} tiene una distribución normal multivariante $N(\mu, \Sigma)$, la marginal de cualquier subconjunto de componentes de \mathbf{X} es normal multivariante con medias, varianzas y covarianzas obtenidas tomando las correspondientes componentes de μ y Σ respectivamente.

Consideremos ahora la transformación

$$Z = DX$$

donde Z tiene q componentes y D es una matriz real $q \times n$. El rango esperado de Z viene dado por

$$E(Z) = D\mu$$

y su matriz de covarianzas es

$$E(Z - D\mu)(Z - D\mu)' = D \Sigma D'$$

El caso $q=n$ y D no singular ya lo estudiamos al principio de este párrafo. Si $q \leq n$ y D es de rango q , podemos encontrar una matriz E , $(n-q) \times n$ de tal forma que

$$\begin{pmatrix} Z \\ W \end{pmatrix} = \begin{pmatrix} D \\ E \end{pmatrix} X$$

sea una transformación no singular. En este caso Z y W tienen distribución conjunta normal multivariante y Z de acuerdo con el teorema anterior tendrá una marginal normal cuyos parámetros ya conocemos. Tenemos en consecuencia lo siguiente teorema:

TEOREMA.- Si X es un vector normal multivariante $N(\mu, \Sigma)$, entonces $Z = DX$ se distribuye $N(D\mu, D \Sigma D')$, donde D es una matriz real $q \times n$ de rango $q \leq n$.

Podría ser posible llevar a cabo una generalización de este teorema para aquellos casos en que D sea una matriz rectangular, no necesariamente de rango completo. En cualquier caso no ocuparemos de ello en este momento.

Distribuciones condicionales

Vamos a tratar de encontrar las distribuciones condicionales derivadas de una normal multivariante.

Sea X un vector normal multivariante $N(\mu, \Sigma)$ con Σ no singular. Si dividimos a cabo la partición de los subvectores

$$X = \begin{bmatrix} X^{(1)} \\ X^{(2)} \end{bmatrix}$$

de q y $n-q$ componentes respectivamente. Recordemos que en el párrafo anterior obteníamos la distribución conjunta de los vectores $Y^{(1)} = X^{(1)} - \Sigma_{12} \Sigma_{22}^{-1} X^{(2)}$ e $Y^{(2)} = X^{(2)}$, cuya expresión era un producto de normales, a saber

$$n(Y^{(1)}) | \mu^{(1)} - \Sigma_{11} \Sigma_{22}^{-1} \mu^{(2)}, \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \cdot n(Y^{(2)} | \mu^{(2)}, \Sigma_{22})$$

La densidad conjunta de $X^{(1)}$ y $X^{(2)}$ podemos obtenerla a partir de la anterior expresión, substituyendo $Y^{(1)}$ por $X^{(1)} - \Sigma_{12} \Sigma_{22}^{-1} X^{(2)}$ e $Y^{(2)}$ por $X^{(2)}$ teniendo en cuenta que el Jacobiano de la transformación es 1. El resultado de esta substitución será:

$$f(x^{(1)}, x^{(2)}) = \frac{1}{(2\pi)^{q/2} |\Sigma_{11.2}|^{1/2}} \exp \left\{ -\frac{1}{2} \left[(x^{(1)} - \mu^{(1)}) - \Sigma_{12} \Sigma_{22}^{-1} (x^{(2)} - \mu^{(2)}) \right]' \Sigma_{11.2}^{-1} \left[(x^{(1)} - \mu^{(1)}) - \Sigma_{12} \Sigma_{22}^{-1} (x^{(2)} - \mu^{(2)}) \right] \right\} \\ \cdot \frac{1}{(2\pi)^{\frac{n-q}{2}} |\Sigma_{22}|^{1/2}} \exp \left\{ -\frac{1}{2} (x^{(2)} - \mu^{(2)})' \Sigma_{22}^{-1} (x^{(2)} - \mu^{(2)}) \right\}$$

$$\text{donde } \Sigma_{11.2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}.$$

Expresión ésta que debe dividirse, para obtener, en la densidad de X , a saber $n(x | \mu, \Sigma)$. El interés de la factorización anterior quedará de manifiesto a continuación. En efecto, la densidad condicional de $X^{(1)}$ dado que $X^{(2)} = x^{(2)}$ será el cociente de la anterior expresión y la densidad marginal de $X^{(2)}$, que es precisamente el segundo factor.

$$f(x^{(1)}/x^{(2)}) = \frac{1}{(2\pi)^{q/2} |\Sigma_{11.2}|^{1/2}} \exp \left\{ -\frac{1}{2} [(x^{(1)} - \mu^{(1)}) - \Sigma_{11} \Sigma_{22}^{-1} (x^{(2)} - \mu^{(2)})]' \Sigma_{11.2}^{-1} [(x^{(1)} - \mu^{(1)}) - \Sigma_{11} \Sigma_{22}^{-1} (x^{(2)} - \mu^{(2)})] \right\}$$

claramente, la densidad $f(x^{(1)}/x^{(2)})$ es una densidad normal q -variante con vector media

$$E(x^{(1)}/x^{(2)}) = \mu^{(1)} + \Sigma_{12} \Sigma_{22}^{-1} (x^{(2)} - \mu^{(2)}) = \nu(x^{(2)})$$

y matriz de covarianzas

$$E[(x^{(1)} - \nu(x^{(2)})) (x^{(1)} - \nu(x^{(2)}))'] = \Sigma_{11.2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

Hay que observar en estas expresiones que la media de $x^{(1)}$ dado $x^{(2)}$ es simplemente una función lineal de $x^{(2)}$ intentamos que la matriz de covarianzas correspondiente no dependa para nada de $x^{(2)}$.

DEFINICION.- la matriz $\Sigma_{12} \Sigma_{22}^{-1}$ es la matriz de regresión de $x^{(1)}$ sobre $x^{(2)}$.

El i, j -ésimo elemento de $\Sigma_{12} \Sigma_{22}^{-1}$ viene a menudo representado por

$$\beta_{ij}, i=1, \dots, q, j=1, \dots, n$$

Al vector $\mu^{(1)} + \Sigma_{12} \Sigma_{22}^{-1} (x^{(2)} - \mu^{(2)})$ se le denomina función de regresión.

Si por $\sigma_{ij}, i=1, \dots, q, j=1, \dots, n$ designamos el elemento i, j -ésimo de $\Sigma_{11.2}$, llamaremos a esta cantidad covarianza parcial.

DEFINICION.-

$$\rho_{ij}, i=1, \dots, q, j=1, \dots, n = \frac{\sigma_{ij}, i=1, \dots, q, j=1, \dots, n}{\sqrt{\sigma_{ii}, i=1, \dots, q, j=1, \dots, n} \cdot \sqrt{\sigma_{jj}, i=1, \dots, q, j=1, \dots, n}}$$

es el coeficiente de correlación parcial entre x_i y x_j cuando x_{q+1}, \dots, x_n permanecen fijas.

DEFINICION.- Al vector aleatorio

$$x_{1.2} = x^{(1)} - \mu^{(1)} - \Sigma_{12} \Sigma_{22}^{-1} (x^{(2)} - \mu^{(2)})$$

se le llama el conjunto de las variables residuales y representan las discrepancias de los elementos de $x^{(1)}$ y de la predicción de sus valores hecha a partir de la relación lineal del vector media de la distribución condicional con las variables de $x^{(2)}$.

Como ya vimos anteriormente

$$\text{Cor}(x^{(1)}, y^{(2)}) = 0 \rightarrow \text{Cor}(x_{1.2}, x^{(1)}) = 0$$

mientras que

$$\text{Cor}(y^{(1)}, y^{(2)}) = \Sigma_{11.2} \rightarrow \text{Cor}(x_{1.2}, x^{(1)}) = \Sigma_{11.2}$$

Como la numeración de los componentes de x es arbitraria y q también lo es, todo lo anterior sirve para definir la distribución condicional de cualquier conjunto de q componentes de x , fijando las $n-q$ restantes. Análogamente, dado que la marginal de cualquiera r componentes de x es normal, podemos también definir la distribución condicional de cualquiera q componentes ($q \leq r$) dado que las restantes $r-q$ permanecen fijas. Podemos enunciar el siguiente teorema.

TEOREMA.- Sea x un vector aleatorio n -dimensional dividido en dos grupos ~~componentes~~ que dan lugar a los subvectores $x^{(1)}$ y $x^{(2)}$. Supongamos que dividimos a la media μ de forma similar $\mu^{(1)}$ y $\mu^{(2)}$ y que la matriz de covarianzas Σ aparece dividida en las submatrices $\Sigma_{11}, \Sigma_{12}, \Sigma_{22}$ que son las matrices de covarianzas de $x^{(1)}$, de $x^{(1)}, x^{(2)}$ y de $x^{(2)}$ respectivamente. Entonces, si la distribución de x es normal, la distribución condicional de $x^{(1)}$ dado $x^{(2)} = x^{(2)}$ es también normal con media $\mu^{(1)} + \Sigma_{12} \Sigma_{22}^{-1} (x^{(2)} - \mu^{(2)})$ y matriz de covarianzas $\Sigma_{11.2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$.

Consideremos, como un ejemplo de todo lo expuesto el caso de una normal bivariable y vamos a encontrar la condicional de X_1 dado $X_2 = x_2$. En este caso $\mu^{(1)} = \mu_1$, $\mu^{(2)} = \mu_2$, $\Sigma_{11} = \sigma_1^2$, $\Sigma_{12} = \rho \sigma_1 \sigma_2$, $\Sigma_{22} = \sigma_2^2$. La matriz de los coeficientes de regresión vendrá dada por $\Sigma_{12} \Sigma_{22}^{-1} = \rho \sigma_1 / \sigma_2$ y la matriz de varianzas

$$\Sigma_{11.2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} = \sigma_1^2 - \rho \sigma_1 \sigma_2 \cdot \frac{1}{\sigma_2^2} \cdot \rho \sigma_1 \sigma_2^2 = \sigma_1^2 - \rho^2 \sigma_1^2 = \sigma_1^2 (1 - \rho^2)$$

entonces la densidad de Σ_1 , dado $\Sigma_2 = x_2$ será

$$f(x_1/x_2) = \frac{1}{[2\pi\sigma_1^2(1-\rho^2)]^{1/2}} \exp \left\{ -\frac{1}{2} (x_1 - \mu_1 - \rho \sigma_1/\sigma_2 (x_2 - \mu_2))^2 / \sigma_1^2(1-\rho^2) \right\}$$

Hay que hacer notar que para valores de ρ positivos, la media condicional aumenta a medida que ρ ~~aumenta~~ y disminuye si ρ es negativo.

Algebra considering geometrical

Una intuición geométrica de la teoría expuesta puede ayudar a comprenderla mejor. La densidad $f(x_1, x_2)$ puede ser considerada como una superficie $z = f(x_1, x_2)$ sobre el plano x_1, x_2 . Si intersectamos esta superficie mediante el plano $x_2 = c$ obtenemos una curva $z = f(x_1, c)$ sobre la línea $x_2 = c$ en el plano x_1, x_2 . Las ordenadas de esta curva son proporcionales a la densidad condicional de X_1 , dado $X_2 = x_2$; es decir, son proporcionales a las ordenadas de la curva de una distribución normal univariante. En el caso más general conviene considerar elipsoides de densidad constante en el espacio n -dimensional. Las superficies de densidad constante de $f(x_1, \dots, x_n)$ en la intersección de las superficies de densidad constante de $f(x_1, \dots, x_n)$, los hiperplanos $x_{n+1} = c_{n+1}, \dots, x_n = c_n$. Estos son también elipsoides.

Para una mayor claridad de las ideas poseas, todavía, hace algunas consideraciones acerca de una población real idealizada mediante una distribución normal. Consideremos por ejemplo una población de parejas padre-hijo. Si la población es esencialmente homogénea, las alturas de los padres y las alturas de los hijos tendrán aproximadamente una distribución normal. Una distribución condicional puede obtenerse considerando a los hijos de aquellos padres cuya altura sea, por ejemplo, 1.75 mts. Las alturas de estos hijos tendrán una distribución aproximadamente normal ~~multivariante~~ univariante. La media de esta distribución difiere de la media de las alturas de los hijos cuyos padres tienen una altura de 1.65 mts, por ejemplo, pero las varianzas serán las mismas.

Consideremos también las triptetas constituidas por las acturas de un padre y sus dos hijos mayores. La colección de acturas de ambos hijos para padres de 175 ans. de actura es la distribución condicional de dos variables; la correlación entre las acturas de los dos hijos mayores es un coeficiente de correlación parcial. El hecho de mantener constante la actura del padre elimina el efecto hereditario debido a los mismos. No obstante, cabría esperar una correlación positiva entre ambas acturas, puesto que los efectos hereditarios de la madre y los ambientes tenderían a causar en las acturas de los hermanos ~~efectos~~ variaciones aciliares (de la misma índole).

Como hemos visto anteriormente, una distribución condicional obtenida a partir de una distribución normal y normal con media una función lineal de las variables fijadas y matriz de covarianzas constante. En el caso de distribuciones no normales la distribución condicional de un conjunto de variables sobre otro no tiene necesariamente esta propiedad. No obstante, pueden construirse distribuciones no normales tal que algunas distribuciones condicionales tengan la citada propiedad. Esto puede hacerse tomando como densidad para X el producto $\pi(x/\mu^{(1)} + \beta(x^{(1)} - \mu^{(1)}), \Sigma_{11}) f(x^{(1)})$ donde $f(x^{(1)})$ es una densidad cualquiera.

El coeficiente de correlación múltiple

2 (7)

Consideremos nuevamente X particionado en $X^{(1)}$, $X^{(2)}$. Vamos a estudiar algunas propiedades de $\Sigma_{12} \Sigma_{22}^{-1} X^{(2)}$. Como solo estamos interesados, ahora, en funciones de las covarianzas y estas son invariantes por cambios de origen o de escala vamos a suponer $\mu=0$. Elijamos X_i una componente de $X^{(1)}$. Entonces

$$E(X_i / X^{(2)}) = \beta X^{(2)}$$

$$\text{donde } \beta = \sigma_{(i)} \Sigma_{22}^{-1}$$

viendo $\sigma_{(i)}$ la i -ésima fila de Σ_{12} definida mediante

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

Consideremos ahora la función lineal de las $X^{(1)}$, $(\beta X^{(2)})$. Hemos anteriormente que la covarianza entre $X^{(1)} - \Sigma_{12} \Sigma_{22}^{-1} X^{(2)}$ y $X^{(2)}$ es nula, entonces las variables $X_i - \beta X^{(2)}$ y $X^{(2)}$ son independientes.

Tratemos de buscar una función lineal $\alpha X^{(2)}$ para la que $(X_i - \alpha X^{(2)})$ tenga mínima variancia. Puesto que $E(Z^2) = E(ZZ')$ cuando Z es un escalar, la variancia será

$$\begin{aligned} E(X_i - \alpha X^{(2)})^2 &= E[(X_i - \beta X^{(2)}) + (\beta - \alpha) X^{(2)}]^2 = E[(X_i - \beta X^{(2)}) + (\beta - \alpha) X^{(2)}][(X_i - \beta X^{(2)}) + (\beta - \alpha) X^{(2)}]' \\ &= E[(X_i - \beta X^{(2)})(X_i - \beta X^{(2)})'] + (\beta - \alpha) E(X_i - \beta X^{(2)})(X^{(2)})' + (\beta - \alpha) E(X^{(2)})(X_i - \beta X^{(2)})' \\ &= E[(X_i - \beta X^{(2)})(X_i' - \beta X^{(2)'})] + (\beta - \alpha) E(X^{(2)} X^{(2)'}) (\beta - \alpha)' = \\ &= E(X_i X_i') - \beta E(X^{(2)} X_i') - E(X_i X^{(2)'}) \beta' + \beta E(X^{(2)} X^{(2)'}) \beta' + (\beta - \alpha) E(X^{(2)} X^{(2)'}) (\beta - \alpha)' = \\ &= \sigma_{ii} - \beta \sigma_{(i)}' - \sigma_{(i)} \beta' + \beta \Sigma_{22} \beta' + (\beta - \alpha) \Sigma_{22} (\beta - \alpha)' = \\ &= \sigma_{ii} - \sigma_{(i)} \Sigma_{22}^{-1} \sigma_{(i)}' - \sigma_{(i)} \Sigma_{22}^{-1} \sigma_{(i)}' + \sigma_{(i)} \Sigma_{22}^{-1} \sigma_{(i)}' + (\beta - \alpha) \Sigma_{22} (\beta - \alpha)' = \\ &= (\sigma_{ii} - \sigma_{(i)} \Sigma_{22}^{-1} \sigma_{(i)}') + (\beta - \alpha) \Sigma_{22} (\beta - \alpha)' \end{aligned}$$

Puesto Σ_{22} es definida positiva el segundo término de la última línea es no negativo y alcanzará su mínimo cuando sea nulo, es decir $\alpha = \beta$. Así, la función de regresión es la función de $X^{(2)}$ tal que $(X_i - \alpha X^{(2)})$ tiene mínima variancia.

Vamos ahora a demostrar que la máxima correlación entre X_i y $X^{(2)}$ se obtiene para $\alpha = \beta$. Sabemos que

$$E(X_i - \beta X^{(2)})^2 \leq E(X_i - \alpha X^{(2)})^2 \quad \forall \alpha \neq \beta.$$

Por tanto

$$\sigma_{ii} + E(\beta X^{(2)})^2 - 2E(X_i \beta X^{(2)}) \leq \sigma_{ii} + c^2 E(\alpha X^{(2)})^2 - 2cE(X_i \alpha X^{(2)})$$

De aquí

$$2 \frac{E(X_i \beta X^{(2)})}{\sqrt{\sigma_{ii}} \sqrt{E(\beta X^{(2)})^2}} - \frac{E(\beta X^{(2)})^2}{\sqrt{\sigma_{ii}} \sqrt{E(\beta X^{(2)})^2}} \geq 2c \frac{E(X_i \alpha X^{(2)})}{\sqrt{\sigma_{ii}} \sqrt{E(\alpha X^{(2)})^2}} - c^2 \frac{E(\alpha X^{(2)})^2}{\sqrt{\sigma_{ii}} \sqrt{E(\alpha X^{(2)})^2}}$$

$$\text{eligiendo } c^2 = \frac{E(\beta X^{(2)})^2}{E(\alpha X^{(2)})^2}$$

tenemos

$$\frac{E(X_i \beta X^{(2)})}{\sqrt{\sigma_{ii}} \sqrt{E(\beta X^{(2)})^2}} \geq \frac{E(X_i \alpha X^{(2)})}{\sqrt{\sigma_{ii}} \sqrt{E(\alpha X^{(2)})^2}}$$

Podemos ahora enunciar el siguiente teorema:

TEOREMA .- Sea \mathbf{X} un vector aleatorio normal multivariante $N(\mu, \Sigma)$. Sean $\mathbf{X}' = [\mathbf{X}^{(1)'} \mathbf{X}^{(2)'}]$, $\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$ y sea β la i -ésima fila de $\Sigma_{12} \Sigma_{22}^{-1}$, $i=1, \dots, q$. De todas las combinaciones lineales $\alpha \mathbf{X}^{(2)}$, la combinación que minimiza la varianza de $\mathbf{X}_i - \alpha \mathbf{X}^{(2)}$ y que maximiza la correlación entre \mathbf{X}_i y $\alpha \mathbf{X}^{(2)}$ es la combinación lineal $\beta \mathbf{X}^{(2)}$.

DEFINICION .- la máxima correlación entre \mathbf{X}_i y la combinación lineal $\alpha \mathbf{X}^{(2)}$ se denomina coeficiente de correlación múltiple entre \mathbf{X}_i y $\mathbf{X}^{(2)}$.

Se sigue que dicho coeficiente es

$$\bar{R}_{i, q+1, \dots, n} = \frac{E(\beta \mathbf{X}^{(2)} \mathbf{X}_i)}{\sqrt{\sigma_{ii}} \sqrt{E(\beta \mathbf{X}^{(2)} \mathbf{X}^{(2)'} \beta)}} = \frac{\sigma_{ii} \Sigma_{22}^{-1} \sigma'_{ii}}{\sqrt{\sigma_{ii}} \sqrt{\sigma_{ii} \Sigma_{22}^{-1} \sigma'_{ii}}} = \frac{\sqrt{\sigma_{ii} \Sigma_{22}^{-1} \sigma'_{ii}}}{\sqrt{\sigma_{ii}}}$$

Una fórmula útil es la siguiente

$$1 - \bar{R}_{i, q+1, \dots, n}^2 = \frac{\sigma_{ii} - \sigma_{ii} \Sigma_{22}^{-1} \sigma'_{ii}}{\sigma_{ii}} = \frac{|\Sigma^*|}{\sigma_{ii} |\Sigma_{22}|}$$

donde $\Sigma^* = \begin{bmatrix} \sigma_{ii} & \sigma_{ii} \\ \sigma'_{ii} & \Sigma_{22} \end{bmatrix}$ (recordemos que el determinante puede obtenerse en función de los determinantes de las submatrices mediante la expresión $|\Sigma^*| = |\Sigma_{22}| \cdot |\sigma_{ii} - \sigma_{ii} \Sigma_{22}^{-1} \sigma'_{ii}|$)

Puesto que

$$\sigma_{ii, q+1, \dots, n} = \sigma_{ii} - \sigma_{ii} \Sigma_{22}^{-1} \sigma'_{ii} \quad (\text{elemento de la diagonal de } \Sigma_{11,2})$$

se sigue que

$$1 - \bar{R}_{i, q+1, \dots, n}^2 = \frac{\sigma_{ii, q+1, \dots, n}}{\sigma_{ii}} \rightarrow \sigma_{ii, q+1, \dots, n} = \sigma_{ii} (1 - \bar{R}_{i, q+1, \dots, n}^2)$$

Lo que demuestra que cualquier varianza condicional de un componente de \mathbf{X} no puede ser mayor que la varianza. De efecto, cuanto mayor es $\bar{R}_{i, q+1, \dots, n}$, mayor es la reducción en varianza para la distribución condicional.

~~Recordemos también que el coeficiente de correlación múltiple puede verse como el coeficiente~~

Recordemos lo referente al coeficiente de correlación múltiple demostrando la siguiente propiedad:

PROPIEDAD .- El coeficiente de correlación múltiple es invariante bajo transformaciones no singulares de las variables originales.

$$\begin{aligned} \text{Supongamos } Y_i &= a \mathbf{X}_i + b \\ Y^{(2)} &= C \mathbf{X}^{(2)} + d \end{aligned} \quad (1)$$

donde a, b son escalares, C es una matriz no singular y d es un vector de constantes. Entonces el cuadrado del coeficiente de correlación múltiple de Y_i respecto de $Y^{(2)}$ es

$$\begin{aligned} \bar{R}_{i, q+1, \dots, n}^2 &= \frac{\{E(\beta^* Y^{(2)} Y_i)\}^2}{\sigma_{ii}^* \cdot E(\beta^* Y^{(2)} Y^{(2)} \beta^*)} = \frac{\{E(\beta^* C \mathbf{X}^{(2)} a \mathbf{X}_i)\}^2}{a^2 \sigma_{ii} \cdot E(\beta^* C \mathbf{X}^{(2)} (C \mathbf{X}^{(2)})' \beta^*)} \\ &= \frac{\{a \sigma_{ii} C' (C \Sigma_{22} C')^{-1} C \sigma'_{ii}\}^2}{a^2 \sigma_{ii} \cdot (C \sigma_{ii} C' (C \Sigma_{22} C')^{-1} C \sigma'_{ii})} = \frac{\sigma_{ii} C' (C \Sigma_{22} C')^{-1} C \sigma'_{ii}}{\sigma_{ii}} = \frac{\sigma_{ii} \Sigma_{22}^{-1} \sigma'_{ii}}{\sigma_{ii}} \end{aligned}$$

Esta propiedad implica que la misma correlación puede obtenerse indistintamente partiendo de la matriz de correlación o de la matriz de covarianzas. Los coeficientes de regresión de las variables transformadas (1)

$$y = a [\sigma_{ii} C'] [C \Sigma_{22} C']^{-1} = a \sigma_{ii} C' C^{-1} \Sigma_{22}^{-1} C^{-1} = a (\sigma_{ii} \Sigma_{22}^{-1}) C^{-1} = a \beta C^{-1}$$

Si C es una matriz diagonal de factores de escala, por ejemplo los inversos de las desviaciones standard σ_i de las variables en $X^{(2)}$, el efecto de este cambio de escala consiste en el producto de cada elemento de β por σ_i . Por ejemplo, la regresión bivariente transformaría sus coeficientes en la cantidad adimensional ρ haciendo $a = \frac{1}{\sigma_1}$, $c = \frac{1}{\sigma_2}$.

Finalmente, si Σ y Σ_{22} tienen igual rango, q , esto significa que X_1 puede ser expresada exactamente como una combinación lineal de las q variables que componen $X^{(2)}$. En este caso la correlación múltiple es exactamente la unidad y β es el vector de los coeficientes de la combinación lineal. Si Σ_{22} tiene un rango menor que q , la expresión del coeficiente de correlación múltiple, vector univariado y la correlación múltiple puede ser redefinida como la máxima correlación de X_1 con una combinación lineal de un subconjunto de $X^{(2)}$ cuyo rango coincide con el de Σ_{22} .

4. FUNCIÓN CARACTERÍSTICA. MOMENTOS

Definición.- La función característica de un vector aleatorio X es

$$\phi(t) = E(e^{it'X}) \quad \text{para cualquier vector real } t.$$

Para darle un significado a esta definición definiremos la esperanza de una función vectorial aleatoria a valores complejos, mediante

Definición.- Sea $g(x)$ una función del vector aleatorio X a valores complejos, que podemos escribir de la forma $g(x) = g_1(x) + i g_2(x)$, donde $g_1(x)$ y $g_2(x)$ son funciones reales. Entonces la esperanza de $g(x)$ viene dada por

$$E(g(X)) = E(g_1(X)) + i E(g_2(X))$$

En particular

$$E(e^{it'X}) = E(\cos t'X) + i E(\sin t'X).$$

No es nuestro objetivo entrar en detalles acerca de las propiedades de la función característica pero sí mencionar algunas de gran importancia y utilidad.

Recordemos en primer lugar que el teorema de factorización, cuando se trata de variables independientes, también es verificado por las funciones características, de manera que si $X' = (X^{(1)}, X^{(2)})$ (donde $X^{(1)}, X^{(2)}$ independientes entre sí)

$$\phi(t) = E(e^{it'X}) = E(e^{it^{(1)'}X^{(1)}}) \cdot E(e^{it^{(2)'}X^{(2)}}) = \phi_1(t^{(1)}) \cdot \phi_2(t^{(2)}) \quad \text{con } t' = (t^{(1)'}, t^{(2)'})$$

Así, si X es tal que está constituido por n variables independientes tendremos

$$\phi(t) = E(e^{it'X}) = \prod_{j=1}^n E(e^{it_j'X_j}) = \prod_{j=1}^n \phi_j(t_j)$$

Esta propiedad nos va a permitir obtener la función característica para un vector X normal multivariante.

Teorema.- Sea X un vector aleatorio $N(\mu, \Sigma)$, entonces su función característica viene dada por

$$\phi(t) = E(e^{it'X}) = e^{it'\mu - \frac{1}{2} t'\Sigma t}, \quad \forall t, \text{ vector real.}$$

Demostración.- Sabemos de la existencia de una matriz C no singular, tal que

$$C'XC = I$$

y de aquí

$$\Sigma^{-1} = C^{-1'} \cdot C^{-1} = (C'C)^{-1}$$

Hagamos

$$X - \mu = CY$$

entonces Y se distribuirá $N(0, I)$, como ya sabemos. La función característica de Y , vendrá dada por

$$\Psi(u) = E(e^{iu'x}) = \prod_{j=1}^n E(e^{iu_j y_j})$$

y como $y_j \sim N(0,1)$, tenemos

$$\Psi(u) = \prod_{j=1}^n E(e^{iu_j y_j}) = \prod_{j=1}^n e^{-\frac{1}{2} u_j^2} = e^{-\frac{1}{2} u'u}$$

A partir de aquí

$$\phi(t) = E(e^{it'x}) = E(e^{it'(Gx+\mu)}) = e^{it'\mu} E(e^{it'Gx}) = e^{it'\mu} e^{-\frac{1}{2} (t'G)(t'G)'} =$$

para $t'G = u'$. Esta cadena de igualdades se obtiene fácilmente escribiendo los sumandos como integrales y teniendo en cuenta las reglas que rigen los cambios de variable en una integral. (~~función característica de una variable aleatoria es la esperanza de su exponencial complejo~~). Tendremos finalmente

$$\phi(t) = e^{it'\mu - \frac{1}{2} t'G G' t} = e^{it'\mu - \frac{1}{2} t' \Sigma t}$$

Con el auxilio de los dos teoremas que anunciaremos a continuación y del anterior resultado podrían haberse obtenido también más rápidamente algunos de los anteriores resultados.

Teorema (de Levy). - Si la variable X tiene una densidad $f(x)$ y una función característica $\phi(t)$, entonces

$$f(x) = \frac{1}{(2\pi)^n} \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} e^{-it'x} \phi(t) dt_1 \dots dt_n.$$

Teorema (de la continuidad de Levy-Cramer). - Sea $\{F_j(x)\}$ una sucesión de funciones de distribución y sea $\{\phi_j(t)\}$ la correspondiente sucesión de funciones características. Una condición necesaria y suficiente para que $F_j(x)$ converja a una distribución $F(x)$ es que, para cualquier t , $\phi_j(t)$ converja a un límite $\phi(t)$, continua en $t=0$. Entonces el límite $\phi(t)$ es además la función característica de la distribución límite $F(x)$.

Obsérvese que el primero de estos teoremas lo que afirma es que las funciones de densidad y características de un vector aleatorio están en correspondencia biunívoca. Entonces:

Si $Z = DX$, con $X \sim N(\mu, \Sigma)$, la función característica de Z tendrá por expresión

$$\begin{aligned} \phi(t) &= E(e^{it'Z}) = E(e^{it'DX}) = E(e^{i(D't)'X}) = e^{i(D't)'\mu - \frac{1}{2} (D't)'\Sigma(D't)} = \\ &= e^{it'(D\mu) - \frac{1}{2} t'(D\Sigma D')t} \end{aligned}$$

aplicando el teorema de Levy llegamos a la conclusión de que Z es normal con media $D\mu$ y matriz de covarianzas $D\Sigma D'$.

Análogamente al uso de las funciones características y sus propiedades nos permite demostrar el siguiente

Teorema. - Si cualquier combinación lineal de los componentes de un vector X se distribuye normalmente, entonces X se distribuye normalmente.

Demostración. - Supongamos la combinación lineal $Z = u'X$ que se distribuye normalmente, tendremos que su función característica vendrá dada por

$$E(e^{it'Z}) = E(e^{it'u'X}) = e^{it'u'\mu - \frac{1}{2} t'u'\Sigma u}$$

donde μ y Σ la media y matriz de covarianzas de X . Si hacemos ahora $t=1$, tendremos que

$$E(e^{iu'X}) = e^{iu'\mu - \frac{1}{2} u'\Sigma u}, \text{ por tanto } X \text{ se distribuirá normalmente.}$$

Es muy importante insistir en el hecho de que han de ser todas las combinaciones lineales de las componentes de \mathbf{X} . Puede verse un ejemplo en Anderson pags. 37, 38. 2 (9)

Momentos de un vector aleatorio normal

Los momentos de $\mathbf{Z}_1, \dots, \mathbf{Z}_p$ en una distribución conjunta normal pueden obtenerse a partir de la función característica. A saber:

$$E(\mathbf{Z}_j) = \frac{1}{i} \frac{\partial \phi}{\partial t_j} \Big|_{t=0} = \frac{1}{i} \left(- \sum_k \sigma_{jk} t_k + i \mu_j \right) \phi(t) \Big|_{t=0} = \mu_j$$

$$\begin{aligned} E(\mathbf{Z}_k \mathbf{Z}_j) &= \frac{1}{i^2} \frac{\partial^2 \phi}{\partial t_k \partial t_j} \Big|_{t=0} = \frac{1}{i^2} \left\{ \left(- \sum_k \sigma_{kk} t_k + i \mu_k \right) \left(- \sum_k \sigma_{kj} t_k + i \mu_j \right) - \sigma_{kj} \right\} \phi(t) \Big|_{t=0} = \\ &= \sigma_{kj} + \mu_k \mu_j \end{aligned}$$

análogamente.

5. MUESTRAS DE UNA NORMAL MULTIVARIANTE

En todo cuanto precede hemos estudiado las propiedades de una distribución normal multivariante como si los valores de sus parámetros fueran conocidos y estuvieran bajo control. Esta situación solo suelta raras veces. Desgraciadamente en biología y ciencias del comportamiento no ocurre así generalmente, y lo que pretendemos en este apartado es considerar métodos que permitan estimar los diferentes parámetros de la normal multivariante a partir de muestras relativamente pequeñas, estudiar, después, las propiedades muestrales de estas estimaciones.

Antes de entrar en detalles acerca de la obtención de las estimaciones insistir en la imperiosa necesidad de que la aleatorización de la muestra considerada sea cierta, ni siquiera que nuestras observaciones tengan alguna validez. Es decir, se trata de que las unidades muestrales hayan sido tomadas independientemente unas de otras, de una población homogénea. Estas unidades muestrales no pueden tener características comunes o rasgos que pudieran indicar alguna independencia entre ellas. Por ejemplo, una investigación de los niveles medios de cuatro componentes bioquímicos en el cerebro de cierta raza de ratas, no puede estar basada en muestras extraídas de unas pocas camadas de ratas, por grandes que estas sean, sino en muestras extraídas aleatoriamente de ratas de la misma especie. De la misma manera, obtenemos estimaciones vagas y poco fiables del efecto general de una droga tranquilizante si la muestra está basada en las respuestas diarias de un pequeño número de pacientes psiquiátricos estudiados durante varias semanas. De la misma forma que elementos de una misma camada de ratas tienen probablemente rasgos biológicos y genéticos comunes, parece razonable pensar que una medición afectiva del paciente tiene rasgos asociados tan ligeros que harán las contribuciones de días sucesivos altamente dependientes. Además, variaciones importantes en el comportamiento de una persona en la sala del hospital donde se encuentran los enfermos ocasionados indicarán probablemente cambios en los otros sujetos del estudio.

Estimaciones máximo-verosímiles del vector media y de la matriz de covarianzas

Quisiéramos obtener las estimaciones máximo-verosímiles de μ y Σ a partir de una muestra de N observaciones procedentes de una población normal p -variante cuyos parámetros son precisamente μ y Σ .

La decisión, precisamente, de este tipo de estimados se debe a que son los o algunos simples transformaciones de los mismos poseen, usualmente, algunas propiedades óptimas de los estimados. En el caso particular que ahora nos ocupa, los estimados son asintóticamente eficientes.

Supongamos que nuestra muestra de tamaño N está constituida por los vectores p -variante $\mathbf{x}_1, \dots, \mathbf{x}_N$, $N > p$. La función de verosimilitud de la muestra viene dada por

$$L = \frac{1}{(2\pi)^{\frac{1}{2}pN} |\Sigma|^{\frac{1}{2}N}} \exp \left[-\frac{1}{2} \sum_{\alpha=1}^N (x_{\alpha} - \mu)' \Sigma^{-1} (x_{\alpha} - \mu) \right]$$

Como el exponente aparece en términos de Σ^{-1} obtendremos primero los estimados de μ y $\Sigma^{-1} = \Psi$. Recordemos que en la función L , las variables son μ y Σ , mientras que x_{α} son valores muestrales y por tanto perfectamente conocidos.

Para maximizar L , lo haremos a través de su log, por cuanto el log es una función creciente y el máximo coincidirá. Así pues

$$\log L = -\frac{1}{2} p N \log(2\pi) + \frac{1}{2} \log |\Psi| - \frac{1}{2} \sum_{\alpha=1}^N (x_{\alpha} - \mu)' \Psi (x_{\alpha} - \mu).$$

Definamos la media muestral como

$$\bar{x} = \frac{1}{N} \sum_{\alpha=1}^N x_{\alpha} = \begin{bmatrix} \frac{1}{N} \sum_{\alpha=1}^N x_{1\alpha} \\ \vdots \\ \frac{1}{N} \sum_{\alpha=1}^N x_{p\alpha} \end{bmatrix} = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{bmatrix}$$

y la matriz de los sumas de cuadrados y productos cruzados de las desviaciones respecto a la media

$$A = \sum_{\alpha=1}^N (x_{\alpha} - \bar{x})(x_{\alpha} - \bar{x})' = \begin{bmatrix} \sum_{\alpha=1}^N (x_{1\alpha} - \bar{x}_1)(x_{1\alpha} - \bar{x}_1) & \dots & \sum_{\alpha=1}^N (x_{1\alpha} - \bar{x}_1)(x_{p\alpha} - \bar{x}_p) \\ \vdots & \ddots & \vdots \\ \sum_{\alpha=1}^N (x_{p\alpha} - \bar{x}_p)(x_{1\alpha} - \bar{x}_1) & \dots & \sum_{\alpha=1}^N (x_{p\alpha} - \bar{x}_p)(x_{p\alpha} - \bar{x}_p) \end{bmatrix} \quad i, j = 1, \dots, p$$

hemos a tener de dos otras expresiones a $\log L$, para ello necesitaremos algunas propiedades que expusimos anteriormente.

LEMA.- Sean x_1, \dots, x_N N vectores y sea \bar{x} la media antes definida. Entonces para cualquier vector b

$$\sum_{\alpha=1}^N (x_{\alpha} - b)(x_{\alpha} - b)' = \sum_{\alpha=1}^N (x_{\alpha} - \bar{x})(x_{\alpha} - \bar{x})' + N(\bar{x} - b)(\bar{x} - b)'$$

Demostración.-

$$\begin{aligned} \sum_{\alpha=1}^N (x_{\alpha} - b)(x_{\alpha} - b)' &= \sum_{\alpha=1}^N [(x_{\alpha} - \bar{x}) + (\bar{x} - b)][(x_{\alpha} - \bar{x}) + (\bar{x} - b)]' = \\ &= \sum_{\alpha=1}^N [(x_{\alpha} - \bar{x})(x_{\alpha} - \bar{x})' + (x_{\alpha} - \bar{x})(\bar{x} - b)' + (\bar{x} - b)(x_{\alpha} - \bar{x})' + (\bar{x} - b)(\bar{x} - b)'] = \\ &= \sum_{\alpha=1}^N (x_{\alpha} - \bar{x})(x_{\alpha} - \bar{x})' + \left[\sum_{\alpha=1}^N (x_{\alpha} - \bar{x}) \right] (\bar{x} - b)' + (\bar{x} - b) \left[\sum_{\alpha=1}^N (x_{\alpha} - \bar{x})' \right] + N(\bar{x} - b)(\bar{x} - b)' = \\ &= \sum_{\alpha=1}^N (x_{\alpha} - \bar{x})(x_{\alpha} - \bar{x})' + N(\bar{x} - b)(\bar{x} - b)' \end{aligned}$$

$$\text{pues } \sum_{\alpha=1}^N (x_{\alpha} - \bar{x}) = N\bar{x} - N\bar{x} = 0.$$

Si hacemos $b = \mu$, tendremos

$$\sum_{\alpha=1}^N (x_{\alpha} - \mu)(x_{\alpha} - \mu)' = \sum_{\alpha=1}^N (x_{\alpha} - \bar{x})(x_{\alpha} - \bar{x})' + N(\bar{x} - \mu)(\bar{x} - \mu)' = A + N(\bar{x} - \mu)(\bar{x} - \mu)'$$

Por otra parte haciendo uso de la propiedad de la traza de una matriz que dice: $\text{tr}(CD) = \text{tr}(DC) = \sum_{i,j} c_{ij}d_{ji}$, tenemos

$$\begin{aligned} \sum_{\alpha=1}^N (x_{\alpha} - \mu)' \Psi (x_{\alpha} - \mu) &= \text{tr} \sum_{\alpha=1}^N (x_{\alpha} - \mu)' \Psi (x_{\alpha} - \mu) = \text{tr} \sum_{\alpha=1}^N \Psi (x_{\alpha} - \mu)(x_{\alpha} - \mu)' = \text{tr} \sum_{\alpha=1}^N \Psi [A + N(\bar{x} - \mu)(\bar{x} - \mu)'] = \\ &= \text{tr} \Psi A + \text{tr} \Psi N(\bar{x} - \mu)(\bar{x} - \mu)' = \text{tr} \Psi A + N(\bar{x} - \mu)' \Psi (\bar{x} - \mu) \end{aligned}$$

Podemos entonces escribir el $\log L$ de la siguiente manera

$$\log L = -\frac{1}{2} p \log(2\pi) + \frac{1}{2} N \log |\Psi| - \frac{1}{2} b' \Psi A - \frac{1}{2} N (\bar{x} - \mu)' \Psi (\bar{x} - \mu)$$

La matriz Ψ es definida positiva, por tanto $N(\bar{x} - \mu)' \Psi (\bar{x} - \mu) \geq 0$, $\forall (\bar{x} - \mu) \neq 0$, es decir que así nada cuando $\bar{x} - \mu = 0$, o sea $\mu = \bar{x}$. Para maximizar el logaritmo y tener término de la expresión anterior tendremos presente el siguiente resultado

LEMA.- Sea

$$f(A) = \frac{1}{2} N \log |A| - \frac{1}{2} \sum_{i,j=1}^p c_{ij} d_{ij}$$

donde $A = [c_{ij}]$ es definida positiva y donde $D = [d_{ij}]$ es definida positiva. Entonces el máximo de $f(A)$ se alcanza para $A = ND^{-1}$ y este máximo es

$$f(ND^{-1}) = \frac{1}{2} p N \log N - \frac{1}{2} N \log |D| - \frac{1}{2} p N$$

Demostración.- Obsérvese que $f(A)$ tiende a $-\infty$ cuando A se aproxima a una matriz singular. Igualmente se demuestra (véase Anderson pag. 47) que como lo mínimo cuando algún o algunos elementos de A se aproximan a ∞ y/o $-\infty$. Los máximos de $f(A)$ tendrán por tanto sus derivadas respecto a los elementos de A iguales a cero.

$$\frac{\partial f}{\partial c_{kk}} = \frac{1}{2} \frac{N}{|A|} \frac{\partial |A|}{\partial c_{kk}} - \frac{1}{2} d_{kk} = \frac{1}{2} N \frac{G_{kk}}{|A|} - \frac{1}{2} d_{kk}$$

donde G_{kk} es el cofactor de c_{kk} en A . Para $k \neq l$

$$\frac{\partial f}{\partial c_{kl}} = \frac{1}{2} \frac{N}{|A|} \frac{\partial |A|}{\partial c_{kl}} - \frac{1}{2} (d_{kl} + d_{lk})$$

y teniendo en cuenta la simetría de A y de D , tenemos

$$\frac{\partial f}{\partial c_{kl}} = N \frac{G_{kl}}{|A|} - d_{kl}$$

donde G_{kl} es el cofactor de c_{kl} en A . Igualando a cero estas parciales y teniendo en cuenta que $G_{kl}/|A|$ es el k -ésimo elemento de A^{-1} , obtenemos $NA^{-1} = D$, y de aquí $A = ND^{-1}$. El correspondiente valor máximo para $f(ND^{-1})$ es

$$\begin{aligned} f(ND^{-1}) &= \frac{1}{2} N \log |ND^{-1}| - \frac{1}{2} b' ND^{-1} D = \frac{1}{2} N \log N^p |D^{-1}| - \frac{1}{2} b' N I = \\ &= \frac{1}{2} p N \log N - \frac{1}{2} N \log |D| - \frac{1}{2} N p \end{aligned}$$

Aplicando este lema a la función $\log L$ podemos maximizar para $\Psi = NA^{-1} = \left(\frac{1}{N} A\right)^{-1}$. Supondremos que A no es singular para que tenga sentido lo dicho (no tiene verosimilitud que esto ocurra en probabilidad 1). En resumen los estimadores máximo verosimiles de μ y Ψ son $\hat{\mu} = \bar{x}$ y $\hat{\Psi} = NA^{-1}$. Para encontrar ahora el estimador de Σ necesitamos el siguiente lema y su recíproco.

LEMA.- Sea $f(\theta)$ una función real definida en un cierto conjunto S y sea ϕ una función real a valores simples, en una universo de iguales características, de S a algún otro conjunto S^* ; es decir, para cada $\theta \in S$ \exists un único $\theta^* \in S^*$ e inversamente para cada $\theta^* \in S^*$. Sea

$$g(\theta^*) = f[\phi^{-1}(\theta^*)]$$

Entonces si $f(\theta)$ alcanza un máximo en $\theta = \theta_0$, $g(\theta^*)$ alcanza un máximo en $\theta^* = \theta_0^* = \phi(\theta_0)$. Si el máximo de $f(\theta)$ es único en θ_0 , igual ocurre en el máximo de $g(\theta^*)$ en θ_0^* .

Demostración: Por hipótesis

$$f(\theta_0) \geq f(\theta) \quad \forall \theta \in S$$

Entonces para $\forall \theta^* \in S^*$, tenemos

$$g(\theta^*) = f[\phi^{-1}(\theta^*)] = f(\theta) \leq f(\theta_0) = g[\phi(\theta_0)] = g(\theta_0^*)$$

Resulta que $g(\theta^*)$ alcanza un máximo en θ_0^* . Si el máximo de $f(\theta)$ en θ_0 es único, la desigualdad anterior es estricta para $\theta \neq \theta_0$ y por tanto el máximo para $g(\theta^*)$ es único.

Tenemos el siguiente corolario.

COROLARIO - Si sobre la base de una muestra dada $\hat{\theta}_1, \dots, \hat{\theta}_m$ son los estimados máximo verosimiles de los parámetros $\theta_1, \dots, \theta_m$ de una distribución, entonces $\phi_1(\hat{\theta}_1, \dots, \hat{\theta}_m), \dots, \phi_m(\hat{\theta}_1, \dots, \hat{\theta}_m)$ son los estimados máximo verosimiles de $\phi_1(\theta_1, \dots, \theta_m), \dots, \phi_m(\theta_1, \dots, \theta_m)$ si la transformación de $\theta_1, \dots, \theta_m$ a ϕ_1, \dots, ϕ_m es uno a uno. Si los estimados de $\theta_1, \dots, \theta_m$ son únicos, lo son también en los ϕ_1, \dots, ϕ_m .

Aplicando el corolario a una muestra aleatoria podemos afirmar que el estimador máximo verosimil de Σ viene dado por $\hat{\Sigma} = \phi^{-1} = (1/N)A$. Resumiendo:

TEOREMA - Si x_1, \dots, x_N constituye una muestra de una $N(\mu, \Sigma)$ ($p < N$), los estimados máximo verosimiles de μ y Σ son $\hat{\mu} = \bar{x} = (1/N) \sum_{\alpha} x_{\alpha}$ y $\hat{\Sigma} = (1/N) \sum_{\alpha} (x_{\alpha} - \bar{x})(x_{\alpha} - \bar{x})'$.

COROLARIO - Si x_1, \dots, x_N constituye una muestra de una $N(\mu, \Sigma)$, donde $\sigma_{ij} = \sigma_i \sigma_j \rho_{ij}$ ($\rho_{ii} = 1$), el estimador máximo verosimil de μ es $\bar{x} = (1/N) \sum_{\alpha} x_{\alpha}$, la estimación máximo verosimil de $\sigma_i^2 = \hat{\sigma}_i^2 = (1/N) \sum_{\alpha} (x_{i\alpha} - \bar{x}_i)^2 = (1/N) \left(\sum_{\alpha} x_{i\alpha}^2 - N \bar{x}_i^2 \right)$, donde $x_{i\alpha}$ es la i -ésima componente de x_{α} y \bar{x}_i es la i -ésima componente de \bar{x} , el estimador máximo verosimil de ρ_{ij} es

$$\hat{\rho}_{ij} = \frac{\sum_{\alpha} (x_{i\alpha} - \bar{x}_i)(x_{j\alpha} - \bar{x}_j)}{\sqrt{\sum_{\alpha} (x_{i\alpha} - \bar{x}_i)^2 \cdot \sum_{\alpha} (x_{j\alpha} - \bar{x}_j)^2}}$$

Demostración - El conjunto de parámetros $\mu_i = \mu_i, \sigma_i^2 = \sigma_i^2$ y $\rho_{ij} = \frac{\sigma_{ij}}{\sqrt{\sigma_i^2 \sigma_j^2}}$ es una transformación uno a uno del conjunto de parámetros μ_i, σ_{ij} . Aplicando el corolario anterior los estimados máximo verosimiles son las correspondientes transformadas.

Vamos a finalizar este apartado de los estimados máximo verosimiles de μ, Σ dando una interpretación geométrica de algunos de los ~~elementos~~ elementos que aparecen en estos estimados.

Una interesante interpretación de la muestra desde un punto de vista geométrico, viene en terminos de las filas de la matriz de datos $X = (x_1, \dots, x_N)$, a saber

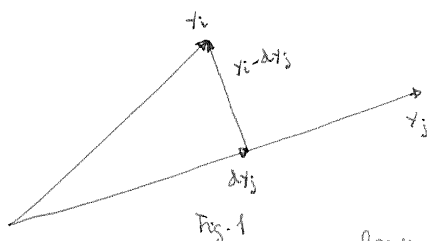


Fig. 1

$$X = \begin{bmatrix} x_{11} & \dots & x_{1N} \\ \vdots & & \vdots \\ x_{p1} & \dots & x_{pN} \end{bmatrix} = \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix}$$

El vector x_i puede ser considerado como un elemento en el espacio euclideo d -dimensional. Más concretamente como el extremo de un vector en R^d , cuyo origen coincide con el origen de coordenadas.

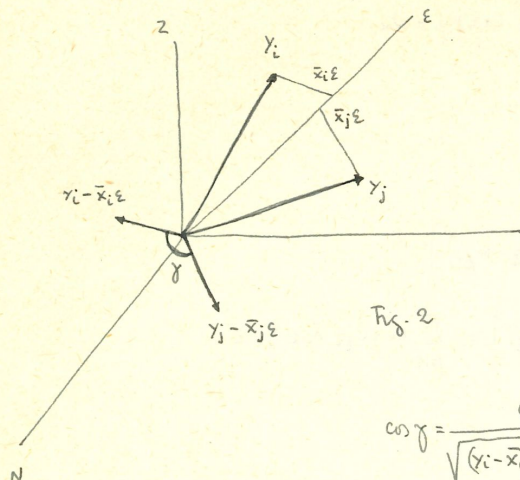
Comprobemos en primer lugar que el coseno del ángulo entre x_i e x_j es $x_i x_j' / \sqrt{x_i x_i' x_j x_j'}$.

En efecto, elegiremos un escalar d de tal forma que $dy_j \in x_i - dy_j$ sean ortogonales. Entonces

$dy_j (x_i - dy_j)' = 0$ y de aquí $d(x_i x_j' - d x_j x_j') = 0 \rightarrow d = x_i x_j' / x_j x_j'$. Si descomponemos ahora x_i en dy_j e $x_i - dy_j$ (en Fig. 1) (tenemos que el valor absoluto del seno del ángulo que forman x_i e x_j vendrá dado por el cociente entre la longitud de dy_j y la longitud de x_i , es decir

$$\sqrt{\frac{(dy_j)(dy_j)'}{x_i x_i'}} = \sqrt{\frac{d x_j x_j' d}{x_i x_i'}} = \sqrt{\frac{[x_i x_j'] [x_j x_j'] (x_i x_j')}{(x_j x_j') (x_i x_i') (x_j x_j')}} = \frac{x_i x_j'}{\sqrt{(x_i x_i') (x_j x_j')}}.$$

Podemos ahora dar una interpretación geométrica a a_{ij} y $a_{ij}/\sqrt{a_{ii}a_{jj}}$, donde a_{ij} son los elementos de la matriz A . Introduciremos la recta equiangular que es una línea cuya dirección es paralela a la



de vector $(1,1,1)$ y pasa por el origen. La proyección del vector x_i sobre la dirección que $E = (1,1,1)$ determina una línea dada por $\{(x_i E')/(E E')\} E = (\sum_{\alpha} x_{i\alpha} / \sum_{\alpha} 1) E = \bar{x}_i E = (\bar{x}_i, \dots, \bar{x}_i)$

Podemos ahora descomponer x_i en dos vectores, $\bar{x}_i E$ en la dirección de E y otro $x_i - \bar{x}_i E$ proyección de x_i sobre el plano perpendicular a E por el extremo de $\bar{x}_i E$. El cuadrado de la longitud de $x_i - \bar{x}_i E$ es $(x_i - \bar{x}_i E)(x_i - \bar{x}_i E)' = \sum_{\alpha} (x_{i\alpha} - \bar{x}_i)^2$

$= N \hat{\sigma}_{ii} = a_{ii}$. Si trasladamos ahora (fig. 2) los vectores $x_i - \bar{x}_i E$ y $x_j - \bar{x}_j E$ de manera que su origen coincida con el origen de coordenadas el coseno del ángulo entre ambos vectores vendrá determinado por

$$\cos \gamma = \frac{(x_i - \bar{x}_i E)(x_j - \bar{x}_j E)'}{\sqrt{(x_i - \bar{x}_i E)(x_i - \bar{x}_i E)' (x_j - \bar{x}_j E)(x_j - \bar{x}_j E)'}} = \frac{\sum_{\alpha} (x_{i\alpha} - \bar{x}_i)(x_{j\alpha} - \bar{x}_j)}{\sqrt{\sum_{\alpha} (x_{i\alpha} - \bar{x}_i)^2 \cdot \sum_{\alpha} (x_{j\alpha} - \bar{x}_j)^2}} = r_{ij}$$

6.- DISTRIBUCION DEL VECTOR MEDIA MUESTRAL. INFERENCIAS ACERCA DE μ CUANDO Σ ES CONOCIDA

la distribución conjunta de μ y Σ .

Recordemos que en el caso univariante la media muestral y la varianza muestral se distribuyen independientemente una de otra y además, la media, en particular, se distribuye normalmente, cuando la población de origen es normal. El resultado similar puede demostrarse para el caso multivariante, previamente a ello haremos al conjunto de las observaciones una transformación.

TEOREMA .- Supongamos que X_1, \dots, X_N son independientes, con X_{α} distribuyéndose $N(\mu_{\alpha}, \Sigma)$. Sea $C = (C_{\alpha\beta})$ una matriz ortogonal. Entonces $Y_{\alpha} = \sum_{\beta=1}^N C_{\alpha\beta} X_{\beta}$ se distribuye según una $N(\mu_{\alpha}, \Sigma)$, en $V_{\alpha} = \sum_{\beta=1}^N C_{\alpha\beta} \mu_{\beta}$ y Y_1, \dots, Y_N son independientes.

Demostración .- El conjunto de los vectores $\{Y_{\alpha}\}$ tiene una distribución conjunta normal puesto que el conjunto de todas sus componentes es un conjunto de combinaciones lineales de las componentes de $\{X_{\alpha}\}$ que tienen una distribución conjunta normal. De hecho, se ve que

$$f(x_1, x_2, \dots, x_N) = \frac{1}{(2\pi)^{\frac{NP}{2}} |\Sigma|^{\frac{N}{2}}} \exp \left\{ -\frac{1}{2} \sum_{\alpha} (x_{\alpha} - \mu_{\alpha})' \Sigma^{-1} (x_{\alpha} - \mu_{\alpha}) \right\}$$

que puede escribirse de la forma

$$f(x_{11}, x_{12}, \dots, x_{1N}, x_{21}, x_{22}, \dots, x_{2N}, \dots, x_{N1}, x_{N2}, \dots, x_{NN}) = \frac{1}{(2\pi)^{\frac{NP}{2}} |\Sigma|^{\frac{N}{2}}} \exp \left\{ -\frac{1}{2} ((x_1 - \mu_1)', \dots, (x_N - \mu_N)') \Sigma^{*-1} ((x_1 - \mu_1), \dots, (x_N - \mu_N))' \right\}$$

donde Σ^* es una matriz de submatrices de la forma

$$\Sigma^* = \begin{bmatrix} \Sigma & 0 & \dots & 0 \\ 0 & \Sigma & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \Sigma \end{bmatrix} \quad \text{de dimensión } N_p \times N_p$$

y cuyo determinante vale $|\Sigma^*| = |\Sigma|^N$.

El paso de los vectores X_{α} a los Y_{α} puede llevarse a cabo mediante una transformación del tipo

$$\begin{bmatrix} Y_{11} \\ Y_{12} \\ \vdots \\ Y_{N1} \\ \vdots \\ Y_{1p} \\ \vdots \\ Y_{Np} \end{bmatrix} = C^* \begin{bmatrix} X_{11} \\ X_{12} \\ \vdots \\ X_{N1} \\ \vdots \\ X_{1p} \\ \vdots \\ X_{Np} \end{bmatrix} \quad \text{con } C^* = \begin{bmatrix} C & 0 & \dots & 0 \\ 0 & C & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & C \end{bmatrix}$$

la esperanza de Y_α es

$$E(Y_\alpha) = E\left(\sum_{\beta} c_{\alpha\beta} X_{\beta}\right) = \sum_{\beta} c_{\alpha\beta} E(X_{\beta}) = \sum_{\beta} c_{\alpha\beta} \mu_{\beta} = \nu_{\alpha}.$$

la matriz de covarianzas entre Y_α e Y_β viene dada por

$$\begin{aligned} \text{Cov}(Y_\alpha, Y_\beta) &= E[(Y_\alpha - \nu_\alpha)(Y_\beta - \nu_\beta)] = \\ &= E\left[\left(\sum_{\gamma} c_{\alpha\gamma}(X_{\gamma} - \mu_{\gamma})\right)\left(\sum_{\delta} c_{\beta\delta}(X_{\delta} - \mu_{\delta})\right)\right] = \\ &= \sum_{\gamma, \delta} c_{\alpha\gamma} c_{\beta\delta} E[(X_{\gamma} - \mu_{\gamma})(X_{\delta} - \mu_{\delta})] = \sum_{\gamma, \delta} c_{\alpha\gamma} c_{\beta\delta} \delta_{\gamma\delta} \Sigma = \\ &= \sum_{\gamma} c_{\alpha\gamma} c_{\beta\gamma} \Sigma = \delta_{\alpha\beta} \Sigma \end{aligned}$$

donde $\delta_{\alpha\beta}$ es la delta de Kronecker ($= 1$ si $\alpha = \beta$ e $= 0$ si $\alpha \neq \beta$). Este resultado muestra que Y_α es independiente de Y_β y que Y_α tiene matriz de covarianzas Σ .

LEMA. - Si $G = [c_{\alpha\beta}]$ es ortogonal, entonces $\sum_{\alpha=1}^N X_{\alpha} X_{\alpha}' = \sum_{\alpha=1}^N Y_{\alpha} Y_{\alpha}'$, donde $Y_{\alpha} = \sum_{\beta} c_{\alpha\beta} X_{\beta}$.

Demostración.

$$\sum_{\alpha} Y_{\alpha} Y_{\alpha}' = \sum_{\alpha} \left(\sum_{\beta} c_{\alpha\beta} X_{\beta}\right) \left(\sum_{\gamma} c_{\alpha\gamma} X_{\gamma}'\right) = \sum_{\beta, \gamma} \left(\sum_{\alpha} c_{\alpha\beta} c_{\alpha\gamma}\right) X_{\beta} X_{\gamma}' = \sum_{\beta, \gamma} \delta_{\beta\gamma} X_{\beta} X_{\gamma}' = \sum_{\beta} X_{\beta} X_{\beta}'.$$

Sean ahora X_1, \dots, X_N independientes, cada uno de ellos distribuido de acuerdo con una $N(\mu, \Sigma)$. Entonces existe una matriz ortogonal $B = [b_{\alpha\beta}]$ cuya última fila es

$$(1/\sqrt{N}, \dots, 1/\sqrt{N}).$$

Esta transformación es una rotación en el espacio N dimensional en la que la línea equiangular descrita anteriormente se hace coincidir con el N -ésimo eje de coordenadas. Sea A la matriz de la suma de los cuadrados de las derivadas repetidas de la media muestral y los productos cruzados de las mismas, y sea $A = N \hat{\Sigma}$, y sea

$$Z_{\alpha} = \sum_{\beta} b_{\alpha\beta} X_{\beta}.$$

Entonces

$$Z_N = \sum_{\beta} b_{N\beta} X_{\beta} = \sum_{\beta} \frac{1}{\sqrt{N}} X_{\beta} = \sqrt{N} \bar{X}$$

Por el lema anterior tendremos:

$$A = \sum_{\alpha=1}^N Z_{\alpha} Z_{\alpha}' - N \bar{X} \bar{X}' = \sum_{\alpha=1}^N Z_{\alpha} Z_{\alpha}' - Z_N Z_N' = \sum_{\alpha=1}^{N-1} Z_{\alpha} Z_{\alpha}'$$

Puesto que Z_N es independiente de Z_1, \dots, Z_{N-1} , \bar{X} será independiente de A . Por otra parte

$$E[Z_N] = \sum_{\beta} b_{N\beta} E(X_{\beta}) = \sum_{\beta} \frac{1}{\sqrt{N}} \mu_{\beta} = \sum_{\beta} \frac{1}{\sqrt{N}} \mu = \sqrt{N} \mu,$$

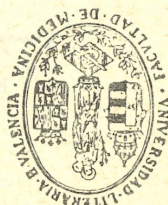
con lo que Z_N se distribuye $N(\sqrt{N} \mu, \Sigma)$ y por tanto $\bar{X} = \frac{1}{\sqrt{N}} Z_N$ será $N(\mu, \frac{1}{N} \Sigma)$.

Además

$$E(Z_{\alpha}) = \sum_{\beta} b_{\alpha\beta} E(X_{\beta}) = \sum_{\beta} b_{\alpha\beta} \mu = \sum_{\beta} b_{\alpha\beta} b_{N\beta} \sqrt{N} \mu = 0, \text{ para } \alpha \neq N.$$

Podemos reunir todos estos resultados en el siguiente teorema

EXAMENES



TEOREMA - La media muestral de una muestra de tamaño N de una población $N(\mu, \Sigma)$ se distribuye de acuerdo con una $N(\mu, (1/N)\Sigma)$ e independientemente de $\hat{\Sigma}$, estimador máximo verosímil de Σ . $N\hat{\Sigma}$ se distribuye igual que $\sum_{\alpha=1}^{N-1} Z_{\alpha} Z_{\alpha}'$ donde Z_{α} se distribuye $N(0, \Sigma)$ e independientemente de Z_{β} ($\alpha \neq \beta$).

Observemos por otra parte que

$$E(\hat{\Sigma}) = \frac{1}{N} \sum_{\alpha=1}^{N-1} E(Z_{\alpha} Z_{\alpha}') = \frac{N-1}{N} \Sigma$$

Resulta pues que $\hat{\Sigma}$ es un estimador sesgado de Σ . Lo podemos convertir en un estimador insesgado definiendo

$$S = \frac{1}{N-1} A = \frac{1}{N-1} \sum_{\alpha=1}^N (x_{\alpha} - \bar{x})(x_{\alpha} - \bar{x})'$$

matriz conocida como la matriz de varianzas muestral. Sus elementos ^{diagonales} son los estimadores insesgados habituales de las varianzas de las componentes de Σ .

Inferencia acerca de μ cuando Σ es conocida

Un problema estadístico de considerable importancia es el de contrastar la hipótesis de que el vector media de una distribución normal, con vector dado, y un problema directamente relacionado con ste el dar una región de confianza para dicho vector media. Vamos a estudiar estos problemas, bajo el supuesto de conocer la matriz de varianzas Σ . Más adelante consideraremos estos problemas para el caso en que suente desconocida Σ .

Recordemos que en el caso univariante y bajo planteamientos análogos a los anteriores, los tests intervalos de confianza para la media se basan en el hecho de que la diferencia entre la media muestral y la poblacional se distribuye normalmente en media 0 y varianza conocida. Un hecho similar utilizaremos en nuestro caso, también ahora esta diferencia se distribuye normalmente en vector media nulo y matriz de varianzas conocida. Podríamos elegir límites para cada componente sobre la base de esta distribución, pero este procedimiento tiene la desventaja de que la elección de límites es arbitraria y en el caso de tests ~~que inducen a~~ tests que pueden ser muy pobres frente a algunas alternativas, so con añadir la dificultad de calcular los referidos límites por cuanto sólo se dispone de fórmulas para la normal bivariente. El procedimiento que vamos a describir es sencillo en cuanto a la obtención de resultados y además pueden darse justificaciones teóricas e intuitivas generales.

El procedimiento se basa en el siguiente teorema:

TEOREMA - Si las m -componentes de un vector Y se distribuyen de acuerdo con una $N(0, T)$ (no singular), entonces $Y'T^{-1}Y$ se distribuye según una χ^2 con m grados de libertad.

Demostremos - Sea C una matriz no singular tal que $CTC' = I$ y definamos $Z = CY$. Entonces Z se distribuye normalmente con media 0 y matriz de varianzas $E(ZZ') = E(CYY'C) = CTC' = I$. Entonces $Y'T^{-1}Y = Z'(C')^{-1}T^{-1}C^{-1}Z = Z'(CTC')^{-1}Z = Z'Z$ que es la suma de los cuadrados de las componentes de Z . Como cada componente es $N(0, 1)$ e independiente de las restantes $Z'Z = Y'T^{-1}Y$ es una χ^2 con m grados de libertad.

Ahora bien, $\sqrt{N}(\bar{X} - \mu)$ se distribuye $N(0, \Sigma)$, aplicando el teorema

$$N(\bar{X} - \mu)' \Sigma^{-1} (\bar{X} - \mu)$$

será una χ^2 con p grados de libertad. Este resultado es fundamental para determinar regiones de confianza para μ .

Sea $\chi_p^2(\alpha)$, el número tal que

$$P(\chi_p^2 \geq \chi_p^2(\alpha)) = \alpha$$

Entonces

$$P(N(\bar{X} - \mu)' \Sigma^{-1} (\bar{X} - \mu) \geq \chi_p^2(\alpha)) = \alpha$$

Para contrastar la hipótesis $\mu = \mu_0$, utilizaremos como región crítica (o de rechazo)

$$N(\bar{x} - \mu_0)' \Sigma^{-1} (\bar{x} - \mu_0) \geq \chi_p^2(\alpha)$$

Intuitivamente podemos observar que la probabilidad de rechazar μ_0 sea mayor que α si μ es muy diferente de μ_0 , puesto que en el espacio de los \bar{x} , la expresión anterior representa una elipse centrada en μ_0 y cuando μ está alejado de μ_0 la densidad de \bar{x} se concentrará en un punto fuera de la citada elipse, lo que será difícil encontrar valores de \bar{x} cercanos a μ_0 .

Una demostración análoga a la del teorema anterior nos llevará a demostrar que $N(\bar{x} - \mu_0)' \Sigma^{-1} (\bar{x} - \mu_0)$ se distribuye como una χ^2 con p g.l. y parámetro de no centralización $N(\mu - \mu_0)' \Sigma^{-1} (\mu - \mu_0)$ cuando \bar{x} es la media de una muestra de tamaño N de una $N(\mu, \Sigma)$. La primera demostración del teorema anterior fue hecha por Pearson en 1900.

Consideremos ahora la siguiente afirmación hecha sobre la base de una muestra en media \bar{x} : "La media de la distribución satisface

$$N(\bar{x} - \mu^*)' \Sigma^{-1} (\bar{x} - \mu^*) \leq \chi_p^2(\alpha)$$

considerada como una desigualdad sobre μ^* ". De acuerdo con lo anterior la probabilidad de que extraigamos una muestra tal que sea cierta la afirmación es $1 - \alpha$, puesto que el suceso $\{N(\bar{x} - \mu^*)' \Sigma^{-1} (\bar{x} - \mu^*) \geq \chi_p^2(\alpha)\}$ es equivalente a negar la afirmación. Así, el conjunto de valores μ^* que satisfacen la última desigualdad están constituyen una región de confianza para μ con un nivel de confianza de $1 - \alpha$.

En un espacio p -dimensional de \bar{x} , $N(\bar{x} - \mu_0)' \Sigma^{-1} (\bar{x} - \mu_0) \geq \chi_p^2(\alpha)$ es la superficie y el exterior de un elipsoide centrado en μ_0 , la forma del elipsoide depende de Σ^{-1} y el tamaño de $(1/N) \chi_p^2(\alpha)$ para Σ^{-1} dado. En el espacio p -dimensional de μ^* , $N(\bar{x} - \mu^*)' \Sigma^{-1} (\bar{x} - \mu^*) \leq \chi_p^2(\alpha)$ es la superficie y el interior de un elipsoide centrado en \bar{x} . Si $\Sigma^{-1} = I$, la $\Pr\{N(\bar{x} - \mu)' \Sigma^{-1} (\bar{x} - \mu) \geq \chi_p^2(\alpha)\} = \alpha$, significa que la probabilidad de que la distancia entre \bar{x} y μ sea mayor que $[(1/N) \chi_p^2(\alpha)]^{1/2}$ es igual a α .

TEOREMA .- Si \bar{x} es la media de una muestra de N elementos extraída de una $N(\mu, \Sigma)$ y si Σ es conocida, entonces $N(\bar{x} - \mu_0)' \Sigma^{-1} (\bar{x} - \mu_0) \geq \chi_p^2(\alpha)$ es una región crítica de tamaño α para contrastar la hipótesis $\mu = \mu_0$ y $N(\bar{x} - \mu^*)' \Sigma^{-1} (\bar{x} - \mu^*) \leq \chi_p^2(\alpha)$ es una región de confianza para μ de nivel de confianza $1 - \alpha$. $\chi_p^2(\alpha)$ es el percentil $1 - \alpha$ de la distribución χ_p^2 .

Caso de dos medias

Una técnica similar puede ser utilizada para los problemas que surgen en dos muestras. Supongamos que tenemos una muestra $\{X_{\alpha}^{(1)}\}$ ($\alpha = 1, \dots, N_1$) de la distribución $N(\mu^{(1)}, \Sigma)$ y una muestra $\{X_{\alpha}^{(2)}\}$ ($\alpha = 1, \dots, N_2$) de una segunda población normal $N(\mu^{(2)}, \Sigma)$ ambas con la misma matriz de varianzas. Entonces las dos medias muestrales

$$\bar{x}^{(1)} = \frac{1}{N_1} \sum_{\alpha=1}^{N_1} X_{\alpha}^{(1)} \quad \bar{x}^{(2)} = \frac{1}{N_2} \sum_{\alpha=1}^{N_2} X_{\alpha}^{(2)}$$

se distribuirán independientemente $N(\mu^{(1)}, (1/N_1)\Sigma)$ y $N(\mu^{(2)}, (1/N_2)\Sigma)$ respectivamente. La diferencia entre ambas medias $y = \bar{x}^{(1)} - \bar{x}^{(2)}$ se distribuye $N[\nu, (1/N_1 + 1/N_2)\Sigma]$, donde $\nu = \mu^{(1)} - \mu^{(2)}$. Así

$$\frac{N_1 N_2}{N_1 + N_2} (y - \nu)' \Sigma^{-1} (y - \nu) \leq \chi_p^2(\alpha)$$

es una región de confianza para la diferencia ν entre las dos verdaderas medias y una región crítica para contrastar la hipótesis $\mu^{(1)} = \mu^{(2)}$ viene dada por

$$\frac{N_1 N_2}{N_1 + N_2} (\bar{x}^{(1)} - \bar{x}^{(2)})' \Sigma^{-1} (\bar{x}^{(1)} - \bar{x}^{(2)}) \geq \chi_p^2(\alpha).$$

Mahalanobis (1930) sugiere la utilización de $(\mu^{(1)} - \mu^{(2)})' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)})$ como una medida de la distancia entre las dos poblaciones.

Heur demostrado que

$$\sum_{\alpha} (x_{\alpha} - \mu)(x_{\alpha} - \mu)' = A + N(\bar{x} - \mu)(\bar{x} - \mu)' \quad \gamma$$

$$\sum_{\alpha} (x_{\alpha} - \mu)' \Sigma^{-1} (x_{\alpha} - \mu) = \text{tr}(\Sigma^{-1} A) + N(\bar{x} - \mu)' \Sigma^{-1} (\bar{x} - \mu)$$

Osi la densidad de X_1, \dots, X_N podemos escribir como

$$K \exp \left\{ -\frac{1}{2} \left[N(\bar{x} - \mu)' \Sigma^{-1} (\bar{x} - \mu) + \text{tr}(\Sigma^{-1} A) \right] \right\} = K_1 \exp \left\{ -\frac{1}{2} N(\bar{x} - \mu)' \Sigma^{-1} (\bar{x} - \mu) \right\} K_2 \exp \left\{ -\frac{1}{2} \text{tr}(\Sigma^{-1} A) \right\}$$

Por tanto \bar{x} y $(1/N)A$ forman un conjunto suficiente de estadísticas para μ y Σ . Si Σ es conocido, \bar{x} es un estadístico suficiente para μ . Sin embargo, si μ es conocido, $(1/N)A$ no es un estadístico suficiente para Σ , pero sí lo es $(1/N) \sum_{\alpha=1}^N (x_{\alpha} - \mu)(x_{\alpha} - \mu)'$.

Recordemos que t es un estadístico suficiente para θ si

$$\prod_{\alpha=1}^N f(x_{\alpha}; \theta) = g(t; \theta) \cdot h(x_1, \dots, x_N)$$

donde $f(x_{\alpha}; \theta)$ es la densidad de la α -ésima observación; $g(t; \theta)$ es una función de θ y $h(x_1, \dots, x_N)$ no depende de θ (criterio de factorización de Fisher-Neyman).

Si un vector aleatorio Y q -dimensional tiene media ν y matriz de varianzas Ψ , entonces

$$(Y - \nu)' \Psi^{-1} (Y - \nu) = q + 2$$

se denomina elipse de concentración de Y (véase Cramer, 344). La densidad definida mediante una distribución uniforme sobre el interior de este elipse tiene el mismo vector media y la misma matriz de varianzas que Y .

Por otra parte, si θ es un vector de q parámetros y t es un vector de estadísticas no sesgadas de θ basados en N observaciones de una población con matriz de varianzas Ψ , entonces el elipse

$$N(t - \theta)' E \left[\frac{\partial \log f}{\partial \theta} \right] \left[\frac{\partial \log f}{\partial \theta} \right]' (t - \theta) = q + 2$$

está enteramente contenido en el elipse de concentración de t (véase Cramer p. 566). Si este último elipse ~~coincide~~ coincide en el de concentración de t , entonces se dice que t es eficiente. En general el cuadrado de la raíz del volumen de este último elipse y el del de concentración se define como la eficiencia relativa de t . En el caso de una distribución normal multivariante, si $\theta = \mu$, el vector de las medias muestrales \bar{x} es eficiente. Si θ incluye a μ y a Σ , entonces \bar{x} y S tienen una eficiencia relativa igual a $[(N-1)/N]^{p(p+1)/2}$.

En efecto

$$\log f = \log \left[\frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \right] - \frac{1}{2} (x - \mu)' \Sigma^{-1} (x - \mu)$$

$$\begin{aligned} \frac{\partial \log f}{\partial \mu} &= -\frac{1}{2} [-2 \Sigma^{-1} (x - \mu)] \quad \gamma \text{ de aquí } E \left[\frac{\partial \log f}{\partial \theta} \right] \left[\frac{\partial \log f}{\partial \theta} \right]' = E \left[(\Sigma^{-1} (x - \mu)) (\Sigma^{-1} (x - \mu))' \right] = \\ &= \Sigma^{-1} E[(x - \mu)(x - \mu)'] \Sigma^{-1} = \Sigma^{-1} \end{aligned}$$

entonces el elipse correspondiente es

$$N(\bar{x} - \mu)' \Sigma^{-1} (\bar{x} - \mu) = q + 2$$

pero $\bar{x} \sim N(\mu, \frac{1}{N} \Sigma)$, entonces el elipse de concentración viene dado por

$$N(\bar{x} - \mu)' \frac{\Sigma^{-1}}{N} (\bar{x} - \mu) = N(\bar{x} - \mu)' \Sigma^{-1} (\bar{x} - \mu) = q + 2. \text{ Ambos coinciden.}$$

Con frecuencia posemos experimentos con matrices de datos en las que al faltar algunas observaciones pone en entredicho la utilización de las técnicas usuales para la obtención de las correspondientes estimaciones de μ y Σ . Situaciones de este tipo pueden presentarse por ejemplo, cuando algunos de los animales en estudio en una experimentación de laboratorio fallece por causas ajenas a la experimentación, cuando en grandes investigaciones interdisciplinarias surgen nuevas variables a medida que progresa el estudio. En cualquier caso, es esencial que las causas que originan las pérdidas de datos sean completamente independientes de la naturaleza o de los valores de las variables sujeta a examinación. Si el número de lecturas de observaciones completas es pequeño y el costo de la obtención de datos es apreciable, preferiremos utilizar métodos de estimación y contrastes de hipótesis que optimicen el uso de la información disponible.

Las estimaciones máximo-verosímiles del vector media y de la matriz de covarianzas de una distribución multivariante mediante una matriz de datos incompleta con un modelo general aleatorio conduce a sistemas de ecuaciones no lineales cuya solución requiere técnicas de análisis numérico. Pueden consultarse en este caso los trabajos de Wilks (1932), Elashoff y Alfifi (1966), Tinn (1970), Mantley and Hocking (1971) y Basilek (1972).

Un caso especial de pérdida de datos permitirá el cálculo directo de las estimaciones máximo-verosímiles de μ y Σ . Se trata del caso en que la matriz de datos obedece al modelo conocido con el nombre de "monótono" o "anidado" y que es de la forma

$$X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1,k-1} & x_{1,k} \\ x_{21} & x_{22} & \dots & x_{2,k-1} & - \\ \dots & \dots & \dots & \dots & \dots \\ x_{k-1,1} & x_{k-1,2} & \dots & - & - \\ x_{k,1} & - & \dots & - & - \end{bmatrix}$$

en la que los guiones indican bloque de observaciones perdidas. x_{ij} es una submatriz $N_i \times r_j$, donde $\sum_{i=1}^k N_i = N$, número total de unidades muestrales independientes y $\sum_{j=1}^k r_j = p$, el número de repuestas. Consecuentemente el vector media aparece fraccionado de la forma $\mu' = [\mu'_1, \mu'_2, \dots, \mu'_k]$ y la matriz de varianzas aparece también fraccionada en matrices Σ_{ij} de dimensiones $r_i \times r_j$, $i, j = 1, \dots, k$. Como consecuencia de la monotonía del modelo de la matriz de datos posemos escribir la similitud de la muestra como

$$L(\mu, \Sigma) = f(X) = \left[\prod_{i=1}^k f(x_{i1}) \right] \cdot \left[\prod_{i=1}^{k-1} f(x_{i2}/x_{i1}) \right] \cdot \left[\prod_{i=1}^{k-2} f(x_{i3}/x_{i1}, x_{i2}) \right] \dots f(x_{k1}/x_{11}, \dots, x_{k-1,1})$$

en la que la parte derecha de la igualdad aparece enteramente de densidades condicionales a fin de evitar una notación excesivamente complicada. $f(x_{ij}/x_{i1}, \dots, x_{i,j-1})$ es la verosimilitud condicional de la i -ésima submatriz, cuya expresión puede obtenerse de la forma usual mediante el producto de N_i densidades multivariantes r_i -dimensionales e independientes cuyos parámetros se obtienen aplicando las expresiones correspondientes a las densidades condicionales. Se procede mediante la maximización sucesiva de los factores que aparecen en los corchetes: la primera de estas operaciones conduce a las estimaciones usuales $\hat{\mu}_1$ y $\hat{\Sigma}_{11}$ utilizando para ello las N unidades muestrales. El segundo factor proporciona estimaciones de $\mu_2 - \Sigma_{12}' \Sigma_{11}^{-1} \mu_1$, de los parámetros de regresión $\Sigma_{12}' \Sigma_{11}^{-1}$ y de la matriz de varianzas condicional $\Sigma_{22} - \Sigma_{12}' \Sigma_{11}^{-1} \Sigma_{12}$. Con la ayuda de las estimaciones previas de μ_1 y Σ_{11} podemos obtener, a partir de los anteriores, las estimaciones máximo-verosímiles de μ_2 , Σ_{12} y Σ_{22} . Continuamos de esta forma hasta que hayamos sido estimados los parámetros condicionales del k -ésimo conjunto. Este método de aproximación fue desarrollado por Anderson (1957) y fue extendido a más de dos conjuntos de variables por Bhargava (1962).

A continuación desarrollamos en detalle el método para el caso más sencillo, aquel en el que el modelo está formado tan sólo por dos conjuntos de variables. Sean $X' \equiv (x'_{11}, x'_{21})$ y $X \equiv x_{12}$. Entonces

$$\hat{\mu}_1 = \frac{1}{N} \sum_{i=1}^N x_i, \quad \hat{\Sigma}_{11} = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})'$$

estimados en la forma usual a partir del conjunto completo de datos. Para las otras estimaciones, tenemos

$$\bar{x} = \frac{1}{N_1} \sum_{i=1}^{N_1} x_i \quad \bar{y}_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} y_i \quad \bar{y}_2 = \frac{1}{N_2} \sum_{i=N_1+1}^N y_i$$

$$A_{11}(N_1) = \sum_{i=1}^{N_1} (x_i - \bar{x})(x_i - \bar{x})' \quad A_{11}(N) = \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})'$$

$$A_{22} = \sum_{i=1}^{N_1} (x_i - \bar{x})(x_i - \bar{x})' \quad A_{12} = \sum_{i=1}^{N_1} (y_i - \bar{y}_1)(x_i - \bar{x})'$$

$$B = A_{11}^{-1}(N_1) A_{12}$$

Las estimaciones de los parámetros restantes son

$$\hat{\mu}_2 = \bar{x} - \left(\frac{N_2}{N} \right) B' (\bar{y}_1 - \bar{y}_2)$$

$$\hat{\Sigma}_{22} = \frac{1}{N_1} [A_{22} - A_{12}' A_{11}^{-1}(N_1) A_{12}] + \frac{1}{N} B' A_{11}(N) B$$

$$\hat{\Sigma}_{12} = \frac{1}{N} A_{11}(N) A_{11}^{-1}(N_1) A_{12}$$

Las esperanzas y algunos momentos de segundo orden de estas estimaciones han sido calculados por Morrison (1971). En particular $\hat{\mu}_2$ y $[N/(N-1)] \hat{\Sigma}_{12}$ son insesgados, mientras que es posible redefinir los divisores de los distintos términos de $\hat{\Sigma}_{22}$ para obtener el siguiente estimador insesgado:

$$\hat{\Sigma}_{22}(x) = \frac{1}{N_1 - r_1 - 1} \left[1 - \frac{r_1}{N_1 - r_1 - 2} + \frac{r_1(r_1 + 1)}{(N_1 - r_1 - 2)(N - 1)} \right] [A_{22} - A_{12}' A_{11}^{-1}(N_1) A_{12}] + \frac{1}{N - 1} B' A_{11}(N) B$$

Los valores de las varianzas de las estimaciones de μ_2 , Σ_{22} y Σ_{12} para N_1 y N pequeños, indican que estas estimaciones son menos eficientes que las convencionales, obtenidas a partir, solamente, de las primeras N_1 observaciones completas cuando las correlaciones entre los dos conjuntos de variables son bajas. Antes para la elección del tipo de estimación más conveniente han sido ~~indicadas~~ indicadas, en la referencia anterior de Morrison (1971).

Ejemplo.- Una universidad utiliza una ecuación de regresión múltiple para estimar los índices puntuales de acceso (IPA) de los solicitantes a partir de sus rangos de clase, tests de aptitud scholastica (TAS) y otros tipos de tests de interés general. Las estimaciones de los IPA se incluyen en los formularios resumidos de los solicitantes para uso de los miembros del comité de admisión. Una muestra de $N=34$ solicitudes evaluadas por un miembro de la Facultad contiene cinco de ellas sin el correspondiente IPA. Sin embargo, todas las solicitudes tienen TAS de matemáticas y resultados de tests verbales entre otras medidas de tipo académico. Los datos así presentados son un ~~ejemplo~~ ejemplo de modelo monotónico de observaciones incompletas, con $r_1=2$ respuestas completas y $r_2=1$ incompletas, $N_1=29$ unidades muestrales completas y $N_2=5$ incompletas. Supongamos que las puntuaciones de los tests verbales (y_1), del TAS de matemáticas (y_2) y del IPA (x) son una muestra trivariante de una distribución normal cuyos parámetros sean estimados mediante los métodos que acabamos de exponer.

Comenzaremos preparando una representación gráfica informal de los valores del IPA frente a los resultados del TAS para los 29 vectores completos. Esta representación indica una ligera correlación positiva y parece que, en estas condiciones, las estimaciones basadas en observaciones incompletas hayan de ser más eficientes que las estimaciones basadas en los 29 resultados de los IPA. Continuemos pues

$$\bar{x} = 2.43 \quad \bar{y}_1 = \begin{bmatrix} 56.38 \\ 61.31 \end{bmatrix} \quad \bar{y}_2 = \begin{bmatrix} 55.80 \\ 55.60 \end{bmatrix} \quad \bar{y} = \begin{bmatrix} 56.29 \\ 60.47 \end{bmatrix}$$

$$A_{11}(N_1) = \begin{bmatrix} 1292.83 & 206.59 \\ 206.59 & 1328.21 \end{bmatrix} \quad A_{11}(N) = \begin{bmatrix} 1409.06 & 238.29 \\ 238.29 & 1668.47 \end{bmatrix} \quad a_{22} = 2.9455$$

$$A_{11}^{-1}(N_1) = 10^{-4} \begin{bmatrix} 7.9321 & -1.2338 \\ -1.2338 & 7.9208 \end{bmatrix} \quad B = 10^{-2} \begin{bmatrix} 1.6862 \\ 2.9142 \end{bmatrix}$$

y las estimaciones máximo-verosímiles de los parámetros de IPA con

2 (15)

$$\hat{\mu}_x = 2.110 \quad \hat{\sigma}_x^2 = 0.1034$$

La estimación de la media es ligeramente menor que la estimación usual basada en datos completos, \bar{x} . Para poder establecer comparaciones con σ_x^2 a continuación presentamos distintas estimaciones alternativas de la varianza:

Estimación	valor
MLE univariante sesgada σ^2/n	0.1016
MLE no sesgada datos incompletos	0.1067
MLE univariante no sesgada $\sigma^2/(n-1)$	0.1082

La estimación basada en los datos incompletos, tanto en la sesgada como la no sesgada, es ligeramente mayor que su correspondiente contraparte univariante (basada en los datos completos).

Capítulo 3

Coeficientes de correlación muestrales

1. INTRODUCCION

Se pretende en este capítulo el estudio de los coeficientes de correlación muestrales que no más que las cantidades equivalentes a los coeficientes de correlación poblacionales definidos en el capítulo anterior. Consideraremos sus distribuciones de probabilidad y construiremos tablas de hipótesis y regiones de rechazo para determinadas situaciones de interés.

En el caso de distribuciones mixtas conjuntas, los coeficientes de correlación son la medida natural de la dependencia entre variables. Para las estimaciones de los correlaciones poblacionales utilizaremos los correspondientes muestrales y justificaremos el porqué.

Comencemos recordando aquí la diferencia que existe entre los conceptos de regresión y correlación. En la teoría de la regresión una variable se considera aleatoria o dependiente de otras fijas o independientes. En la teoría de la correlación consideramos varias variables como aleatorias, las tratamos todas ellas simétricamente. Si consideramos una distribución conjunta normal y mantenemos fijas todas las variables menos una, obtenemos el modelo mínimo cuadrático puesto que la esperanza de la variable aleatoria en la distribución condicional es una función lineal de las variables que permanecen fijas. Los coeficientes de regresión muestrales obtenidos mediante el método de los mínimos cuadrados son funciones de las varianzas muestrales y de las correlaciones.

Si contestamos la independencia llegaremos a resultados análogos tanto si partimos de un caso como de otro (es decir; distribución conjunta normal o distribución condicional). La distribución de probabilidad bajo la hipótesis nula es la misma. La distribución de las estadísticas del test cuando la hipótesis nula no es cierta difiere en los dos casos. Si todas las variables pueden considerarse aleatorias utilizaremos la teoría de la correlación tal como aparece a continuación, si sólo una variable es aleatoria, utilizaremos la teoría de los mínimos cuadrados (posteriormente en otro capítulo).

2. COEFICIENTE DE CORRELACION DE UNA MUESTRA BIVARIANTE

Vimos en el capítulo anterior que una estimación máximo-verosímil del coeficiente de correlación entre los componentes X_i y X_j del vector aleatorio \mathbf{X} , cuando utilizamos una muestra de tamaño N (X_1, \dots, X_N) es

$$r_{ij} = \frac{\sum_{\alpha=1}^N (X_{i\alpha} - \bar{X}_i)(X_{j\alpha} - \bar{X}_j)}{\sqrt{\sum_{\alpha=1}^N (X_{i\alpha} - \bar{X}_i)^2} \sqrt{\sum_{\alpha=1}^N (X_{j\alpha} - \bar{X}_j)^2}}$$

Vamos a determinar la distribución de r_{ij} cuando la correlación poblacional entre X_i y X_j es uno y vamos como utilizar el coeficiente de correlación muestral para contrastar que el poblacional es cero. Desarrollaremos la teoría, por comodidad, para r_{12} . Como r_{12} depende sólo de las dos primeras componentes de cada X_α , bastará considerar solamente la distribución conjunta de $(X_{1\alpha}, X_{2\alpha})$, $\alpha = 1, \dots, N$. Podemos entonces reformular el problema anterior de una distribución normal bivalente. Sean X_1^*, \dots, X_N^* observaciones repetidas de

$$N \left[\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{bmatrix} \sigma_1^2 & \sigma_{12} \rho \\ \sigma_{12} \rho & \sigma_2^2 \end{bmatrix} \right]$$

Sea

$$r = r_{12} = \frac{a_{12}}{\sqrt{a_{11}} \sqrt{a_{22}}}$$

$$\text{con } a_{ij} = \sum_{\alpha=1}^N (X_{i\alpha} - \bar{X}_i)(X_{j\alpha} - \bar{X}_j) \quad i, j = 1, 2$$

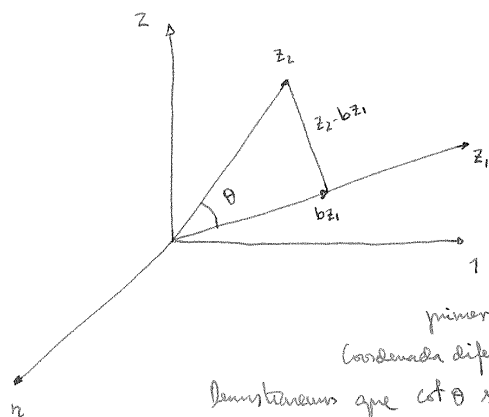
(viendo $X_{i\alpha}$ la i -ésima componente de X_α^*).

Sabemos que a_{ij} se distribuye como $\sum_{\alpha=1}^n Z_{i\alpha} Z_{j\alpha}$ $i, j = 1, 2$, donde $n = N - 1$ y $(Z_{1\alpha}, Z_{2\alpha})$ se distribuye de acuerdo con una normal bivalente

$$N \left[\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{bmatrix} \sigma_1^2 & \sigma_{12} \rho \\ \sigma_{12} \rho & \sigma_2^2 \end{bmatrix} \right]$$

e independientemente de $(Z_{1\beta}, Z_{2\beta})$, $\alpha \neq \beta$.

Sea $Z_i' = (z_{i1}, \dots, z_{in})$, $i=1,2$. Estos dos vectores pueden representarse en un espacio n -dimensional.



El coeficiente de correlación es el coseno del ángulo θ entre Z_1 y Z_2 .

Para encontrar la distribución de $\cos \theta$ buscaremos primero la de la $\cot \theta$. Como $Z_2 = (Z_2 - bZ_1) + bZ_1$, podemos b como una función de Z_1 y Z_2 tal que $Z_2 - bZ_1$ es ortogonal a bZ_1 . Entonces

$$\cot \theta = b \sqrt{Z_1' Z_1 / (Z_2 - bZ_1)' (Z_2 - bZ_1)}$$

Si Z_1 permanece fijo podemos girar los ejes de manera que el primer eje de coordenadas coincida con Z_1 . Entonces bZ_1 tiene solo la primera coordenada diferente de cero y $Z_2 - bZ_1$ tendrá la primera coordenada igual a cero.

Demostremos que $\cot \theta$ es proporcional a una variable t cuando $\rho = 0$.

La distribución condicional de $Z_{2\alpha}$ dado $Z_{1\alpha} = z_{1\alpha}$ es $N(\beta z_{1\alpha}, \sigma^2)$, donde $\beta = \rho \sigma_2 / \sigma_1$ y $\sigma^2 = \sigma_2^2 (1 - \rho^2)$ como vimos en el momento. La distribución conjunta de Z_2 dado $Z_1 = z_1$ es $N(\beta z_1, \sigma^2 I)$ puesto que las $Z_{2\alpha}$ son independientes. Ahora bien, la densidad conjunta de Z_1 y Z_2 es

$$\prod_{\alpha=1}^n n \left[\begin{pmatrix} z_{1\alpha} \\ z_{2\alpha} \end{pmatrix} \middle| \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{bmatrix} \sigma_1^2 & \sigma_1 \sigma_2 \rho \\ \sigma_1 \sigma_2 \rho & \sigma_2^2 \end{bmatrix} \right]$$

y la marginal de Z_1 viene dada por

$$\prod_{\alpha=1}^n n(z_{1\alpha} | 0, \sigma_1^2) = n(Z_1 | 0, \sigma_1^2 I)$$

por tanto la condicional de Z_2 dado $Z_1 = z_1$ más el cociente de la conjunta por la marginal de Z_1 en z_1 , da como

$$\prod_{\alpha=1}^n \left\{ n \left[\begin{pmatrix} z_{1\alpha} \\ z_{2\alpha} \end{pmatrix} \middle| \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{bmatrix} \sigma_1^2 & \sigma_1 \sigma_2 \rho \\ \sigma_1 \sigma_2 \rho & \sigma_2^2 \end{bmatrix} \right] / n(z_{1\alpha} | 0, \sigma_1^2) \right\} = \prod_{\alpha=1}^n n(z_{2\alpha} | \beta z_{1\alpha}, \sigma^2)$$

Sea ahora $b = Z_2' Z_1 / Z_1' Z_1 = (a_{21}/a_{11})$, entonces $bZ_1' (Z_2 - bZ_1) = 0$, y sea $V = (Z_2 - bZ_1)' (Z_2 - bZ_1) = Z_2' Z_2 - b^2 Z_1' Z_1 (= a_{22} - a_{21}^2/a_{11})$. Entonces $\cot \theta = b \sqrt{a_{11}/V}$. La rotación de los ejes coordenados antes descrita impone elegir una matriz $(n \times n)$ ortogonal, C , cuya primera fila sea $(1/c) Z_1'$ donde $c^2 = Z_1' Z_1$.

Podemos ahora aplicar un teorema del capítulo anterior que hacía referencia al muestreo por una matriz ortogonal. Si hacemos $X_\alpha = Z_{2\alpha}$ y $Y_\alpha = \sum_{\beta=1}^n c_{\alpha\beta} Z_{1\beta}$. Entonces los $\{Y_\alpha\}$ son independientes y se distribuyen normalmente con varianza σ^2 y media

$$E(Y_1) = \sum_{\beta=1}^n c_{1\beta} E(Z_{1\beta}) = \sum_{\beta=1}^n c_{1\beta} \beta z_{1\beta} = \frac{\beta}{c} \sum_{\beta=1}^n Z_{1\beta}^2 = \beta c$$

$$E(Y_\alpha) = \sum_{\beta=1}^n c_{\alpha\beta} E(Z_{1\beta}) = \sum_{\beta=1}^n c_{\alpha\beta} \beta z_{1\beta} = \beta c \sum_{\beta=1}^n c_{\alpha\beta} c_{1\beta} = 0, \quad \alpha \neq 1.$$

Tomemos $b = \sum_{\alpha} Z_{2\alpha} Z_{1\alpha} / \sum_{\alpha} Z_{1\alpha}^2 = c \sum_{\alpha} Z_{2\alpha} c_{1\alpha} / \sum_{\alpha} Z_{1\alpha}^2 = c Y_1 / c^2 = Y_1 / c$ y del lema que presentábamos a continuación del teorema antes utilizado podemos deducir

$$V = \sum_{\alpha} Z_{2\alpha}^2 - b^2 \sum_{\alpha} Z_{1\alpha}^2 = \sum_{\alpha} Y_\alpha^2 - Y_1^2 = \sum_{\alpha=2}^n Y_\alpha^2$$

que es independiente de b .

LEMA. Si $(Z_{1\alpha}, Z_{2\alpha})$, $\alpha=1, \dots, n$, son independientes, cada uno con la distribución normal conjunta de media $\mu' = [0, 0]$ y $\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_1 \sigma_2 \rho \\ \sigma_1 \sigma_2 \rho & \sigma_2^2 \end{bmatrix}$, entonces la distribución condicional de $b = \sum_{\alpha} Z_{2\alpha} Z_{1\alpha} / \sum_{\alpha} Z_{1\alpha}^2$ y $V/\sigma^2 = \sum_{\alpha} (Z_{2\alpha} - b Z_{1\alpha})^2 / \sigma^2$ dado $Z_{1\alpha} = z_{1\alpha}$ ($\alpha=1, \dots, n$) es $N(\beta, \sigma^2/c^2)$, $c^2 = \sum_{\alpha} Z_{1\alpha}^2$ y χ^2 con $n-1$ g.l. respectivamente y además b y V son independientes.

Si $\rho=0$, entonces $\beta=0$, por tanto b redistribuye incondicionalmente $N(0, \sigma^2/n)$ y

3 (2)

$$\frac{cb/\sigma}{\sqrt{\frac{V/\sigma^2}{n-1}}} = \frac{cb}{\sqrt{\frac{V}{n-1}}}$$

tiene una distribución condicional t con $n-1$ g.l. Ahora bien, esta variable aleatoria t

$$\sqrt{n-1} \cdot \frac{(a_{11})^{1/2} a_{12}/a_{11}}{\sqrt{a_{22} - a_{12}^2/a_{11}}} = \sqrt{n-1} \cdot \frac{a_{12}/\sqrt{a_{11}a_{22}}}{\sqrt{1 - [a_{12}^2/(a_{11}a_{22})]}} = \sqrt{n-1} \cdot \frac{r}{\sqrt{1-r^2}}$$

Así pues $\sqrt{n-1} r/\sqrt{1-r^2}$ tiene una distribución condicional t con $n-1$ g.l. La densidad de t es

$$\frac{\Gamma(\frac{1}{2}n)}{\sqrt{n-1} \Gamma(\frac{1}{2}(n-1)) \sqrt{\pi}} \left(1 + \frac{t^2}{n-1}\right)^{-\frac{1}{2}n}$$

y la densidad de $r/\sqrt{1-r^2}$ viene dada por

$$\frac{\Gamma(\frac{1}{2}n)}{\Gamma(\frac{1}{2}(n-1)) \sqrt{\pi}} (1+w^2)^{-\frac{1}{2}n}$$

Puesto que $w = r(1-r^2)^{-\frac{1}{2}}$ y $dw/dr = (1-r^2)^{-\frac{3}{2}}$, la densidad de r vendrá dada por (haciendo $n=N-1$)

$$\frac{\Gamma(\frac{1}{2}(N-1))}{\Gamma(\frac{1}{2}(N-2)) \sqrt{\pi}} (1-r^2)^{\frac{1}{2}(N-4)}$$

No olvidemos que esta es la densidad de r para z_1 fijo. Pero dado que no depende para nada de z_1 , será también la marginal de r .

TEOREMA .- Sean X_1, \dots, X_N independientes, cada una con distribución $N(\mu_i, \Sigma)$. Si $\rho_{ij}=0$, la densidad de r_{ij} ~~es la misma que la de r~~ definido como ante, viene dada por la última expresión anterior.

Observamos que la densidad de r es simétrica respecto del origen. Además, para $N \geq 4$, tiene una moda en $r=0$ y presenta valles en los puntos ± 1 cuyo orden es $\frac{1}{2}(N-5)$ para N impar y $\frac{1}{2}N-3$ para N par. Puesto que la densidad es par los momentos impares son nulos, en particular la media es cero. Los momentos pares se pueden encontrar por integración (haciendo el cambio $x=r^2$ y utilizando la función Γ), se obtiene $E[r^{2m}]$ mediante la expresión $E[r^{2m}] = \frac{\Gamma(\frac{1}{2}(N-1)) \Gamma(\frac{1}{2}+m)}{\Gamma(\frac{1}{2}(N-1)+m) \sqrt{\pi}}$, en particular la variancia ($m=1$) es $1/(N-1)$.

El uso más importante del anterior teorema es buscar puntos de significación para contrastar la hipótesis de que un par de variables están vinculadas. Consideremos la hipótesis

$$H_0: \rho_{ij}=0$$

para algún par (i,j) . Parece razonable rechazar H_0 si r_{ij} difiere mucho de cero, pero cómo cuantificar este "mucha".

Supongamos que vamos a contrastar H_0 frente a $H_1: \rho_{ij} > 0$. Rechazaremos H_0 si r_{ij} es mayor que una determinada cantidad r_0 . La probabilidad de rechazar H_0 siendo cierta es

$$\int_{r_0}^1 K_N(r) dr$$

donde $K_N(r)$ es la densidad del coeficiente de correlación muestral basado en N observaciones. Elegimos r_0 de manera que la anterior probabilidad coincida con el nivel de significación deseado. Si la alternativa es de la forma $\rho_{ij} < 0$, entonces rechazaremos para $r_{ij} < -r_0$.

Para el caso de una alternativa bilateral, $\rho_{ij} \neq 0$. Entonces rechazamos H_0 si $r_{ij} > r_1$ o $r_{ij} < -r_1$.
 La probabilidad de rechazo para H_0 , siendo cierta, vendrá dada por

$$\int_{-r_1}^{-r_1} K_N(r) dr + \int_{r_1}^1 K_N(r) dr.$$

Existen tablas (Fisher y Yates, 1942) para la distribución de r , de cualquier forma, teniendo en cuenta que $\sqrt{N-2} r / \sqrt{1-r^2}$ tiene una distribución t en $N-2$ g.l. podemos usar alternativamente esta distribución para la obtención de r_0 y r_1 . Así, para alternativas del tipo $\rho_{ij} \neq 0$, rechazamos H_0 si

$$\sqrt{N-2} \frac{|r_{ij}|}{\sqrt{1-r_{ij}^2}} > t_{N-2}(\alpha)$$

de forma análoga actuaremos en los otros casos.

Por otra parte, en el transcurso de la deducción del test anterior queda claro que el mismo estadístico puede ser utilizado para contrastar la hipótesis de que la regresión de Z_2 sobre Z_1 es cero. Estimemos de los observados originales, tenemos

$$\sqrt{N-2} \frac{r}{\sqrt{1-r^2}} = \frac{b \sqrt{\sum_{\alpha} (x_{1\alpha} - \bar{x}_1)^2}}{\sqrt{\sum_{\alpha} [x_{2\alpha} - \bar{x}_2 - b(x_{1\alpha} - \bar{x}_1)]^2 / (N-2)}}$$

donde $b = \frac{\sum_{\alpha} (x_{2\alpha} - \bar{x}_2)(x_{1\alpha} - \bar{x}_1)}{\sum_{\alpha} (x_{1\alpha} - \bar{x}_1)^2}$ es el coeficiente de regresión mínimo cuadrático de $x_{2\alpha}$ sobre $x_{1\alpha}$. Ya hemos visto que el test $\rho_{12} = 0$ es equivalente al test de que la regresión de Z_2 sobre Z_1 es nula (es decir, que $\rho_{12}/\sigma_1 = 0$).

Distribución cuando $\rho \neq 0$. Test de hipótesis y regiones de confianza.

En este caso deberemos en primer lugar obtener la distribución conjunta de a_{11}, a_{12} y a_{22} . Lo haremos antes que las condicionales, para Z_1 fijo, de $b = a_{12}/a_{11}$ y $v/\sigma^2 = (a_{22} - a_{12}^2/a_{11})/\sigma^2$ son independientes y más concretamente $N(\beta, \sigma^2/c^2)$ y χ^2 en $n-1$ g.l. respectivamente. Denotando la densidad de la χ^2_{n-1} mediante $g_{n-1}(v)$, la conjunta correspondiente de b y v será

$$n(b|\beta^2, \sigma^2/a_{11}) g_{n-1}(v/\sigma^2)/\sigma^2$$

la conjunta de Z_1, b y v vendrá dada por

$$n(z_1|0, \sigma_1^2 I) \cdot n(b|\beta^2, \sigma^2/a_{11}) \cdot g_{n-1}(v/\sigma^2)/\sigma^2.$$

la densidad marginal de $Z_1, Z_1/\sigma_1^2 = a_{11}/\sigma^2$, se puede obtener de mediante

$$\int_{z_1' z_1 = a_{11}} n(z_1|0, \sigma_1^2 I) dW$$

y como resultado de la suma de cuadrados de n normales $(0,1)$, tendremos que más una χ^2 con n g.l., es decir

$$\frac{1}{\sigma_1^2} g_n\left(\frac{a_{11}}{\sigma_1^2}\right) = \int_{z_1' z_1 = a_{11}} n(z_1|0, \sigma_1^2 I) dW.$$

donde dW es el diferencial de volumen adecuado, más concretamente es un elemento de volumen sobre la esfera $z_1' z_1 = a_{11}$. la densidad conjunta de b, v y a_{11} es

$$\int_{z_1' z_1 = a_{11}} n(b|\beta^2, \sigma^2/a_{11}) g_{n-1}(v/\sigma^2) \frac{1}{\sigma^2} n(z_1|0, \sigma_1^2 I) dW =$$

$$\begin{aligned}
 &= g_n(a_{11}/\sigma_1^2) \cdot n(b|\beta, \sigma^2/a_{11}) \cdot g_{n-1}(v/\sigma^2) \cdot (\sigma_1^2 \sigma^2) = \\
 &= \frac{a_{11}^{\frac{1}{2}(n-1)}}{(2\sigma_1^2)^{\frac{1}{2}n} \Gamma(\frac{1}{2}n)} \exp\left(-\frac{1}{2\sigma_1^2} a_{11}\right) \cdot \frac{\sqrt{a_{11}}}{\sqrt{2n\sigma^2}} \exp\left[-\frac{a_{11}}{2\sigma^2} (b-\beta)^2\right] \cdot \frac{1}{(2\sigma^2)^{\frac{1}{2}(n-1)} \Gamma(\frac{1}{2}(n-1))} v^{\frac{1}{2}(n-3)} \exp\left(-\frac{1}{2\sigma^2} v\right)
 \end{aligned}$$

Hagamos ahora $b = a_{12}/a_{11}$, $v = a_{22} - a_{12}^2/a_{11}$. El Jacobiano de la transformación

$$\frac{\partial(b, v)}{\partial(a_{11}, a_{22})} = \begin{vmatrix} \frac{1}{a_{11}} & 0 \\ -2\frac{a_{12}}{a_{11}^2} & 1 \end{vmatrix} = \frac{1}{a_{11}}.$$

Así la densidad para a_{11}, a_{12}, a_{22} será

$$\frac{a_{11}^{\frac{1}{2}(n-3)} \left(\frac{a_{11} a_{22} - a_{12}^2}{a_{11}} \right)^{\frac{1}{2}(n-3)} \cdot e^{-\frac{1}{2}Q}}{2^n \sigma_1^n \left(\frac{\sigma_1^2 \sigma_2^2 - \rho^2 \sigma_1^2 \sigma_2^2}{\sigma_1^2} \right) \sqrt{n} \Gamma(\frac{1}{2}n) \cdot \Gamma(\frac{1}{2}(n-1))}$$

donde

$$\begin{aligned}
 Q &= \frac{a_{11}}{\sigma_1^2} + \frac{a_{11}}{\sigma^2} \left(\frac{a_{12}^2}{a_{11}^2} - 2\rho \frac{\sigma_1 \sigma_2}{\sigma_1^2} \cdot \frac{a_{12}}{a_{11}} + \frac{\rho^2 \sigma_1^2 \sigma_2^2}{\sigma_1^4} \right) + \frac{1}{\sigma^2} \left(a_{22} - \frac{a_{12}^2}{a_{11}} \right) = \\
 &= a_{11} \left[\frac{1}{\sigma^2} + \frac{\rho^2 \sigma_1^2 \sigma_2^2}{\sigma_1^4 \sigma_2^2 (1-\rho^2)} \right] - 2a_{12} \frac{\rho \sigma_2}{\sigma_1 \sigma_2^2 (1-\rho^2)} + \frac{a_{22}}{\sigma_2^2 (1-\rho^2)} = \\
 &= \frac{1}{1-\rho^2} \left(\frac{a_{11}}{\sigma_1^2} - 2\rho \frac{a_{12}}{\sigma_1 \sigma_2} + \frac{a_{22}}{\sigma_2^2} \right).
 \end{aligned}$$

La densidad puede también escribirse

$$\frac{|A|^{\frac{1}{2}(n-3)} e^{-\frac{1}{2}Q}}{2^n |\Sigma|^{\frac{1}{2}n} \sqrt{n} \Gamma(\frac{1}{2}n) \Gamma(\frac{1}{2}(n-1))}$$

que es un caso particular de la distribución de Wishart que estudiaremos en el capítulo posterior.

Ahora la densidad de a_{11}, a_{22} , $r = a_{12}/\sqrt{a_{11}a_{22}}$ ($a_{12} = r\sqrt{a_{11}a_{22}}$) vendrá dada por

$$\frac{a_{11}^{\frac{1}{2}n-1} \cdot a_{22}^{\frac{1}{2}n-1} \cdot (1-r^2)^{\frac{1}{2}(n-3)} e^{-\frac{1}{2}Q}}{2^n [\sigma_1^2 \sigma_2^2 (1-\rho^2)]^{\frac{1}{2}n} \sqrt{n} \Gamma(\frac{1}{2}n) \cdot \Gamma(\frac{1}{2}(n-1))}$$

donde

$$Q = \frac{1}{1-\rho^2} \left(\frac{a_{11}}{\sigma_1^2} - 2\rho r \frac{\sqrt{a_{11}} \sqrt{a_{22}}}{\sigma_1 \sigma_2} + \frac{a_{22}}{\sigma_2^2} \right)$$

para obtener la densidad de r debemos integrar la anterior expresión respecto a a_{11} y a_{22} sobre la semicircunferencia real positiva. Podemos hacerlo de varias maneras, que dan lugar a diferentes expresiones para la densidad. Indicaremos aquí el método más directo. Desarrollaremos en primer lugar parte de la exponencial

$$\exp \left[\frac{\rho r \sqrt{a_{11}} \sqrt{a_{22}}}{(1-\rho^2) \sigma_1 \sigma_2} \right] = \sum_{\alpha=0}^{\infty} \frac{(\rho r \sqrt{a_{11}} \sqrt{a_{22}})^{\alpha}}{\alpha! [\sigma_1 \sigma_2 (1-\rho^2)]^{\alpha}}.$$

La misma expresión para la densidad de a_{11}, a_{22} y r será

$$\frac{(1-r^2)^{\frac{1}{2}(n-3)}}{\sigma_1^2 \sigma_2^2 (1-r^2)^{\frac{1}{2}n} 2^n \sqrt{n} \Gamma(\frac{1}{2}n) \Gamma(\frac{1}{2}(n-1))} \sum_{\alpha=0}^{\infty} \frac{(er)^\alpha}{\alpha! (1-r^2)^\alpha \sigma_1^\alpha \sigma_2^\alpha} \cdot \left\{ \exp \left[-\frac{a_{11}}{2(1-r^2)\sigma_1^2} \right] a_{11}^{\frac{1}{2}(n+\alpha)-1} \right\} \cdot \left\{ \exp \left[-\frac{a_{22}}{2(1-r^2)\sigma_2^2} \right] a_{22}^{\frac{1}{2}(n+\alpha)-1} \right\}.$$

puesto que

$$\int_0^{\infty} a_{11}^{\frac{1}{2}(n+\alpha)-1} \exp \left[-\frac{a_{11}}{2(1-r^2)\sigma_1^2} \right] da_{11} = \Gamma\left(\frac{1}{2}(n+\alpha)\right) \cdot [2\sigma_1^2(1-r^2)]^{\frac{1}{2}(n+\alpha)}$$

integrando la densidad (podemos integrar término a término) llegamos a

$$\begin{aligned} & \frac{(1-r^2)^{\frac{1}{2}(n-3)}}{\sigma_1^2 \sigma_2^2 (1-r^2)^{\frac{1}{2}n} 2^n \sqrt{n} \Gamma(\frac{1}{2}n) \Gamma(\frac{1}{2}(n-1))} \sum_{\alpha=0}^{\infty} \frac{(er)^\alpha}{\alpha! (1-r^2)^\alpha \sigma_1^\alpha \sigma_2^\alpha} \cdot \Gamma^2\left[\frac{1}{2}(n+\alpha)\right] \cdot 2^{(n+\alpha)} \sigma_1^{(n+\alpha)} \sigma_2^{(n+\alpha)} (1-r^2)^{(n+\alpha)} = \\ & = \frac{(1-r^2)^{\frac{1}{2}n} (1-r^2)^{\frac{1}{2}(n-3)}}{\sqrt{n} \Gamma(\frac{1}{2}n) \Gamma(\frac{1}{2}(n-1))} \sum_{\alpha=0}^{\infty} \frac{(2er)^\alpha}{\alpha!} \Gamma^2\left(\frac{1}{2}(n+\alpha)\right) \end{aligned}$$

Utilizando la fórmula $\Gamma(z) \Gamma(z + \frac{1}{2}) = \sqrt{\pi} \Gamma(2z) / 2^{2z-1}$ podemos modificar la constante. Empezamos entonces el siguiente lema:

Teorema.- La coeficiente de correlación en una muestra de tamaño N de una normal bivariate en correlación ρ se distribuye con densidad

$$\frac{2^{n-2} (1-r^2)^{\frac{1}{2}n} (1-r^2)^{\frac{1}{2}(n-3)}}{(n-2)! \pi} \sum_{\alpha=0}^{\infty} \frac{(2er)^\alpha}{\alpha!} \Gamma^2\left(\frac{1}{2}(n+\alpha)\right) \quad \text{donde } n = N-1.$$

La distribución de r fue obtenida por primera vez por Fisher en 1915, quien dio para su densidad la siguiente expresión

$$\frac{(1-r^2)^{\frac{1}{2}n} (1-r^2)^{\frac{1}{2}(n-3)}}{\pi(n-2)!} \left[\frac{d^{n-1}}{dx^{n-1}} \left\{ \frac{w^{-1}(-x)}{\sqrt{1-x^2}} \right\} \right]_{x=r\rho}$$

que se obtiene de las expresiones iniciales haciendo los cambios $a_{11} = u x^2$, $a_{22} = u x^2$.

Hotelling estudió exhaustivamente la distribución de r y recomendó el uso de la expresión

$$\frac{n-1}{\sqrt{2\pi}} \frac{\Gamma(n)}{\pi(n+\frac{1}{2})} (1-r^2)^{\frac{1}{2}n} (1-r^2)^{\frac{1}{2}(n-3)} (1-r\rho)^{-n+\frac{1}{2}} F\left(\frac{1}{2}, \frac{1}{2}; n+\frac{1}{2}; \frac{1+er}{2}\right)$$

donde

$$F(a, b; c; x) = \sum_{j=0}^{\infty} \frac{\Gamma(a+j)}{\Gamma(a)} \cdot \frac{\Gamma(b+j)}{\Gamma(b)} \cdot \frac{\Gamma(c)}{\Gamma(c+j)} \cdot \frac{x^j}{j!} \quad \text{es una función hipergeométrica.}$$

Esta expresión se obtiene también utilizando los cambios de variables utilizados por Fisher. La ventaja de utilizar esta fórmula para la densidad está en que la serie converge más rápidamente que la expresión que dimos en el teorema anterior.

La distribución acumulada de r , es decir $Pr(r \leq r^*) = F(r^* | N, \rho)$ ha sido tabulada por F.N. David (1938) para $\rho = 0, .1$ en incrementos de $.1$, $N = 3(1, 25, 50, 100, 200, 400)$ y $r^* = -.5(.05)1$. Dada la simetría de la densidad para r y ρ , tenemos que $F(r^* | N, \rho) = 1 - F(-r^* | N, -\rho)$. Veamos el uso de las tablas en algunos procedimientos estadísticos, de los que nos ocupamos a continuación.

a) Consideremos el problema de contrastar la hipótesis

3 (4)

$$H_0: \rho = \rho_0$$

a partir de una muestra determinada.

Si la alternativa es del tipo $H_1: \rho > \rho_0$, rechazaremos H_0 si el coeficiente de correlación muestral es tal que supera a r_0 , donde r_0 lo elegimos de manera que $1 - F(r_0 | N, \rho_0) = \alpha$, nivel de significancia.

Si la alternativa es de la forma $H_1: \rho < \rho_0$, rechazaremos para $r_{ij} < r'_0$ en r'_0 tal que $F(r'_0 | N, \rho_0) = \alpha$. Para $H_1: \rho \neq \rho_0$, la regla de rechazo viene dada por $r_{ij} < r'_1$ y $r_{ij} > r_1$ donde r_1 y r'_1 son tales que $F(r'_1 | N, \rho_0) + (1 - F(r_1 | N, \rho_0)) = \alpha$. David asegura que r_1 y r'_1 pueden elegirse de manera que ambas cosas tengan la misma probabilidad e igual a $\frac{1}{2}\alpha$. Demostró en 1937 que para $N \geq 10$ y $|\rho| \leq 0.8$ la región crítica es similar a la región de un test sesgado de H_0 , es decir, un test cuya función de potencia tiene su mínimo en ρ_0 .

Debemos también señalar que cualquier test basado en r es invariante bajo transformaciones del tipo $x_i^* = c_i x_i + d_i$ ($i = 1, \dots, N$, $c_i \neq 0$), que r es el invariante único de la estadística suficiente (lo que significa que cualquier otro test invariante puede expresarse como función de r). El procedimiento anterior para contrastar la hipótesis $H_0: \rho = \rho_0$ frente a alternativas del tipo $\rho > \rho_0$ es uniformemente más potente frente a todos los test invariantes.

Amos de ejemplos, supongamos que deseamos contrastar la hipótesis $H_0: \rho = 0.5$ frente a la alternativa $\rho \neq 0.5$ a un nivel $\alpha = 5\%$ utilizando una correlación muestral obtenida a partir de una muestra de tamaño 15. En las tablas de David encontramos (por interpolación)

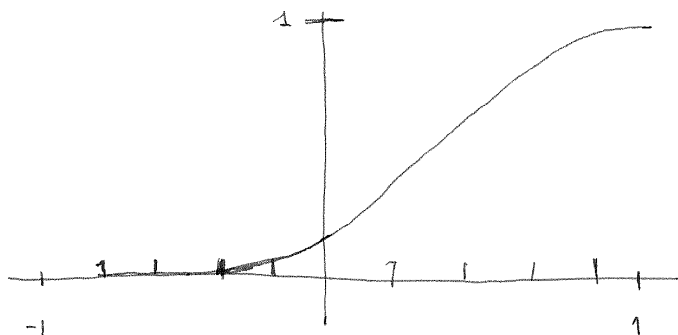
$F(0.027 | 15, 0.5) = 0.025$ y $F(0.805 | 15, 0.5) = 0.975$. Por tanto rechazaremos la hipótesis si r , en nuestra muestra, es menor que 0.027 o mayor que 0.805.

b) Podemos utilizar las tablas de David para calcular la función de potencia de un test de correlación. Si la región de rechazo de H_0 es $r > r_1$ y $r < r'_1$, la potencia del test es una función de la verdadera correlación ρ , $[1 - F(r_1 | N, \rho)] + [F(r'_1 | N, \rho)]$; es decir la probabilidad de rechazar la hipótesis nula cuando la correlación poblacional es ρ .

Consideremos, por ejemplo la función de potencia del test para $\rho = 0$ que obtenimos en el párrafo anterior. La región de rechazo (unilateral) es $r \geq 0.5494$ para $\alpha = 5\%$. Las probabilidades de rechazo son

ρ	-1.0	-0.8	-0.6	-0.4	-0.2	0	0.2	0.4	0.6	0.8	1
Prob	0.0000	0.0000	0.0004	0.0032	0.0147	0.0500	0.1376	0.3215	0.6235	0.9279	1.0000

que gráficamente da lugar a la siguiente representación



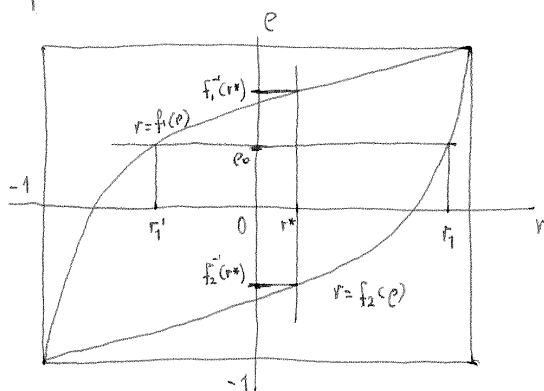
c) Incógnitas de David proporcionan función región de confianza para ρ . Para N dados, r_1' y r_1 , los dos puntos de significación son funciones de ρ , a saber $f_1(\rho)$ y $f_2(\rho)$ respectivamente, de manera que

$$\Pr \{ f_1(\rho) < r < f_2(\rho) | \rho \} = 1 - \alpha$$

Las funciones $f_1(\rho)$ y $f_2(\rho)$ son funciones monótonas crecientes en ρ si r_1 y r_1' se eligen de manera que $1 - F(r_1 | N, \rho) = \frac{1}{2}\alpha = F(r_1' | N, \rho)$. Si $\rho = f_1^{-1}(r)$ es la inversa de $r = f_1(\rho)$, $i=1,2$ entonces la desigualdad $f_i(\rho) < r$ es equivalente a $\rho < f_i^{-1}(r)$ y $r < f_2(\rho)$ es equivalente a $f_2^{-1}(r) < \rho$. De esta forma la anterior probabilidad puede escribirse

$$\Pr \{ f_2^{-1}(r) < \rho < f_1^{-1}(r) | \rho \} = 1 - \alpha$$

Es decir la probabilidad de extraer una muestra y obtener un r de manera que el intervalo de extremos $[f_2^{-1}(r), f_1^{-1}(r)]$ cubra a ρ es precisamente $1 - \alpha$. Se trata pues de un intervalo de confianza para ρ en nivel de confianza $1 - \alpha$. Para valores de N y α dados



Las curvas $r = f_1(\rho)$ y $r = f_2(\rho)$ tienen el aspecto de la figura. Cuando intentamos $H_0: \rho = \rho_0$, la intersección de la línea $\rho = \rho_0$ con las dos curvas dadas determina r_1 y r_1' de manera que el nivel sea el α fijado. Si queremos determinar una región de confianza para ρ sobre la base de r^* , coeficiente de correlación muestral, utilizaremos las curvas anteriores con $r = r^*$ y determinaremos los extremos del intervalo correspondiente. David proporciona estas curvas para $\alpha = 0.1, 0.05, 0.02$ y 0.01 y para varias N . Para el caso de test o regiones unilaterales utilizaremos tan sólo una desigualdad de las anteriores.

También es posible el uso de $F(r | N, \rho)$ para determinar intervalos de confianza. Dado r^* , $f_1^{-1}(r^*)$ es el valor de ρ que hace que $F(r^* | N, \rho) = \frac{1}{2}\alpha$ y análogamente $f_2^{-1}(r^*)$ es aquel valor de ρ para el cual $1 - F(r^* | N, \rho) = \frac{1}{2}\alpha$.

Por ejemplo, un intervalo de confianza con $1 - \alpha = 0.95$ basado en un $r^* = 0.7952$ obtenido a partir de una muestra de tamaño 10. Utilizando la tabla II de David encontramos como límites del intervalo 0.34 y 0.94

Criterio de verosimilitud

Es interesante encontrar el interior de la región de verosimilitud para intentar $H_0: \rho = \rho_0$ a partir de una muestra x_1, \dots, x_N de una $N(\mu, \Sigma)$ con $\mu' = (\mu_1, \mu_2)$ y $\sigma_{11} = \sigma_1^2, \sigma_{12} = \sigma_{21} = \rho\sigma_1\sigma_2$ y $\sigma_{22} = \sigma_2^2$. La función de verosimilitud es

$$K \sigma_1^{-N} \sigma_2^{-N} (1 - \rho^2)^{-N/2} \exp \left\{ -\frac{1}{2(1 - \rho^2)} \left[\frac{a_{11}}{\sigma_1^2} - 2\rho \frac{a_{12}}{\sigma_1\sigma_2} + \frac{a_{22}}{\sigma_2^2} + N \frac{(\bar{x}_1 - \mu_1)^2}{\sigma_1^2} - 2\rho N \frac{(\bar{x}_1 - \mu_1)(\bar{x}_2 - \mu_2)}{\sigma_1\sigma_2} + N \frac{(\bar{x}_2 - \mu_2)^2}{\sigma_2^2} \right] \right\}$$

Si maximizamos la anterior expresión cuando los parámetros varían en Ω (con las restricciones $\sigma_1^2 > 0, \sigma_2^2 > 0, \rho^2 < 1$) obtenemos $\hat{\sigma}_{1\Omega}^2 = a_{11}/N, \hat{\sigma}_{2\Omega}^2 = a_{22}/N, \hat{\rho}_{\Omega} = a_{12}/\sqrt{a_{11}a_{22}}, \hat{\mu}_{1\Omega} = \bar{x}_1, \hat{\mu}_{2\Omega} = \bar{x}_2$.

Si maximizamos en $w \in \Omega$ (restringido para $\rho = \rho_0$) $\hat{\mu}_{1w} = \bar{x}_1$ y $\hat{\mu}_{2w} = \bar{x}_2$ pues la forma cuadrática de $\bar{x}_1 - \mu_1$ y $\bar{x}_2 - \mu_2$ en el exponente es definida negativa (valores por tanto su máximo cuando vale 0). Para maximizar el esto tomamos logaritmos y tenemos

$$\log k - N \log \hat{\sigma}_1 - N \log \hat{\sigma}_2 - \frac{1}{2} N \log (1 - \rho^2) - \frac{1}{2(1-\rho^2)} \left(\frac{a_{11}}{\hat{\sigma}_1^2} - 2\rho_0 \frac{a_{11}}{\hat{\sigma}_1 \hat{\sigma}_2} + \frac{a_{22}}{\hat{\sigma}_2^2} \right)$$

Derivando respecto de $\hat{\sigma}_i$, $i=1,2$ y igualando a 0, tenemos

$$-\frac{N}{\hat{\sigma}_{iw}^2} - \frac{1}{2(1-\rho^2)} \left(-2 \frac{a_{ij}}{\hat{\sigma}_{iw}^3} + 2\rho_0 \frac{a_{ij}}{\hat{\sigma}_{iw}^2 \hat{\sigma}_{jw}} \right) = 0, \quad i \neq j, \quad i, j = 1, 2$$

Esto es

$$\frac{a_{ii}}{\hat{\sigma}_{iw}^2} - \rho_0 \frac{a_{ij}}{\hat{\sigma}_{iw} \hat{\sigma}_{jw}} = N(1-\rho^2)$$

Sumando para $i=1, j=2$ y $j=1, i=2$, llegamos

$$\frac{a_{11}}{\hat{\sigma}_{1w}^2} - 2\rho_0 \frac{a_{12}}{\hat{\sigma}_{1w} \hat{\sigma}_{2w}} + \frac{a_{22}}{\hat{\sigma}_{2w}^2} = 2N(1-\rho^2).$$

Restando para $i=1, j=2$ e $i=2, j=1$, tenemos

$$\frac{a_{11}}{\hat{\sigma}_{1w}^2} = \frac{a_{22}}{\hat{\sigma}_{2w}^2} = \frac{a^2}{\hat{\sigma}^2}. \quad \text{Entonces}$$

$$\frac{a^2}{\hat{\sigma}^2} - \frac{\rho_0 r a^2}{\hat{\sigma}^2} = N(1-\rho^2)$$

Así pues

$$\frac{N(1-\rho^2)}{1-\rho_0 r} = \frac{a^2}{\hat{\sigma}^2} = \frac{\sqrt{a_{11} a_{22}}}{\hat{\sigma}_{1w} \hat{\sigma}_{2w}}$$

y

$$|\hat{\Sigma}_w| = (1-\rho^2) \hat{\sigma}_{1w}^2 \hat{\sigma}_{2w}^2 = \frac{(1-\rho_0 r)^2}{1-\rho^2} \cdot \frac{a_{11}}{N} \cdot \frac{a_{22}}{N}.$$

El determinante máximo para $w \in \Omega$, con

$$\max_w L = \frac{(1-\rho_0^2)^{\frac{1}{2}N} N^N}{(2\pi)^N (1-\rho_0 r)^N a_{11}^{\frac{1}{2}N} a_{22}^{\frac{1}{2}N}} e^{-N}$$

$$\max_{\Omega} L = \frac{N^N}{(2\pi)^N (1-r^2)^{\frac{1}{2}N} a_{11}^{\frac{1}{2}N} a_{22}^{\frac{1}{2}N}} e^{-N}$$

El criterio de la razón de verosimilitud es por tanto

$$\frac{\max_w L}{\max_{\Omega} L} = \frac{(1-\rho^2)^{\frac{1}{2}N} (1-r^2)^{\frac{1}{2}N}}{(1-\rho_0 r)^N} = \left[\frac{(1-\rho^2)(1-r^2)}{(1-\rho_0 r)^2} \right]^{\frac{1}{2}N}$$

El test independiente es $(1-\rho^2)(1-r^2)(1-\rho r)^{-2} < c$, donde c recibe de manera que la probabilidad de la desigualdad, cuando las muestras se extraen de poblaciones normales bivalentes en $\rho = \rho_0$, sea el nivel de significación elegido. La región crítica puede también presentarse

$$(\rho_0^2 c - \rho_0^2 + 1)r^2 - 2\rho_0 c r + c - 1 + \rho_0^2 > 0,$$

o bien

$$r > \frac{\rho_0 c + (1-\rho_0^2)\sqrt{1-c}}{\rho_0^2 c + 1 - \rho_0^2}, \quad r < \frac{\rho_0 c - (1-\rho_0^2)\sqrt{1-c}}{\rho_0^2 c + 1 - \rho_0^2}$$

Así, el test de la unión de similitud de $H_0: \rho = \rho_0$ frente a la alternativa $H_1: \rho \neq \rho_0$ tiene una región crítica de la forma $r > r_1$ y $r < r_1'$, pero r_1 y r_1' no están elegidos de manera que la probabilidad de desigualdad es $\frac{1}{2}\alpha$ cuando H_0 es cierta, sino de la forma descrita en las dos últimas desigualdades y en c teniendo el significado ante dicho (la probabilidad de ambas desigualdades igual a α).

La distribución asintótica de un coeficiente de correlación muestral y de la Z de Fisher

Se trata de demostrar en este apartado que a medida que aumenta el tamaño de la muestra la distribución del coeficiente de correlación muestral tiende a la normal. La distribución de una función particular de la correlación muestral, conocida como la Z de Fisher (obtenida por Fisher en 1921), y cuya variación es aproximadamente independiente de la correlación poblacional, tiende más rápidamente a la normal. En primer lugar demostraremos un teorema central del límite para variables m -dimensionales.

TEOREMA.- Sean X_1, X_2, \dots vectores m -dimensionales independientes e idénticamente distribuidos con medias $E(X_i) = \mu$ y covarianzas $E[(X_i - \mu)(X_i - \mu)'] = T$. Entonces la distribución límite de $(X/\sqrt{n}) \sum_{i=1}^n (X_i - \mu)$ cuando $n \rightarrow \infty$ es $N(0, T)$.

Demostración.- Sea

$$\phi_n(t, u) = E \left[\exp \left[i u' \frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i - \mu) \right] \right],$$

donde u es un vector m -componente. Para valores de t fijos, $\phi_n(t, u)$ pueden considerarse como la función característica de $(X/\sqrt{n}) \sum_{i=1}^n (X_i - \mu)$. Por el teorema central del límite (para variables unidimensionales), la distribución límite correspondiente será $N(0, T)$. Por tanto

$$\lim_{n \rightarrow \infty} \phi_n(t, u) = e^{-\frac{1}{2} u' T u}, \quad \forall u, t$$

Entonces haciendo $u=1$, tenemos

$$\lim_{n \rightarrow \infty} E \left[\exp \left(i t' \frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i - \mu) \right) \right] = e^{-\frac{1}{2} t' T t}, \quad \forall t$$

Como $e^{-\frac{1}{2} t' T t}$ es continua para $t=0$, la convergencia es uniforme en algún entorno de cero. El teorema es cumplido.

Vamos a demostrar a continuación que la distribución asintótica de la matriz de covarianzas muestral es normal.

TEOREMA.- Sea $A(n) = \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})'$, donde X_1, \dots son todas independientes y $N(\mu, \Sigma)$ y $n = N-1$. Entonces la distribución asintótica de $B(n) = (X/\sqrt{n}) [A(n) - n\Sigma]$ es normal con media 0 y covarianzas

$$E(b_{ij}(n), b_{kl}(n)) = \sigma_{ik} \sigma_{jl} + \sigma_{il} \sigma_{jk}.$$

Demostración.- Como demostramos anteriormente $A(n)$ se distribuye como $A(n) = \sum_{i=1}^n Z_i Z_i'$, donde Z_1, Z_2, \dots se distribuyen $N(0, \Sigma)$ e independientes. Podemos arreglar los elementos de $Z_i Z_i'$ de forma vectorial como sigue

$$Y_i' = [Z_{i1}^2, Z_{i1} Z_{i2}, \dots, Z_{i2}^2, \dots, Z_{ip}^2]$$

los momentos X_α pueden deducirse de los momentos de Z_α tal como acostumbra en sumandos. 3 (E)
 Tendremos $E(Z_{i\alpha} Z_{j\alpha}) = \sigma_{ij}$, $E(Z_{i\alpha} Z_{j\alpha} Z_{k\alpha} Z_{l\alpha}) = \sigma_{ij} \sigma_{kl} + \sigma_{ik} \sigma_{jl} + \sigma_{il} \sigma_{jk}$, $E[(Z_{i\alpha} Z_{j\alpha} - \sigma_{ij})(Z_{k\alpha} Z_{l\alpha} - \sigma_{kl})] = \sigma_{ik} \sigma_{jl} + \sigma_{il} \sigma_{jk}$ (se han obtenido mediante la función característica). Así, los vectores X_α satisfacen las condiciones del teorema anterior cuando los elementos de V los elementos de Σ aneclados de forma similar a como lo hemos hecho con X_α y los elementos de T los que acabamos de obtener. Si los n elementos de $A(n)$ los aneclamos también en forma vectorial, $W(n)$, entonces $W(n) - n\mu = \sum_{\alpha=1}^n (X_\alpha - \mu)$. Aplicando ahora el teorema anterior $(1/\sqrt{n})[W(n) - n\mu]$ se distribuye asintóticamente normal con media 0 y la matriz de varianzas la de X_α .

Estamos particularmente interesados en la correlación muestral

$$r(n) = \frac{A_{ij}(n)}{\sqrt{A_{ii}(n) A_{jj}(n)}} \quad , \quad i \neq j$$

Esto puede también escribirse de la forma

$$r(n) = \frac{C_{ij}(n)}{\sqrt{C_{ii}(n) C_{jj}(n)}}$$

donde $C_{gh}(n) = A_{gh}(n) / \sqrt{\sigma_{gg} \sigma_{hh}}$. El conjunto $C_{ii}(n)$, $C_{jj}(n)$ y $C_{ij}(n)$ se distribuye como

$$\sum_{\alpha=1}^n \begin{pmatrix} Z_{i\alpha}^* \\ Z_{j\alpha}^* \end{pmatrix} \begin{pmatrix} Z_{i\alpha}^* & Z_{j\alpha}^* \end{pmatrix} = \sum_{\alpha=1}^n \begin{pmatrix} Z_{i\alpha} / \sqrt{\sigma_{ii}} \\ Z_{j\alpha} / \sqrt{\sigma_{jj}} \end{pmatrix} \begin{pmatrix} Z_{i\alpha} / \sqrt{\sigma_{ii}} & Z_{j\alpha} / \sqrt{\sigma_{jj}} \end{pmatrix}$$

donde $(Z_{i\alpha}^*, Z_{j\alpha}^*)$ son independientes, cada una con distribución $N\left[\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}\right]$ y $\rho = \sigma_{ij} / \sqrt{\sigma_{ii} \sigma_{jj}}$. Sean

$$U(n) = \frac{1}{n} \begin{pmatrix} C_{ii}(n) \\ C_{jj}(n) \\ C_{ij}(n) \end{pmatrix} \quad \text{y} \quad b = \begin{pmatrix} 1 \\ 1 \\ \rho \end{pmatrix}$$

Entonces $\sqrt{n}(U(n) - b)$ es asintóticamente normal con media 0 y matriz de varianzas

$$\begin{pmatrix} 2 & 2\rho^2 & 2\rho \\ 2\rho^2 & 2 & 2\rho \\ 2\rho & 2\rho & 1+\rho^2 \end{pmatrix}$$

Resolvamos ahora el siguiente teorema general

TEOREMA .- Sea $U(n)$ un vector con m componentes y b un vector fijo. Supongamos que $\sqrt{n}(U(n) - b)$ se distribuye asintóticamente $N(0, T)$. Sea $W = f(u)$ una función de un vector u cuyas primera y segunda derivadas existen en un entorno de $u=b$. Sea $\frac{\partial f(u)}{\partial u_i} \Big|_{u=b}$ la i -ésima componente de ϕ_b . Entonces la distribución límite de $\sqrt{n}(f(U(n)) - f(b))$ es $N(0, \phi_b' T \phi_b)$.

La demostración del teorema puede consultarse en el texto de Cramer.

Usar luego el $U(n)$ ante definido, con b y T los correspondientes valores ante introducidos, satisface las condiciones del teorema. La función

$$r = \frac{u_3}{\sqrt{u_1 u_2}} = u_3 u_1^{-1/2} u_2^{-1/2}$$

satisface las condiciones de derivabilidad exigidas. Los elementos de ϕ_b son

$$\frac{\partial r}{\partial u_1} \Big|_{u=b} = -\frac{1}{2} u_3 u_1^{-3/2} u_2^{-1/2} \Big|_{u=b} = -\frac{1}{2} \rho, \quad \frac{\partial r}{\partial u_2} \Big|_{u=b} = -\frac{1}{2} u_3 u_1^{-1/2} u_2^{-3/2} \Big|_{u=b} = -\frac{1}{2} \rho, \quad \frac{\partial r}{\partial u_3} = u_1^{-1/2} u_2^{-1/2} \Big|_{u=b} = 1$$

y $f(\rho) = \rho$. La varianza asintótica de $\sqrt{n}(r(n) - \rho)$ es

$$\begin{pmatrix} -\frac{1}{2}\rho & -\frac{1}{2}\rho & 1 \end{pmatrix} \begin{pmatrix} 2 & 2\rho & 2\rho \\ 2\rho & 2 & 2\rho \\ 2\rho & 2\rho & 1-\rho^2 \end{pmatrix} \begin{pmatrix} -\frac{1}{2}\rho \\ -\frac{1}{2}\rho \\ 1 \end{pmatrix} = (\rho - \rho^3, \rho - \rho^3, 1 - \rho^2) \begin{pmatrix} -\frac{1}{2}\rho \\ -\frac{1}{2}\rho \\ 1 \end{pmatrix} = 1 - 2\rho^2 + \rho^4 = (1 - \rho^2)^2.$$

Hemos obtenido el siguiente teorema:

TEOREMA.- Si $r(n)$ es el coeficiente de correlación muestral de una muestra de tamaño $N = (n+1)$ de una distribución normal con correlación ρ , entonces $\sqrt{n}(r(n) - \rho)/(1 - \rho^2)$ [o $\sqrt{N}(r(n) - \rho)/(1 - \rho^2)$] se distribuye asintóticamente $N(0,1)$.

Apyandanos en el teorema anterior, si $f(x)$ es una función con primera y segunda derivadas en $x = \rho$, entonces $\sqrt{n}(f(r) - f(\rho))$ también es asintóticamente normal con media 0 y varianza $\left(\frac{df}{dx}\bigg|_{x=\rho}\right)^2 (1 - \rho^2)^2$.

Es usual considerar f una función cuya varianza asintótica sea constante (la unidad) independiente del parámetro ρ . Una función de este tipo satisface la ecuación

$$f'(\rho) = \frac{1}{1 - \rho^2} = \frac{1}{2} \left(\frac{1}{1 + \rho} + \frac{1}{1 - \rho} \right)$$

Podemos tomar $f(\rho) = \frac{1}{2} [\log(1 + \rho) - \log(1 - \rho)] = \frac{1}{2} \log \left[\frac{1 + \rho}{1 - \rho} \right]$. La llamada z de Fisher es

$$Z = \frac{1}{2} \log \frac{1+r}{1-r} \quad \text{y sea} \quad \xi = \frac{1}{2} \log \frac{1+\rho}{1-\rho}.$$

TEOREMA.- Sea Z definida como ante, donde r es el coeficiente de correlación muestral de una muestra de tamaño $N = n+1$, de una normal bivariante con correlación ρ , sea ξ el valor ante definido. Entonces $\sqrt{n}(Z - \xi)$ es asintóticamente normal con media 0 y varianza 1.

Se puede demostrar que aproximadamente

$$E(Z) \sim \xi + \frac{\rho}{2n} \quad E(Z - \xi)^2 \sim \frac{1}{n-2} \sim E\left(Z - \xi - \frac{\rho}{2n}\right)^2$$

Esta última aproximación se sigue de la expresión

$$E(Z - \xi)^3 = \frac{1}{n} + \frac{8\rho^2}{4n^2} + \dots$$

Es muy buena para valores de ρ^2/n pequeños. Neyman (1953) estudió los momentos de Z hasta el orden n^{-3} . Una importante propiedad de la z de Fisher es que la aproximación a la normal es mucho más rápida que para r . David (1938) estableció comparaciones entre las probabilidades tabuladas y las probabilidades obtenidas suponiendo que Z se distribuiría normalmente. Recomendó que para $N \geq 25$ se tome Z como distribuido normalmente con media y varianzas los dados anteriormente.

Antes de formular con los coeficientes de correlación daremos algunos puntos para la utilización de la z de Fisher.

a) Supongamos que deseamos contrastar la hipótesis $\rho = \rho_0$ sobre la base de una muestra de tamaño N frente a la alternativa $\rho \neq \rho_0$. Calcularemos r y después Z . Sea

$$\xi_0 = \frac{1}{2} \log \frac{1 + \rho_0}{1 - \rho_0}$$

La región de rechazo al 5% tendrá dada por

$$\sqrt{N-3} |Z - \xi_0| > 1.96$$

Una región mejor es

$$\sqrt{N-3} \left| Z - \xi_0 - \frac{1}{2} \rho_0 / (N-1) \right| > 1.96$$

- b) Supongamos que tenemos una muestra de tamaño N_1 de una población y una muestra N_2 de otra población. Como contrastar la hipótesis de que los dos coeficientes de correlación son iguales, $\rho_1 = \rho_2$?
 Del teorema de la z de Fisher sabemos que si la hipótesis nula es cierta $Z_1 - Z_2$ es asintóticamente normal con media 0 y varianzas $1/(N_1-3) + 1/(N_2-3)$. Como región crítica al 5%, utilizaremos

$$\frac{|Z_1 - Z_2|}{\sqrt{1/(N_1-3) + 1/(N_2-3)}} > 1.96$$

- c) Bajo las indicaciones del apartado anterior supongamos que $\rho_1 = \rho_2 = \rho$. Como utilizar los resultados de ambas muestras para estimar conjuntamente ρ ? Puesto que Z_1 y Z_2 tienen varianzas $1/(N_1-3)$ y $1/(N_2-3)$ respectivamente, podemos estimar ξ mediante

$$\frac{(N_1-3)Z_1 + (N_2-3)Z_2}{N_1 + N_2 - 6}$$

y convertirlo esto en un estimador de ρ mediante la transformación inversa de la z .

- d) Sea r el coeficiente de correlación muestral de N observaciones. Como obtener un intervalo de confianza para ρ ? Sabemos que aproximadamente

$$Pr \{-1.96 < \sqrt{N-3} (z - \xi) < 1.96\} = 0.95$$

De aquí deducimos que $[-1.96/\sqrt{N-3} + \xi, 1.96/\sqrt{N-3} + \xi]$ es un intervalo de confianza para ξ .
 A partir de aquí obtenemos la región para ρ utilizando la transformación inversa $\rho = \tanh \xi = (e^\xi - e^{-\xi}) / (e^\xi + e^{-\xi})$ que es una transformación monótonica. La región es

$$\tanh(z - 1.96/\sqrt{N-3}) < \rho < \tanh(z + 1.96/\sqrt{N-3}).$$

3. COEFICIENTES DE CORRELACION PARCIAL

Recordemos que los coeficientes son coeficientes de correlación en distribuciones condicionales. Como vimos anteriormente si \mathbf{X} se distribuye $N(\mu, \Sigma)$ la distribución condicional del subvector $\mathbf{X}^{(1)}$ dado $\mathbf{X}^{(2)} = \mathbf{x}^{(2)}$, viene dada por $N(\mu^{(1)} + \beta(x^{(2)} - \mu^{(2)}), \Sigma_{11.2})$, donde

$$\beta = \Sigma_{12} \Sigma_{22}^{-1}$$

$$\Sigma_{11.2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

Las correlaciones parciales de $\mathbf{X}^{(1)}$ dado $\mathbf{x}^{(2)}$ son las correlaciones, calculadas de la forma ya establecida, a partir de $\Sigma_{11.2}$. Vamos a considerar a cabo la estimación de estos valores y como, a partir de ellos, podremos efectuar contraste de hipótesis.

En primer lugar nos vamos a ocupar del problema de la estimación. Supongamos que tenemos una muestra de tamaño N de una normal $N(\mu, \Sigma)$. Cuáles serán las estimaciones máximo-verosímiles de los coeficientes parciales de $\mathbf{X}^{(1)}$, $\rho_{1j}, j=1, \dots, p$? Sabemos que el estimador máximo-verosímil de Σ , es

$$\hat{\Sigma} = \frac{1}{N} \sum_{\alpha=1}^N (x_\alpha - \bar{x})(x_\alpha - \bar{x})', \text{ con } \bar{x} = \frac{1}{N} \sum_{\alpha=1}^N x_\alpha.$$

La independencia entre Σ y $\Sigma_{11.2}$, β y Σ_{22} es una a una en virtud de las relaciones (definiciones) anteriores, de

$$\hat{\Sigma}_{12} = \hat{\beta} \hat{\Sigma}_{22}$$

$$\hat{\Sigma}_{11} = \hat{\Sigma}_{11.2} + \hat{\beta} \hat{\Sigma}_{22} \hat{\beta}'$$

por tanto, teniendo en cuenta una propiedad de los estimados máximo-verosímiles podemos afirmar que las estimaciones de $\Sigma_{11.2}$ y β y Σ_{22} vendrán dadas por $\hat{\Sigma}_{11.2} = \hat{\Sigma}_{11} - \hat{\Sigma}_{12} \hat{\Sigma}_{22}^{-1} \hat{\Sigma}_{21}$, $\hat{\beta} = \hat{\Sigma}_{12} \hat{\Sigma}_{22}^{-1}$ y $\hat{\Sigma}_{22}$. Además las estimaciones de los coeficientes de correlación parcial serán:

$$\hat{\rho}_{ij, q+1, \dots, p} = \frac{\hat{\sigma}_{ij, q+1, \dots, p}}{\sqrt{\hat{\sigma}_{ii, q+1, \dots, p} \cdot \hat{\sigma}_{jj, q+1, \dots, p}}} \quad i, j = 1, \dots, q$$

donde $\hat{\sigma}_{ij, q+1, \dots, p}$ es el ij -ésimo elemento de $\hat{\Sigma}_{11,2}$.

TEOREMA.- Sea x_1, \dots, x_N una muestra de tamaño N de una $N(\mu, \Sigma)$. Los estimadores máximo-verosímiles de $\rho_{ij, q+1, \dots, p}$, correlaciones parciales de los primeros q componentes condicionadas a los $p-q$ restantes, vienen dados por

$$\hat{\rho}_{ij, q+1, \dots, p} = \frac{a_{ij, q+1, \dots, p}}{\sqrt{a_{ii, q+1, \dots, p} \cdot a_{jj, q+1, \dots, p}}}$$

donde

$$(a_{ij, q+1, \dots, p}) = A_{11} - A_{12} A_{22}^{-1} A_{21} \quad ,$$

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = A = \sum_{\alpha=1}^N (x_{\alpha} - \bar{x})(x_{\alpha} - \bar{x})'$$

la estimación $\hat{\rho}_{ij, q+1, \dots, p}$ que denotaremos por $\hat{r}_{ij, q+1, \dots, p}$ se interpreta como el coeficiente de correlación muestral entre \bar{X}_i y \bar{X}_j cuando $\bar{X}_{q+1}, \dots, \bar{X}_p$ permanecen constante.

Las interpretaciones geométricas pueden darse para estos resultados. A saber:

a) En un espacio p -dimensional consideremos los N puntos x_1, \dots, x_N , todos muestrales. La función de regresión

$$x^{(i)} = \bar{x}^{(i)} + \hat{\beta} (x^{(i)} - \bar{x}^{(i)})$$

es un hiperplano $p-q$ dimensional que resulta de la intersección de los q hiperplanos $p-1$ dimensionales

$$x_i = \bar{x}_i + \sum_{j=q+1}^p \hat{\beta}_{ij} (x_j - \bar{x}_j) \quad , \quad i = 1, \dots, q$$

Aquí $\hat{\beta}_{ij}$ es un elemento de $\hat{\beta} = \hat{\Sigma}_{12} \hat{\Sigma}_{22}^{-1} = A_{12} A_{22}^{-1}$. La i -ésima fila de $\hat{\beta}$ es $[\hat{\beta}_{i, q+1}, \dots, \hat{\beta}_{i, p}]$. La parte derecha de la anterior igualdad es la función de regresión mínimo-cuadrática de x_i sobre x_{q+1}, \dots, x_p ; es decir, si proyectamos los puntos x_1, \dots, x_N en el hiperplano de las coordenadas x_i, x_{q+1}, \dots, x_p , entonces la expresión anterior es el plano de regresión. El punto de coordenadas

$$x_i = \bar{x}_i + \sum_{j=q+1}^p \hat{\beta}_{ij} (x_{jk} - \bar{x}_j) \quad , \quad i = 1, \dots, q$$

$$x_j = x_{jk} \quad , \quad j = q+1, \dots, p$$

está en el hiperplano anterior. La diferencia entre la i -ésima coordenada de x_{jk} y el punto anterior

$$y_{ik} = x_{ik} - \left[\bar{x}_i + \sum_{j=q+1}^p \hat{\beta}_{ij} (x_{jk} - \bar{x}_j) \right] \quad \text{para } i = 1, \dots, q \quad \text{y } 0 \text{ para las otras coordenadas. Pongamos}$$

$$y_{\alpha} = \begin{bmatrix} y_{1\alpha} \\ \vdots \\ y_{q\alpha} \end{bmatrix}$$

Estos puntos pueden representarse como N puntos en un espacio q -dimensional. Entonces $A_{12} = \sum_{\alpha=1}^N y_{\alpha} x_{\alpha}'$.

b) Podemos, por otra parte, interpretar la muestra como p puntos en un espacio N -dimensional. Sea

$Z_j = (x_{j1}, \dots, x_{jN})$ el j -ésimo punto (las j -ésimas coordenadas de la muestra), y hagamos $Z_{p+1} = (1, \dots, 1)$ para el $p+1$ -ésimo punto. El punto de coordenadas $\bar{x}_1, \dots, \bar{x}_i$ es $\bar{x}_i Z_{p+1}$. La proyección de Z_i sobre el hiperplano definido mediante Z_{q+1}, \dots, Z_{p+1} es

$$Z_i^* = \bar{x}_i Z_{p+1} + \sum_{j=q+1}^p \hat{\beta}_{ij} (Z_j - \bar{x}_j Z_{p+1})$$

y es el punto del hiperplano cuya distancia a Z_i es mínima. Sea \tilde{Z}_i el vector de Z_i^* a Z_i , 3 (8)
 es decir $Z_i - Z_i^*$, o mejor en equivalente trasladado de manera que originen coincide en el origen de
 coordenadas. El conjunto de vectores $\tilde{Z}_1, \dots, \tilde{Z}_p$ son las proyecciones de Z_1, \dots, Z_p en el hiperplano ortogonal
 a Z_{q+1}, \dots, Z_p . Entonces $\tilde{Z}_i^T \tilde{Z}_j = a_{ij} \cdot q + s, \dots, p$ es el cuadrado de la longitud de \tilde{Z}_i , o lo que es lo mismo, el
 cuadrado de la distancia de Z_i^* a Z_i . Así $\tilde{Z}_i^T \tilde{Z}_j / \sqrt{\tilde{Z}_i^T \tilde{Z}_i \tilde{Z}_j^T \tilde{Z}_j} = r_{ij} \cdot q + s, \dots, p$ es el coseno del ángulo entre
 \tilde{Z}_i y \tilde{Z}_j .

Un ejemplo de coeficiente de correlación parcial.

Consideremos algunos datos [(Hooker, 1907) J. Roy. Stat. Soc. vol 70. pp. 1-42] acerca de la cosecha de heno (X_1)
 en cientos de libras de por acre, cantidad de lluvia caída durante la primavera (X_2) y temperaturas acumuladas
 superiores a $42^\circ F$ en la primavera (X_3) para una determinada zona de Inglaterra durante un periodo de 20 años.
 Las estimaciones del μ_i ($= \sqrt{\sigma_{ii}}$) y ρ_{ij} son

$$\hat{\mu} = \bar{x} = \begin{bmatrix} 28.02 \\ 4.91 \\ 59.4 \end{bmatrix}, \quad \begin{bmatrix} \hat{\sigma}_1 \\ \hat{\sigma}_2 \\ \hat{\sigma}_3 \end{bmatrix} = \begin{bmatrix} 4.42 \\ 1.10 \\ 8.5 \end{bmatrix}, \quad \begin{bmatrix} 1 & \hat{\rho}_{12} & \hat{\rho}_{13} \\ \hat{\rho}_{21} & 1 & \hat{\rho}_{23} \\ \hat{\rho}_{31} & \hat{\rho}_{32} & 1 \end{bmatrix} = \begin{bmatrix} 1.00 & 0.80 & -0.40 \\ 0.80 & 1.00 & -0.56 \\ -0.40 & -0.56 & 1.00 \end{bmatrix}$$

A partir de las correlaciones observamos que la cantidad de lluvia caída (x_2) y la cosecha (x_1) están relacionadas
 positivamente, cosecha y temperatura lo están negativamente así como que lluvia y temperatura. Quié interpre-
 tación puede darse a la relación negativa aparente que existe entre la cosecha y la temperatura? ¿Siempre
 las altas temperaturas, a cambio de bajas cosechas o las altas temperaturas están asociadas con bajas cantidades
 de lluvia y por tanto con bajas cosechas? Para responder a estas cuestiones podemos considerar la correlación
 entre cosecha y temperatura cuando la cantidad de lluvia es constante; es decir, se trata de utilizar los
 datos anteriores para estimar el coeficiente de correlación parcial entre X_1 y X_3 . Su valor es

$$\frac{\hat{\sigma}_{13.2}}{\sqrt{\hat{\sigma}_{11.2} \hat{\sigma}_{33.2}}} = 0.097$$

Así, al efecto de la lluvia eliminado, cosecha y temperatura están relacionados positivamente. La conclu-
 sión es que cuando, altas lluvias y altas temperaturas aumentan las cosechas de heno, pero en casi todos
 los años observados las grandes lluvias han venido con bajas temperaturas y viceversa.

La distribución del coeficiente de correlación parcial

Las correlaciones parciales se obtienen a partir de $A_{11.2} = A_{11} - A_{12}A_{22}^{-1}A_{21}$ de la misma forma mediante la
 que obtenimos las correlaciones de A . Para obtener la distribución de las correlaciones demostramos que A se distribuye
 como $\sum_{\alpha=1}^{n-1} Z_\alpha Z_\alpha'$ donde las Z_α eran independientes y cada una de ellas $N(0, I)$. Actuaremos de forma similar
 ahora. Para ello demostraremos a continuación un lema que será también de utilidad posterior.

TEOREMA. - Supongamos que X_1, \dots, X_n son independientes con X_α distribuida como $N(\Gamma W_\alpha, \Phi)$, donde W_α es un
 vector de r componentes. Sea $G = \sum_{\alpha} X_\alpha W_\alpha' H^{-1}$ donde $H = \sum_{\alpha} W_\alpha W_\alpha'$ y es no singular. Entonces $\sum_{\alpha=1}^m X_\alpha X_\alpha' -$
 $- G H G'$ se distribuye como $\sum_{\alpha=1}^{m-r} U_\alpha U_\alpha'$, donde U_α son independientes y $N(0, \Phi)$, y además independiente
 de G .

Demostración. - Sea $W = (W_1, \dots, W_m)$ y sea F una matriz cuadrada tal que $F H F' = I$, entonces $F^{-1} H^{-1} F' = I$. Sea
 $E_2 = F W$, entonces $W = E_2 F'$. Así

$$E_2 E_2' = F W W' F' = F \sum_{\alpha} W_\alpha W_\alpha' F' = F H F' = I.$$

Resulta que las filas de E_2 son vectores ortogonales. Sea posible encontrar una matriz $E_2, (m-r) \times m$ tal que
 $E' = [E_1' E_2']$ sea ortogonal. Sea ahora $Y = (Y_1, \dots, Y_m) = U E$, o $U = Y E'$ (es decir, $U_\alpha = \sum_{\beta} e_{\alpha\beta} Y_\beta$). Apliquen-
 do un conocido lema, las columnas de U , es decir U_α , son independientes y distribuidas normalmente, cada
 una de ellas con covarianza Φ . La media resultará dada por

$$E(U) = E(YE') = FWE' = F'F'E_2(E_1'E_2) = (F'F'E_2E_1' \quad F'F'E_2E_2') = (0 \quad F'F')$$

Para completar la demostración necesitamos demostrar que

$$\sum_{\alpha=1}^{m-r} Y_{\alpha} Y_{\alpha}' - G H G' = \sum_{\alpha=1}^{m-r} U_{\alpha} U_{\alpha}'$$

Sabemos que $\sum_{\alpha=1}^m Y_{\alpha} Y_{\alpha}' = \sum_{\alpha=1}^m U_{\alpha} U_{\alpha}'$ por la transformación que pasa de uno a otros, ortogonal. También verifican

$$\begin{aligned} G H G' &= (Y W H^{-1}) H (H^{-1} W Y') = U E E_2' (F^{-1})' H^{-1} F^{-1} E_2 E' U' = U \begin{pmatrix} E_1' \\ E_2' \end{pmatrix} E_2' E_2 (E_1' E_2') U' = \\ &= U \begin{pmatrix} 0 \\ I \end{pmatrix} (0 \quad I) U' = \sum_{\alpha=m-r+1}^m U_{\alpha} U_{\alpha}' \end{aligned}$$

Así pues

$$\sum_{\alpha=1}^m Y_{\alpha} Y_{\alpha}' - G H G' = \sum_{\alpha=1}^m U_{\alpha} U_{\alpha}' - \sum_{\alpha=m-r+1}^m U_{\alpha} U_{\alpha}' = \sum_{\alpha=1}^{m-r} U_{\alpha} U_{\alpha}'$$

De las anteriores consideraciones se sigue que si $P=0$, $E(U)=0$ y obtenemos el siguiente Corolario

COROLARIO.- Si $P=0$, la matriz $G H G'$ del teorema anterior se distribuye como $\sum_{\alpha=m-r+1}^m U_{\alpha} U_{\alpha}'$ donde las U_{α} son independientes y distribuidas $N(0, \phi)$.

Podemos ahora encontrar la distribución de $A_{11,2}$ de la misma forma. Vimos que A se distribuye como $\sum_{\alpha=1}^{N-1} Z_{\alpha} Z_{\alpha}'$, donde las Z_{α} son independientes, cada una $N(0, \Sigma)$. Si fraccionamos Z_{α} en dos subvectores de dimensiones q y $p-q$ respectivamente,

$$Z_{\alpha} = \begin{pmatrix} Z_{\alpha}^{(q)} \\ Z_{\alpha}^{(p-q)} \end{pmatrix}$$

Entonces $A_{ij} = \sum_{\alpha=1}^{N-1} Z_{\alpha}^{(i)} Z_{\alpha}^{(j)'}.$ La distribución condicional de $Z_1^{(q)} \dots Z_{N-1}^{(q)}$ dado $Z_1^{(q)} = z_1^{(q)}, \dots, Z_{N-1}^{(q)} = z_{N-1}^{(q)}$ es

$$\frac{\prod_{\alpha=1}^{N-1} n(z_{\alpha} | 0, \Sigma)}{\prod_{\alpha=1}^{N-1} n(z_{\alpha}^{(q)} | 0, \Sigma)} = \prod_{\alpha=1}^{N-1} \frac{n(z_{\alpha} | 0, \Sigma)}{n(z_{\alpha}^{(q)} | 0, \Sigma)} = \prod_{\alpha=1}^{N-1} n(z_{\alpha}^{(p-q)} | \beta z_{\alpha}^{(q)}, \Sigma_{11,2})$$

donde $\beta = \Sigma_{12} \Sigma_{22}^{-1}$ y $\Sigma_{11,2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$. Aplicando ahora el teorema con $Z_{\alpha}^{(q)} = Y_{\alpha}$, $Z_{\alpha}^{(p-q)} = W_{\alpha}$, $N-1 = m$, $p-q = r$ $\beta = F$, $\Sigma_{11,2} = \phi$, $A_{11} = \sum Y_{\alpha} Y_{\alpha}'$, $A_{12} A_{22}^{-1} = G$, $A_{22} = H$. Vemos que la distribución condicional de $A_{11} - A_{12} A_{22}^{-1} A_{21} = A_{11,2}$ dado $Z_{\alpha}^{(q)} = z_{\alpha}^{(q)}$ es la de $\sum_{\alpha=1}^{m-r} U_{\alpha} U_{\alpha}'$ donde U_{α} son independientes, cada uno con distribución $N(0, \Sigma_{11,2})$. Como esta distribución no depende de $\{z_{\alpha}^{(q)}\}$ entonces obtenemos finalmente:

TEOREMA.- La matriz $A_{11,2} = A_{11} - A_{12} A_{22}^{-1} A_{21}$ se distribuye como $\sum_{\alpha=1}^{m-r} U_{\alpha} U_{\alpha}'$, donde U_{α} se distribuyen independientes, cada uno con distribución $N(0, \Sigma_{11,2})$

Como corolario obtenemos

COROLARIO.- Si $\Sigma_{12} = 0$ ($\beta = 0$) entonces $A_{11,2}$ se distribuye como $\sum_{\alpha=1}^{m-r} U_{\alpha} U_{\alpha}'$ y $A_{12} A_{22}^{-1} A_{21}$ se distribuye como $\sum_{\alpha=m-r+1}^m U_{\alpha} U_{\alpha}'$ donde las U_{α} son independientes y $N(0, \Sigma_{11,2})$.

Se sigue de estos resultados que la distribución de $\rho_{ij, q+1, \dots, p}$ basado en las N observaciones es la misma que la de un coeficiente de correlación ordinario basado en $N - (p-q)$ observaciones con el correspondiente coeficiente de correlación poblacional dado por $\rho_{ij, q+1, \dots, p}$. Demuéstranoslo por fin el teorema

TEOREMA .- Si la función de distribución de r_{ij} basado en una muestra de tamaño N de una $N(0,1)$ con coeficiente de correlación poblacional ρ_{ij} la denotamos por $F(r|N, \rho_{ij})$, entonces la función de distribución del coeficiente de correlación parcial muestral $r_{ij.q+1...p}$ basado en una muestra de tamaño N de la misma normal es $F(r|N-(p-q), \rho_{ij.q+1...p})$.

Esta función de distribución fue obtenida por Fisher en el año 1924.

Test de hipótesis y regiones de confianza para coeficiente de correlación parcial.

Puesto que la distribución de un coeficiente de correlación parcial $r_{ij.q+1...p}$ basado en una muestra de tamaño N de una distribución con correlación poblacional $\rho_{ij.q+1...p}$ igual a cierto valor ρ , es la misma que la distribución de un coeficiente de correlación ordinario r basado en una muestra de tamaño $N-(p-q)$ de una distribución con el correspondiente coeficiente de correlación poblacional ρ . todos los procedimientos de inferencia estadística antes descritos pueden utilizarse ahora. Se actúa de forma semejante pero reemplazando N por $N-(p-q)$. Veamos dos ejemplos de esta actuación.

Ejemplo 1. Supongamos que sobre la base de una muestra de tamaño N deseamos obtener un intervalo de confianza para $\rho_{12.3}$. El coeficiente de correlación parcial muestral es $r_{12.3}$. El procedimiento será utilizar las tablas de David para $N-(p-q)$. Por ejemplo para la ilustración de la cosecha de heno antes dada, si queremos construir un intervalo para $\rho_{12.3}$ con coeficiente de confianza 0.95, como $r_{12.3} = 0.79$, utilizando la tabla $N-(p-q) = 20 - (2-1) = 19$, tenemos $0.52 < \rho_{12.3} < 0.92$.

Ejemplo 2. Análogamente, si queremos utilizar la z de Fisher para construir un contraste de hipótesis para $\rho_{ij.q+1...p} = \rho_0$, frente a una alternativa bilateral, tenemos

$$Z = \frac{1}{2} \log \frac{1 + r_{ij.q+1...p}}{1 - r_{ij.q+1...p}}$$

$$Z_0 = \frac{1}{2} \log \frac{1 + \rho_0}{1 - \rho_0}$$

Entonces $\sqrt{N-(p-q)-3} (Z - Z_0)$ se compara con los puntos de significación de una normal estandarizada.

En el ejemplo de la cosecha de heno, si quisiéramos contrastar la hipótesis $\rho_{12.3} = 0$, al nivel $\alpha = 0.05$, tendríamos $Z_0 = 0$ y $\sqrt{20-1-3} (0.0973) = 0.3892$. y como $|0.3892| < 1.96$ aceptamos la hipótesis.

4. EL COEFICIENTE DE CORRELACION MULTIPLE

Estimación de la correlación múltiple

El coeficiente de correlación múltiple poblacional lo estudiamos en capítulos anteriores. Por razones de brevedad estenderemos la estimación para el caso X_1 por un lado y X_2, \dots, X_p por otro, no pondremos índice alguno a R . Esta actuación no impone pérdida de generalidad por cuanto las variables pueden numerarse para hacer válido lo obtenido ahora en cualquier otra situación. Recordemos que la expresión para el coeficiente de correlación múltiple poblacional viene dada por

$$\bar{R} = \frac{\beta' \Sigma_{22} \beta}{\sqrt{\beta' \Sigma_{22} \beta}} = \sqrt{\frac{\beta' \Sigma_{22} \beta}{\sigma_{11}}} = \sqrt{\frac{\sigma_{(1)}' \Sigma_{22}^{-1} \sigma_{(1)}}{\sigma_{11}}}$$

donde $\beta, \sigma_{(1)}$ y Σ_{22} son

$$\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{(1)}' \\ \sigma_{(1)} & \Sigma_{22} \end{pmatrix} \quad \beta = \sigma_{(1)}' \Sigma_{22}^{-1}$$

Dada una muestra x_1, \dots, x_N ($N > p$) estimamos Σ mediante $S = (N/(N-1)) \hat{\Sigma}$ o

$$\hat{\Sigma} = \frac{1}{N} A = \frac{1}{N} \sum_{\alpha=1}^N (x_\alpha - \bar{x})(x_\alpha - \bar{x})' = \begin{pmatrix} \hat{\sigma}_{11} & \hat{\sigma}_{(1)}' \\ \hat{\sigma}_{(1)} & \hat{\Sigma}_{22} \end{pmatrix}$$

y estimamos β mediante $\hat{\beta} = \hat{\sigma}_{(1)}' \hat{\Sigma}_{22}^{-1} = a_{(1)}' A_{22}^{-1}$. Definimos entonces el coeficiente de correlación múltiple muestral como

$$R = \sqrt{\frac{\hat{\beta}' \hat{\Sigma}_{22} \hat{\beta}}{\hat{\sigma}_{11}}} = \sqrt{\frac{\hat{\sigma}_{(1)}' \hat{\Sigma}_{22}^{-1} \hat{\sigma}_{(1)}}{\hat{\sigma}_{11}}} = \sqrt{\frac{a_{(1)}' A_{22}^{-1} a_{(1)}}{a_{11}}}$$

Que resulta de la estimación máximo-variante de \bar{R} puede justificarse aplicando un teorema de inferencia presentado en capítulos anteriores referente a las estimaciones de funciones de los parámetros, y teniendo en cuenta que \bar{R} , $\hat{\sigma}_{(1)}$ y $\hat{\Sigma}_{22}$ pueden definirse mediante una transformación uno a uno de Σ . Otra expresión para R es

$$1 - R^2 = \frac{|\hat{\Sigma}|}{\hat{\sigma}_{11} |\hat{\Sigma}_{22}|} = \frac{|A|}{a_{11} |A_{22}|}$$

R y $\hat{\beta}$ tienen en la muestra las mismas propiedades que \bar{R} y β poseen en la población. Por ejemplo, de todos los vectores de dim. $(p-1)$ que definen una combinación lineal de $x_{\alpha}^{(2)}$ de las componentes de $x_{\alpha}^{(2)}$, el vector $d = \hat{\beta}$ es el que minimiza $\sum_{\alpha=1}^N [(x_{1\alpha} - \bar{x}_1) - d(x_{\alpha}^{(2)} - \bar{x}^{(2)})]^2$. En efecto, en primer lugar observemos que

$$\sum_{\alpha=1}^N [(x_{1\alpha} - \bar{x}_1) - \beta(x_{\alpha}^{(2)} - \bar{x}^{(2)})] [x_{\alpha}^{(2)} - \bar{x}^{(2)}]' = a_{(1)} - \hat{\beta} A_{22} = 0$$

por $\hat{\beta} = a_{(1)} A_{22}^{-1} A_{21}$

$$\begin{aligned} \sum_{\alpha=1}^N [(x_{1\alpha} - \bar{x}_1) - d(x_{\alpha}^{(2)} - \bar{x}^{(2)})]^2 &= \sum_{\alpha=1}^N \{ [(x_{1\alpha} - \bar{x}_1) - \hat{\beta}(x_{\alpha}^{(2)} - \bar{x}^{(2)})] + (\hat{\beta} - d)(x_{\alpha}^{(2)} - \bar{x}^{(2)}) \}^2 = \\ &= \sum_{\alpha=1}^N [(x_{1\alpha} - \bar{x}_1) - \hat{\beta}(x_{\alpha}^{(2)} - \bar{x}^{(2)})]^2 + (\hat{\beta} - d)' A_{22} (\hat{\beta} - d) \end{aligned}$$

y como A_{22} es definida positiva (excepto para muestras que surtan con probabilidad cero), el mínimo ocurrirá cuando $\hat{\beta} - d = 0$. Este mínimo es

$$\begin{aligned} \sum_{\alpha=1}^N [(x_{1\alpha} - \bar{x}_1) - \hat{\beta}(x_{\alpha}^{(2)} - \bar{x}^{(2)})]^2 &= a_{11} - 2\hat{\beta}' \sum_{\alpha=1}^N (x_{\alpha}^{(2)} - \bar{x}^{(2)})(x_{1\alpha} - \bar{x}_1)' + \hat{\beta}' A_{22} \hat{\beta}' = \\ &= a_{11} - 2\hat{\beta}' a_{(1)} + a_{(1)}' A_{22}^{-1} A_{22} a_{(1)} = a_{11} - 2a_{(1)}' A_{22}^{-1} a_{(1)} + a_{(1)}' A_{22}^{-1} a_{(1)} = \\ &= a_{11} - \hat{\beta}' A_{22} \hat{\beta}' = a_{11} - a_{(1)}' A_{22}^{-1} a_{(1)} = a_{11.2} \quad (\text{varianza condicional por } q=1) \end{aligned}$$

De este resultado puede darse una interesante interpretación geométrica. El vector de N componentes cuya α -ésima componente es $x_{1\alpha} - \bar{x}_1$ es, como en su momento vimos, la proyección del vector de las componentes i -ésimas de la muestra sobre el plano ortogonal a la línea equiangular. Tenemos por este vector $d(x_{\alpha}^{(2)} - \bar{x}^{(2)})$ la α -ésima componente de un vector en el hiperplano determinado por los $p-1$ últimos vectores. Como $\sum_{\alpha=1}^N [(x_{1\alpha} - \bar{x}_1) - d(x_{\alpha}^{(2)} - \bar{x}^{(2)})]^2$ es la distancia entre el primer vector y la combinación lineal de los $p-1$ últimos vectores, $\hat{\beta}(x_{\alpha}^{(2)} - \bar{x}^{(2)})$ es una combinación de un vector que minimiza esta distancia. La interpretación del hecho que $\sum_{\alpha=1}^N [(x_{1\alpha} - \bar{x}_1) - \hat{\beta}(x_{\alpha}^{(2)} - \bar{x}^{(2)})] [x_{\alpha}^{(2)} - \bar{x}^{(2)}]' = 0$ es que el vector con componentes $(x_{1\alpha} - \bar{x}_1) - \hat{\beta}(x_{\alpha}^{(2)} - \bar{x}^{(2)})$ es perpendicular a cada uno de los $p-1$ últimos vectores. Así resulta que el vector $\hat{\beta}(x_{\alpha}^{(2)} - \bar{x}^{(2)})$ es la proyección del primer vector sobre el hiperplano citado. El cuadrado de la longitud del vector proyección es

$$\sum_{\alpha=1}^N [\hat{\beta}(x_{\alpha}^{(2)} - \bar{x}^{(2)})]^2 = \hat{\beta}' A_{22} \hat{\beta}' = a_{(1)}' A_{22}^{-1} a_{(1)}$$

y la longitud del primer vector $\sum_{\alpha=1}^N (x_{1\alpha} - \bar{x}_1)^2 = a_{11}$. Así R es el coseno del ángulo formado por el primer vector con su proyección.

Vimos anteriormente que el coeficiente de correlación ordinario es el coseno del ángulo formado por los vectores involucrados en el plano ortogonal a la línea equiangular. Otra propiedad de R es que representa la máxima correlación entre $x_{1\alpha}$ y la combinación lineal de las componentes de $x_{\alpha}^{(2)}$. Esto responde en la propiedad geométrica de que R es el coseno del ángulo más pequeño entre el vector con componentes $(x_{1\alpha} - \bar{x}_1)$ y un vector en el hiperplano determinado por los otros $p-1$ vectores. Este resultado puede comprobarse sobre cualquier punto de vista, el geométrico citado o el analítico, apoyándose en el resultado precedente y de forma análoga a como lo hicimos con R .

Las interpretaciones geométricas vienen dadas todas ellas, entiendo de los vectores en el hiperplano $N-1$ dimensional ortogonal a la línea equiangular. Ya vimos anteriormente que el vector $(x_{i1} - \bar{x}_1, \dots, x_{iN} - \bar{x}_N)$ en este hiperplano podía ser designado por $(z_{i1}, \dots, z_{iN-1})$ donde las z_{ix} son las coordenadas referidas a un sistema de coordenadas $N-1$ dimensional en el hiperplano. Demostraremos que las nuevas coordenadas se obtienen a partir de las antiguas mediante la transformación $x_{ix} = \sum_{\alpha=1}^N z_{i\alpha} b_{\alpha x}$ donde B es una matriz ortogonal cuya última fila es $(1/\sqrt{N}, \dots, 1/\sqrt{N})$. Entonces

$$a_{ij} = \sum_{\alpha=1}^N (x_{i\alpha} - \bar{x}_i)(x_{j\alpha} - \bar{x}_j) = \sum_{\alpha=1}^{N-1} z_{i\alpha} z_{j\alpha}$$

En adelante cuando nos referamos a R definido anteriormente de z_{ix} lo haremos como al "coeficiente de correlación múltiple sin sustraer las medias".

El cálculo de R supone la extracción de la raíz cuadrada del cociente de $a_{(1)} A_{22}^{-1} a'_{(1)}$ por a_{11} . Puesto que A no depende directamente de los observaciones solo, probablemente, el cálculo de $a_{(1)} A_{22}^{-1} a'_{(1)}$ requiera alguna técnica especial.

Distribución del coeficiente de correlación múltiple muestral cuando el poblacional es cero.

Recordemos que

$$R^2 = \frac{a_{(1)} A_{22}^{-1} a'_{(1)}}{a_{11}}$$

entonces

$$1 - R^2 = 1 - \frac{a_{(1)} A_{22}^{-1} a'_{(1)}}{a_{11}} = \frac{a_{11} - a_{(1)} A_{22}^{-1} a'_{(1)}}{a_{11}} = \frac{a_{11.2}}{a_{11}}$$

7

$$\frac{R^2}{1 - R^2} = \frac{a_{(1)} A_{22}^{-1} a'_{(1)}}{a_{11.2}}$$

Para $g=1$ el resultado de la sección anterior (la referente al coeficiente de correlación parcial) establece que cuando $\beta=0$, es decir, cuando $R=0$, $a_{11.2}$ se distribuye como $\sum_{\alpha=1}^{N-p} V_{\alpha}^2$ y $a_{(1)} A_{22}^{-1} a'_{(1)}$ se distribuye como $\sum_{\alpha=N-p+1}^{N-1} V_{\alpha}^2$, donde las V_{α} son independientes y cada una con distribución $N(0, \sigma_{11.2})$. Entonces $a_{11.2}/\sigma_{11.2}$ y $a_{(1)} A_{22}^{-1} a'_{(1)}/\sigma_{11.2}$ se distribuyen independientemente como variables χ^2 con $N-p$ y $p-1$ grados de libertad respectivamente. Entonces

$$\frac{R^2}{1 - R^2} \cdot \frac{N-p}{p-1} = \frac{a_{(1)} A_{22}^{-1} a'_{(1)}/\sigma_{11.2}}{a_{11.2}/\sigma_{11.2}} \cdot \frac{N-p}{p-1} = \frac{\chi_{p-1}^2}{\chi_{N-p}^2} \cdot \frac{N-p}{p-1} = F_{p-1, N-p}$$

es decir, una F de Snedecor con $p-1$ y $N-p$ grados de libertad. La densidad de F viene dada por

$$\frac{\Gamma(\frac{1}{2}(N-1))}{\Gamma(\frac{1}{2}(p-1)) \Gamma(\frac{1}{2}(N-p))} \left(\frac{p-1}{N-p}\right)^{\frac{1}{2}(p-1)} \cdot x^{\frac{1}{2}(p-1)-1} \cdot \left(1 + \frac{p-1}{N-p} x\right)^{-\frac{1}{2}(N-1)}$$

la densidad de R , definido como

$$R = \sqrt{\frac{\frac{p-1}{N-p} x}{1 + \frac{p-1}{N-p} x}} \quad \text{con } x \text{ una } F_{p-1, N-p}$$

es

$$2 \frac{\Gamma(\frac{1}{2}(N-1))}{\Gamma(\frac{1}{2}(p-1)) \Gamma(\frac{1}{2}(N-p))} R^{p-2} (1 - R^2)^{\frac{1}{2}(N-p)-1}$$

Podemos unir todos estos resultados en el siguiente teorema.

TEOREMA - Sea R el coeficiente de inclusión múltiple múltiple entre X_1 y $X^{(2)'} = (X_2, \dots, X_p)$ basado en una muestra de tamaño N tomada de una $N(\mu, \Sigma)$. Si $\bar{R} = 0$ (es decir, si $(\sigma_{12}, \dots, \sigma_{1p}) = \sigma_{(1)} = 0 = \beta$), entonces $[R^2/(1-R^2)] \cdot [(N-p)/(p-1)]$ se distribuye como una F de Snedecor con $p-1$ y $N-p$ grados de libertad.

Hay que señalar que $p-1$ es el número de componentes de $X^{(2)'}$ y $N-p = N-(p-1)-1$. Si el número de componentes de $X^{(2)'}$ fuera q , las cantidades serían q y $N-q-1$ respectivamente.

Por otra parte la cantidad $R^2/(1-R^2)$ es la que se obtiene en teoría de la regresión para contrastar hipótesis acerca de los coeficientes de regresión, más exactamente, si la regresión de X_1 sobre X_2, \dots, X_p es nula.

Si $\bar{R} \neq 0$ la distribución de R es mucho más difícil de obtener. Lo veremos de esta manera más tarde.

Contrastes de hipótesis acerca de \bar{R} . ($H_0: \bar{R} = 0$)

Vamos ahora a considerar el problema estadístico de contrastar la hipótesis $\bar{R} = 0$ sobre la base de una muestra de tamaño N extraída de una población $N(0, \Sigma)$. Puesto que $\bar{R} \geq 0$ las alternativas consideradas son $\bar{R} > 0$.

Ostendremos en primer lugar el test de la razón de verosimilitud de esta hipótesis. La función de verosimilitud es

$$L(\mu^*, \Sigma^*) = \frac{1}{(2\pi)^{\frac{1}{2}Np} |\Sigma^*|^{\frac{1}{2}N}} \exp \left[-\frac{1}{2} \sum_{\alpha} (x_{\alpha} - \mu^*)' \Sigma^{*-1} (x_{\alpha} - \mu^*) \right].$$

Sea ω la región del espacio paramétrico especificado por la hipótesis nula. El criterio de la razón de verosimilitud se basa en

$$\lambda = \frac{\max_{\mu^*, \Sigma^* \in \omega} L(\mu^*, \Sigma^*)}{\max_{\mu^*, \Sigma^* \in \Omega} L(\mu^*, \Sigma^*)}.$$

Aquí Ω es el espacio de μ^*, Σ^* definida positiva y ω es la región de este espacio donde $\bar{R} = \sqrt{\sigma_{(1)} \Sigma_{22}^{-1} \sigma_{(1)}'} / \sqrt{\sigma_{11}} = 0$, es decir donde $\sigma_{(1)} \Sigma_{22}^{-1} \sigma_{(1)}' = 0$. Como Σ_{22} es definida positiva esta condición equivale a $\sigma_{(1)} = 0$. El máximo de $L(\mu^*, \Sigma^*)$ sobre Ω se alcanza para $\mu^* = \hat{\mu} = \bar{x}$ y $\Sigma^* = \hat{\Sigma} = (1/N) A$, y viene dado por

$$\max_{\mu^*, \Sigma^* \in \Omega} L(\mu^*, \Sigma^*) = \frac{N^{\frac{1}{2}pN} e^{-\frac{1}{2}pN}}{(2\pi)^{\frac{1}{2}pN} |A|^{\frac{1}{2}N}}.$$

En ω , la función de verosimilitud es

$$L(\mu^*, \Sigma^* | \sigma_{(1)} = 0) = \frac{1}{(2\pi)^{\frac{1}{2}N} |\sigma_{11}^*|^{\frac{1}{2}N}} \exp \left[-\frac{1}{2} \sum (x_{1\alpha} - \mu_1^*)^2 / \sigma_{11}^* \right] \cdot \frac{1}{(2\pi)^{\frac{1}{2}(p-1)N} |\Sigma_{22}^*|^{\frac{1}{2}N}} \exp \left[-\frac{1}{2} \sum (x_{\alpha}^{(2)} - \mu^{(2)*})' \Sigma_{22}^{*-1} (x_{\alpha}^{(2)} - \mu^{(2)*}) \right]$$

~~(2\pi)^{\frac{1}{2}N} |\sigma_{11}^*|^{\frac{1}{2}N} \exp \left[-\frac{1}{2} \sum (x_{1\alpha} - \mu_1^*)^2 / \sigma_{11}^* \right] \cdot \frac{1}{(2\pi)^{\frac{1}{2}(p-1)N} |\Sigma_{22}^*|^{\frac{1}{2}N}} \exp \left[-\frac{1}{2} \sum (x_{\alpha}^{(2)} - \mu^{(2)*})' \Sigma_{22}^{*-1} (x_{\alpha}^{(2)} - \mu^{(2)*}) \right]~~

El primer factor se maximiza para $\mu_1^* = \hat{\mu}_1 = \bar{x}_1$ y $\sigma_{11}^* = \hat{\sigma}_{11} = (1/N) a_{11}$, y el segundo factor se maximiza para $\mu^{(2)*} = \hat{\mu}^{(2)} = \bar{x}^{(2)}$ y $\Sigma_{22}^* = \hat{\Sigma}_{22} = (1/N) A_{22}$. El valor máximo de la función es

$$\max_{\mu^*, \Sigma^* \in \omega} L(\mu^*, \Sigma^*) = \frac{N^{\frac{1}{2}N} e^{-\frac{1}{2}N}}{(2\pi)^{\frac{1}{2}N} a_{11}^{\frac{1}{2}N}} \cdot \frac{N^{\frac{1}{2}(p-1)N} e^{-\frac{1}{2}(p-1)N}}{(2\pi)^{\frac{1}{2}(p-1)N} |A_{22}|^{\frac{1}{2}N}}.$$

Así, el criterio de la razón de verosimilitud será

$$\lambda = \frac{|A|^{\frac{1}{2}N}}{a_{11}^{\frac{1}{2}N} |A_{22}|^{\frac{1}{2}N}} = (1-R^2)^{\frac{1}{2}N}$$

La región crítica del test viene dada por $\lambda < \lambda_0$, donde λ_0 se elige de manera que la probabilidad de la desigualdad, cuando $\bar{R} = 0$, es precisamente α , nivel de significación elegido. Un test equivalente será

$$1 - \lambda^{\frac{2}{N}} = R^2 > 1 - \lambda_0^{\frac{2}{N}}.$$

Podemos construir un test a partir de $[R^2/(1-R^2)] \cdot [(N-p)/(p-1)]$, que es una función monótona de R . (11)
 Para ello, teniendo en cuenta que si $\bar{R}=0$ el variante anterior redistribuye como $F_{p-1, N-p}$, determinaremos la región crítica mediante la desigualdad

$$\frac{R^2}{1-R^2} \cdot \frac{N-p}{p-1} > F_{p-1, N-p}(\alpha)$$

donde $F_{p-1, N-p}(\alpha)$ es el punto de significación correspondiente a un nivel α . Este test es equivalente al obtenido mediante la razón de verosimilitud precisamente por ser una función monótona de R .

TEOREMA.- Dada una muestra x_1, \dots, x_N tomada de una $N(\mu, \Sigma)$, el test de la razón de verosimilitud a un nivel de significación α para la hipótesis $\bar{R}=0$, donde \bar{R} es el coeficiente de correlación múltiple entre X_1 y (X_2, \dots, X_p) , viene dado por

$$\frac{R^2}{1-R^2} \cdot \frac{N-p}{p-1} > F_{p-1, N-p}(\alpha)$$

donde R es el coeficiente de correlación múltiple muestral definido como antes.

Del hecho de que la densidad de R es monótona creciente en \bar{R} , como vemos a continuación, podemos asegurar que el test definido en el teorema es uniformemente más potente para contrastar $\bar{R}=0$ en la clase de los test que dependen de R . Como además R es invariante bajo transformaciones del tipo $x_{1\alpha}^* = Cx_{1\alpha} + d$ y $x_{\alpha}^{(2)*} = Cx_{\alpha}^{(2)} + d$ y es la única función del estadístico suficiente que es invariante, podemos concluir que el test es el test invariante uniformemente más potente.

Ejemplo.- Consideremos nuevamente el ejemplo de la cosecha de heno. El coeficiente de correlación múltiple muestral será, para los datos

$$1-R^2 = \frac{\begin{vmatrix} 1 & r_{12} & r_{13} \\ r_{12} & 1 & r_{23} \\ r_{13} & r_{23} & 1 \end{vmatrix}}{\begin{vmatrix} 1 & r_{23} \\ r_{23} & 1 \end{vmatrix}} = 0.357$$

y de aquí $R=0.802$. Para contrastar la hipótesis a un nivel $\alpha=0.01$ de que la cosecha de heno es independiente de la cantidad de lluvia y de la temperatura, compararemos $[R^2/(1-R^2)] \cdot [(20-3)/(3-1)] = 15.3$ con $F_{2,17}(0.01) = 6.11$ y concluiremos que no hay dependencia, pues rechazaremos la hipótesis nula $H_0: \bar{R}=0$.

Distintamente veremos que el test de independencia entre X_1 y $(X_2, \dots, X_p) = X^{(2)}$ es equivalente al test de si la regresión de X_1 sobre $X^{(2)}$ (es decir, el valor condicional esperado de X_1 dado $X^{(2)} = x^{(2)}$), que es $\mu_1 + \beta(x^{(2)} - \mu^{(2)})$, es nula en el sentido de que el vector de regresión $\beta=0$. $\hat{\beta} = a_{11}^{-1} A_{12}'$ es la estimación usual mínimo cuadrática de β , cuyo valor esperado es β y matriz de covarianzas $\sigma_{11}^{-1} A_{12} A_{12}'$ (cuando los $\Sigma_{\alpha}^{(2)}$ son fijos) y $a_{11}/(N-p)$ es la estimación usual de σ_{11} . Entonces

$$\frac{R^2}{1-R^2} \cdot \frac{N-p}{p-1} = \frac{\hat{\beta}' A_{12} \hat{\beta}}{a_{11}} \cdot \frac{N-p}{p-1}$$

es el estadístico F usual para contrastar la hipótesis de nulidad para el vector de coeficientes de regresión de X_1 sobre X_2, \dots, X_p .

Es interesante señalar, por último, que \bar{R} es la única función de μ y Σ que es invariante bajo cambios de posición y escala de las variables y bajo cambios que impliquen transformaciones lineales no singulares de $X^{(2)}$. Análogamente R es la única función de \bar{x} y $\hat{\Sigma}$, estadístico suficiente para μ y Σ , que es invariante bajo transformaciones similares.

Distribución del coeficiente de correlación múltiple muestral cuando el poblacional es diferente de 0.

Antes de encontrar la distribución de R cuando la hipótesis nula no es cierta. Enunciaremos que esta distribución depende sólo del coeficiente de correlación múltiple poblacional, \bar{R} .

Consideremos en primer lugar la distribución condicional de $R^2/(1-R^2) = a_{11}^{-1} A_{12}' A_{12} / a_{11}$ dado $Z_{\alpha}^{(2)} = \bar{Z}_{\alpha}^{(2)}$, $\alpha = 1, 2, \dots, n$. Bajo estas condiciones las $Z_{1\alpha}$ redistribuyen independientemente como $N(\beta Z_{\alpha}^{(2)}, \sigma_{11,2})$, donde $\beta = \sigma_{11}^{-1} \Sigma_{12}'$ y $\sigma_{11,2} = \sigma_{11} - \sigma_{11}^{-1} \Sigma_{12}' \Sigma_{12}^{-1} \Sigma_{12}$. Los condicionados son los de un teorema enunciado cuando estudiábamos las correlaciones parciales (véase párrafo 3 de este mismo capítulo), con $\gamma_{\alpha} = Z_{1\alpha}$, $\Gamma = \beta$, $W_{\alpha} = Z_{\alpha}^{(2)}$ ($r=p-1$), $\phi = \sigma_{11,2}$, $m=n$.

Entonces $A_{11.2} = A_{11} - A_{11}A_{22}^{-1}A_{11}'$ se distribuye con $\sum_{i=1}^n Y_i Y_i' - GHG'$ y en consecuencia $A_{11.2}/\sigma_{11.2}$ tiene una distribución χ^2 con $n-(p-1)$ grados de libertad. Por otra parte $A_{11}A_{22}^{-1}A_{11}' = (A_{01}A_{22}^{-1})A_{22}(A_{22}^{-1}A_{01}')$ se distribuye con GHG' y se distribuye como $\sum U_\alpha^2$ ($\alpha = n-(p-1)+1, \dots, n$) donde $\text{var}(U_\alpha) = \sigma_{11.2}$ y

$$E(U_{n-p+2}, \dots, U_n) = \Gamma \Gamma',$$

donde, a su vez, $\Gamma H \Gamma' = I$ ($H = F^{-1}(F')^{-1}$). Entonces $A_{11}A_{22}^{-1}A_{11}'/\sigma_{11.2}$ se distribuye como $\sum \alpha (U_\alpha/\sqrt{\sigma_{11.2}})^2$, donde $\text{var}(U_\alpha/\sqrt{\sigma_{11.2}}) = 1$

y

$$\sum_{\alpha} \left(\frac{E(U_\alpha)}{\sqrt{\sigma_{11.2}}} \right)^2 = \frac{1}{\sigma_{11.2}} \Gamma F^{-1} (\Gamma F^{-1})' = \frac{\Gamma H \Gamma'}{\sigma_{11.2}} = \frac{\beta A_{22} \beta'}{\sigma_{11.2}}$$

Así pues (condicionalmente) $A_{11}A_{22}^{-1}A_{11}'/\sigma_{11.2}$ tiene una distribución χ^2 no centrada con $p-1$ grados de libertad y parámetro de no centralización $\beta A_{22} \beta'/\sigma_{11.2}$. Podemos reconocer estos resultados en el siguiente teorema.

TEOREMA.- Sea R el coeficiente de relación múltiple múltiple entre X_1 y $X^{(2)} = (X_2, \dots, X_p)$ basado en N observaciones $(X_{11}, X_{11}^{(2)}), \dots, (X_{1N}, X_{1N}^{(2)})$. La distribución condicional de $[R^2/(1-R^2)] [(N-p)/(p-1)]$ dado $X_1^{(2)}$ fijo es una F de Student no centrada con $p-1$ y $N-p$ grados de libertad y parámetro de no centralización $(\beta A_{22} \beta')/\sigma_{11.2}$. La densidad condicional de $F = [R^2/(1-R^2)] [(N-p)/(p-1)]$ es de la forma

$$\frac{(p-1) e^{-\frac{1}{2} \beta A_{22} \beta' / \sigma_{11.2}}}{(N-p) \Gamma(\frac{1}{2}(N-p))} \sum_{\alpha=0}^{\infty} \frac{(\beta A_{22} \beta')^\alpha}{\alpha! \Gamma(\frac{1}{2}(p-1) + \alpha)} \left[\frac{(p-1)x}{N-p} \right]^{\frac{1}{2}(p-1) + \alpha - 1} \Gamma(\frac{1}{2}(N-1) + \alpha)$$

y la densidad condicional de $W = R^2$ es $(dx = [(N-p)/(p-1)] (1-w)^{-2} dw)$

$$\frac{e^{-\frac{1}{2} \beta A_{22} \beta' / \sigma_{11.2}}}{\Gamma(\frac{1}{2}(N-p))} (1-w)^{\frac{1}{2}(N-p) + 1} \sum_{\alpha=0}^{\infty} \frac{(\beta A_{22} \beta')^\alpha}{\alpha! \Gamma(\frac{1}{2}(p-1) + \alpha)} w^{\frac{1}{2}(p-1) + \alpha - 1} \Gamma(\frac{1}{2}(N-1) + \alpha)$$

Para obtener la densidad incondicional necesitamos multiplicar la expresión anterior por la marginal de $Z_1^{(1)} \dots Z_n^{(2)}$. Atendiendo a la marginal de W y $Z_1^{(2)}$, y luego integrar respecto al último conjunto de variables, lo que nos dará, finalmente, la marginal de W . Tenemos

$$\frac{\beta A_{22} \beta'}{\sigma_{11.2}} = \frac{\beta \sum_{\alpha=1}^n Z_\alpha^{(1)} Z_\alpha^{(2)} \beta'}{\sigma_{11.2}} = \sum_{\alpha=1}^n \left(\frac{\beta Z_\alpha^{(2)}}{\sqrt{\sigma_{11.2}}} \right)^2.$$

Puesto que la distribución de $Z_\alpha^{(1)}$ es $N(0, \sigma_{22})$, la distribución de $(\beta Z_\alpha^{(2)})/\sqrt{\sigma_{11.2}}$ es normal con media cero y varianzas

$$E \left(\frac{\beta Z_\alpha^{(2)}}{\sqrt{\sigma_{11.2}}} \right)^2 = \frac{E(\beta Z_\alpha^{(2)} Z_\alpha^{(2)} \beta')}{\sigma_{11.2}} = \frac{\beta \sigma_{22} \beta'}{\sigma_{11.2} - \beta \sigma_{22} \beta'} = \frac{\beta \sigma_{22} \beta' / \sigma_{11}}{1 - \beta \sigma_{22} \beta' / \sigma_{11}} = \frac{\bar{R}^2}{1 - \bar{R}^2}.$$

Así pues $(\beta A_{22} \beta' / \sigma_{11.2}) / [\bar{R}^2 / (1 - \bar{R}^2)]$ es una χ^2 con n grados de libertad. Ponemos $\bar{R}^2 / (1 - \bar{R}^2) = \phi$. Entonces $\beta A_{22} \beta' / \sigma_{11.2} = \phi \chi_n^2$. Calcularemos

$$E \left(e^{-\frac{1}{2} \phi \chi_n^2} \cdot \left(\frac{\phi \chi_n^2}{2} \right)^\alpha \right) = \frac{\phi^\alpha}{2^\alpha} \int_0^\infty u^\alpha e^{-\frac{1}{2} \phi u} \frac{1}{2^{\frac{1}{2}n} \Gamma(\frac{1}{2}n)} u^{\frac{1}{2}n-1} e^{-\frac{1}{2} \phi u} du = \frac{\phi^\alpha}{2^\alpha} \int_0^\infty \frac{1}{2^{\frac{1}{2}n} \Gamma(\frac{1}{2}n)} u^{\frac{1}{2}n+\alpha-1} e^{-\frac{1}{2} \phi u} du =$$

$$= \frac{\phi^\alpha}{(1+\phi)^{\frac{1}{2}n+\alpha}} \frac{\Gamma(\frac{1}{2}n+\alpha)}{\Gamma(\frac{1}{2}n)} \int_0^\infty \frac{1}{2^{\frac{1}{2}n+\alpha} \Gamma(\frac{1}{2}n+\alpha)} v^{\frac{1}{2}n+\alpha-1} e^{-\frac{1}{2} v} dv = \frac{\phi^\alpha}{(1+\phi)^{\frac{1}{2}n+\alpha}} \cdot \frac{\Gamma(\frac{1}{2}n+\alpha)}{\Gamma(\frac{1}{2}n)}.$$

Aplicando este resultado para nuestro desarrollo obtenemos una densidad para R^2 (Fisher 1928)

$$\frac{(1-R^2)^{\frac{1}{2}(n-p+1)} (1-\bar{R}^2)^{\frac{1}{2}n}}{\Gamma(\frac{1}{2}(n-p+1)) \Gamma(\frac{1}{2}n)} \sum_{\mu=0}^{\infty} \frac{(\bar{R}^2)^\mu (R^2)^{\frac{1}{2}(p-1)+\mu-1} \Gamma^2(\frac{1}{2}n+\mu)}{\mu! \Gamma(\frac{1}{2}(p-1)+\mu)}$$

- ① Se sospecha que los niveles de reneción de dos compuestos bioquímicos durante una situación de stress están correlacionados. Los propios biólogos que están implicados en este tipo de situaciones sugieren que la correlación, de existir, debe ser positiva. Para contrastar la hipótesis nula $H_0: \rho=0$ a un nivel $\alpha=0.01$ frente a la alternativa $H_1: \rho>0$ debemos diseñar un experimento basado en un número suficiente, N , de ensayos independientes de manera que una correlación poblacional tan pequeña como $\rho=0.20$ pueda ser detectada con una potencia 0.95. Los valores asimétricos de las probabilidades α y β reflejan las consecuencias relativas de tomar en cuenta una correlación falsa o pasar por alto una correlación real, aunque pequeña.

Utilizaremos como estadístico la z de Fisher. En estas condiciones la región crítica para la hipótesis nula viene dada por aquellos valores de z tales que

$$(z - \xi_0) \sqrt{N-3} \geq z_\alpha \rightarrow z \geq z_\alpha / \sqrt{N-3} + \xi_0$$

donde

$$z = \frac{1}{2} \log \frac{1+r}{1-r} \quad \gamma \quad \xi_0 = \frac{1}{2} \log \frac{1+\rho_0}{1-\rho_0} \quad \gamma \quad z_\alpha / P(z \geq z_\alpha) = \alpha.$$

en nuestro caso $\rho_0=0$ y además podemos llevar a cabo una aproximación mediante la normal, tal que z_α se calcule a partir de las tablas de la normal tipificada.

la potencia del test es la probabilidad de ~~rechazar~~ rechazar la hipótesis nula cuando es falsa, pero lo tanto, en nuestro caso y recurriendo nuevamente a la aproximación mediante la normal, dicha potencia viene dada por

$$1 - \beta(\rho) = 1 - \phi[z_\alpha + (\xi_0 - \xi) \sqrt{N-3}] \rightarrow \phi\left(\frac{z_\alpha / \sqrt{N-3} + \xi_0 - \xi}{1/\sqrt{N-3}}\right) = \phi(z_\alpha + (\xi_0 - \xi) \sqrt{N-3})$$

donde $\beta(\rho)$ es la probabilidad del error de tipo II, es decir, aceptar H_0 cuando es falsa y el verdadero valor es ρ .

Así las cosas, $z_{0.01} = 2.33$ (obtenido de unas tablas de la $N(0,1)$), $\xi_0 = 0$, $\xi = \frac{1}{2} \log \frac{1+\rho}{1-\rho}$ que para $\rho=0.20$, valor a partir del cual queremos detectar correlaciones con dicha potencia (0.95), vale $\xi = 0.2027$. Así pues la potencia del test debe de ratificar la condición

$$0.95 = 1 - \phi(2.33 - 0.2027 \sqrt{N-3})$$

es decir,

$$\phi(2.33 - 0.2027 \sqrt{N-3}) = 0.05$$

y como ϕ es la función de distribución de una $N(0,1)$, entonces

$$2.33 - 0.2027 \sqrt{N-3} = -1.645$$

y de aquí $N = 389$. Es decir, al menos 389 pares de observaciones deben recogerse en el estudio para que verifiquen las condiciones exigidas. El número es elevado, pero es el precio que debemos pagar para alcanzar los niveles α y β de probabilidad, tan rigurosos, deseados.

- ② Se ha aplicado determinado test de inteligencia a 933 individuos. Dicho test admite cuatro tipos de respuestas que se nos pueden considerar tantas variables a las que designaremos por G, comportamiento, V, comprensión verbal, A, edad y E, educación. Para estas cuatro variables se ha obtenido la siguiente matriz de correlaciones

$$\begin{array}{c} \begin{matrix} G & V & A & E \\ G & \begin{bmatrix} 1 & 0.72 & -0.44 & 0.60 \\ V & & 1 & -0.13 & 0.68 \\ A & & & 1 & -0.29 \\ E & & & & 1 \end{bmatrix} \end{matrix} \end{array} \xrightarrow{\text{matriz de variancias}} \begin{bmatrix} 126.07 & 116.44 & -53.98 & 20.85 \\ & 267.47 & -20.46 & 30.31 \\ & & 119.38 & -9.80 \\ & & & 9.58 \end{bmatrix}$$

Las desviaciones standard para cada uno de ellas fueron 11.228, 14.404, 10.926 y 3.095, respectivamente. Las correlaciones parciales cuando mantenemos constantes la edad y la educación vienen dadas por las siguientes matrices:

$$\begin{array}{c} \text{Edad Constante} \\ G \\ V \\ E \end{array} \begin{bmatrix} 1 & 0.74 & 0.55 \\ & 1 & 0.68 \\ & & 1 \end{bmatrix}$$

$$\begin{array}{c} \text{Educación constante} \\ G \\ V \\ A \end{array} \begin{bmatrix} 1 & 0.53 & -0.33 \\ & 1 & 0.10 \\ & & 1 \end{bmatrix}$$

Análogamente, la correlación parcial $r_{12.34}$ entre G y V cuando ambas, A y E permanecen constantes vale

$$r_{12.34} = 0.62.$$

Si queremos llevar a cabo contrastes de hipótesis del tipo $H_0: \rho = 0$ frente a $H_1: \rho \neq 0$, para $\alpha = 0.05$, podemos utilizar la t de student (ya sabemos que $\sqrt{N-2} r / \sqrt{1-r^2}$ es una t de student con $N-2$ grados de libertad). Ahora bien, como $N = 933$ supone una muestra de gran tamaño, podemos aproximar muy bien la t mediante la normal, haciendo caso omiso de la pérdida de grados de libertad (1 o 2 según los casos) ocasionada cuando mantenemos constantes 1 o 2 variables en el cálculo de las correlaciones parciales. En estas condiciones los valores críticos para r vienen dados por $r = \pm 0.064$, lo que supone que la hipótesis $H_0: \rho = 0$ sea rechazada cuando la apliquemos a cada uno de los coeficientes de correlación parcial obtenidos.

Calculemos ahora el coeficiente de correlación múltiple del comportamiento frente las otras tres variables y los coeficientes de regresión correspondientes, tenemos

$$\bar{R}_{1.234}^2 = 0.644 \quad \hat{\beta}' = [0.485, -0.346, 0.289]$$

los coeficientes de regresión vienen dados en las unidades de las variables originales.

Si queremos contrastar la hipótesis $H_0: \bar{R} = 0$ podemos hacer uso del hecho de que $(R^2 / (1-R^2)) \cdot ((N-P)/(P-1))$, bajo la hipótesis nula se distribuye como una F con $P-1$ y $N-P$ grados de libertad. Entonces para $\alpha = 0.05$, tenemos

$$F_{3,929}(0.05) \approx 2.60$$

Como

$$\frac{R^2}{1-R^2} \cdot \frac{N-P}{P-1} = \frac{0.644}{0.356} \cdot \frac{929}{3} = 560.18 > 2.60$$

rechazamos la hipótesis H_0 . Recordemos que este test es equivalente al efectuado para contrastar la hipótesis $\beta = 0$, aceptar que $R \neq 0$ implica aceptar que $\beta \neq 0$, lo que es lo mismo, admitir una dependencia de tipo lineal entre G y las otras variables, ~~es decir~~ cuando $\hat{\beta}'$ una estimación de los independientes coeficientes.

Capítulo 4

El estadístico T^2 generalizado

1. INTRODUCCION

Uno de los grupos de problemas, en estadística univariante, más importante es aquel que se ocupa de las cuestiones concernientes a la media de una distribución dada cuando la variancia de la distribución es desconocida. Sobre la base de una muestra se puede desear decidir si la media es igual a un número determinado de antemano, o bien dar un intervalo que contenga dicha media. El estadístico utilizado usualmente en estadística univariante es la diferencia entre la media muestral, \bar{x} , y la media hipotética poblacional dividida por la desviación típica muestral, s . Si la distribución muestral es $N(\mu, \sigma)$, entonces

$$t = \sqrt{N} \frac{\bar{x} - \mu}{s}$$

tiene una distribución t de Student con $N-1$ g.l., donde N es el número de observaciones en la muestra. Sobre la base de este hecho podemos llevar a cabo test de hipótesis $\mu = \mu_0$, donde μ_0 viene especificado, o bien podemos construir un intervalo de confianza para el parámetro desconocido μ .

El estadístico multivariante equivalente al análogo de t viene dado por

$$T^2 = N(\bar{x} - \mu)' S^{-1} (\bar{x} - \mu),$$

donde \bar{x} es el vector media muestral y S es la matriz de variancias muestrales. Los supuestos de este estadístico y ~~con~~ de su utilización para resolver los problemas, equivalentes en estadística multivariante, así como de otras utilidades del mismo. El primero en introducir el estadístico T^2 fue Hotelling en 1931 quien además estudió su distribución teórica bajo la hipótesis nula que se había planteado.

2. OBTENCIÓN Y DISTRIBUCIÓN DEL ESTADISTICO T^2 GENERALIZADO

Derivación del estadístico T^2 como una función del criterio de la razón de verosimilitud

Aunque el estadístico T^2 tiene varios usos, comenzaremos demostrando que el test de la razón de verosimilitud de la hipótesis $\mu = \mu_0$ sobre la base de una muestra de una población $N(\mu, \Sigma)$ se apoya precisamente en dicho estadístico T^2 . Supongamos que tenemos N observaciones x_1, \dots, x_N ($N > p$). La función de verosimilitud es

$$L(\mu, \Sigma^{-1}) = \frac{|\Sigma^{-1}|^{\frac{1}{2}N}}{(2\pi)^{\frac{1}{2}pN}} \exp \left[-\frac{1}{2} \sum_{\alpha=1}^N (x_\alpha - \mu)' \Sigma^{-1} (x_\alpha - \mu) \right]$$

El criterio de la razón de verosimilitud es

$$\lambda = \frac{\max_{\Sigma^{-1}} L(\mu_0, \Sigma^{-1})}{\max_{\mu, \Sigma^{-1}} L(\mu, \Sigma^{-1})}$$

Como ya sabemos el máximo para el denominador se alcanza substituyendo μ y Σ^{-1} por sus correspondientes estimaciones máximo-verosímiles

$$\hat{\mu}_0 = \bar{x}, \quad \hat{\Sigma}_0 = \frac{1}{N} \sum_{\alpha} (x_\alpha - \bar{x})(x_\alpha - \bar{x})'$$

Cuando $\mu = \mu_0$, la función de verosimilitud se maximiza en

$$\hat{\Sigma}_0 = \frac{1}{N} \sum_{\alpha} (x_\alpha - \mu_0)(x_\alpha - \mu_0)'$$

Como ya vimos en su momento. Aplicando un lema utilizado en el capítulo 2 referente a la obtención del máximo de la función de verosimilitud podemos afirmar que

$$\max_{\Sigma^{-1}, \mu} L(\mu, \Sigma^{-1}) = \frac{1}{(2\pi)^{\frac{1}{2}pN} |\hat{\Sigma}_0|^{\frac{1}{2}N}} e^{-\frac{1}{2}pN} \quad \gamma \quad \max_{\Sigma^{-1}} L(\mu_0, \Sigma^{-1}) = \frac{1}{(2\pi)^{\frac{1}{2}pN} |\hat{\Sigma}_0|^{-\frac{1}{2}N}} e^{-\frac{1}{2}pN}$$

Así, el criterio de la barra de similitud es

$$\lambda = \frac{|\hat{\Sigma}_x|^{1/2N}}{|\hat{\Sigma}_0|^{1/2N}} = \frac{|\sum (x_i - \bar{x})(x_i - \bar{x})'|^{1/2N}}{|\sum (x_i - \mu_0)(x_i - \mu_0)'|^{1/2N}} = \frac{|A|^{1/2N}}{|A + N(\bar{x} - \mu_0)(\bar{x} - \mu_0)'|^{1/2N}}$$

con $A = \sum (x_i - \bar{x})(x_i - \bar{x})' = (N-1)S$.

Para $|B| \neq 0$, utilizando una propiedad de las matrices, tenemos

$$\begin{vmatrix} B & C \\ D & E \end{vmatrix} = \begin{vmatrix} B & C \\ D & E \end{vmatrix} \cdot \begin{vmatrix} I & -B^{-1}C \\ 0 & I \end{vmatrix} = \begin{vmatrix} B & 0 \\ D & E - DB^{-1}C \end{vmatrix} = |B| \cdot |E - DB^{-1}C|$$

Aplicando dos veces este resultado a $\lambda^{2/N}$, tenemos

$$\begin{aligned} \lambda^{2/N} &= \frac{|A|}{|A + [\sqrt{N}(\bar{x} - \mu_0)][\sqrt{N}(\bar{x} - \mu_0)]'|} = \frac{|A|}{\begin{vmatrix} 1 & \sqrt{N}(\bar{x} - \mu_0)' \\ -\sqrt{N}(\bar{x} - \mu_0) & A \end{vmatrix}} = \frac{|A|}{\begin{vmatrix} A & \sqrt{N}(\bar{x} - \mu_0) \\ -\sqrt{N}(\bar{x} - \mu_0)' & 1 \end{vmatrix}} \\ &= \frac{|A|}{|A| \cdot |1 + N(\bar{x} - \mu_0)' A^{-1}(\bar{x} - \mu_0)|} = \frac{1}{|1 + N(\bar{x} - \mu_0)' A^{-1}(\bar{x} - \mu_0)|} = \frac{1}{1 + T^2/(N-1)} \end{aligned}$$

donde $T^2 = N(\bar{x} - \mu_0)' S^{-1}(\bar{x} - \mu_0) = (N-1) N(\bar{x} - \mu_0)' A^{-1}(\bar{x} - \mu_0)$.

El test de la barra de similitud se define mediante la región crítica (o región de rechazo)

$$\lambda \leq \lambda_0$$

donde λ_0 se elige de manera que la probabilidad de que $\lambda \leq \lambda_0$ cuando $\mu = \mu_0$ (hipótesis nula) sea igual al nivel de significación. Si operamos adecuadamente la anterior expresión es equivalente a

$$T^2 \geq T_0^2,$$

donde $T_0^2 = (N-1)(\lambda_0^{-2/N} - 1)$.

TEOREMA.- El test de la barra de similitud de la hipótesis $\mu = \mu_0$ para la distribución $N(\mu, \Sigma)$ viene dado por $T^2 \geq T_0^2$, donde T^2 se define como anteriormente y T_0^2 se elige de manera que $\Pr(T^2 \geq T_0^2)$, bajo la hipótesis nula, sea igual al nivel de significación elegido.

El test basado en la t de Student, para el caso univariante, tiene la propiedad de que cuando se contrasta la hipótesis $\mu = 0$ es invariante respecto a transformaciones de escala. Si la variable aleatoria escalar X se distribuye como $N(\mu, \sigma^2)$ entonces $X^* = cX$ se distribuye $N(c\mu, c^2\sigma^2)$ que está en la misma clase y la hipótesis $E(X) = 0$ equivale a $E(X^*) = cE(X) = 0$. Si las observaciones x_i se transforman análogamente, $x_i^* = cx_i$, entonces, para $c > 0$, t^* calculada a partir de x_i^* es la misma que la t calculada a partir de las x_i . Así, cualquiera que sea la unidad de medida el estadístico resulta ser el mismo.

El test T^2 generalizado tiene la misma propiedad. Si el vector aleatorio X se distribuye $N(\mu, \Sigma)$, entonces $X^* = GX$ (para $|G| \neq 0$) se distribuye $N(G\mu, G\Sigma G')$ que está en la misma clase. La hipótesis $E(X) = 0$ equivale a $E(X^*) = GE(X) = 0$. Si las observaciones x_i se transforman de la misma manera, $x_i^* = Gx_i$, entonces T^2 y T^{*2} son iguales.

Esta consecuencia de que $\bar{x}^* = G\bar{x}$, $A^* = GAG'$ y del siguiente lema:

LEMA.- Para cualquier matriz G , $p \times p$, no singular, H de k ejes característicos y K un vector cualquiera

$$K'H^{-1}K = (CK)'(CHK')^{-1}(CK).$$

Demostración.- $(CK)'(CHK')^{-1}(CK) = K'C'(C)^{-1}H^{-1}C'C K = K'H^{-1}K$.

demostraremos más adelante que de todos los test invariantes bajo estas transformaciones, este es uniformemente más potente. 4 (2)

Podemos dar ahora una interpretación geométrica de la raíz χ^2/N -sima del criterio de la razón de verosimilitud

$$\chi^2/N = \frac{|\sum (x_{\alpha} - \bar{x})(x_{\alpha} - \bar{x})'|}{|\sum (x_{\alpha} - \mu_0)(x_{\alpha} - \mu_0)'|}$$

entendidos de paralelepípedos. En una representación p -dimensional el numerador de χ^2/N es la suma de los cuadrados de los volúmenes de todos los paralelepípedos con ejes principales p vectores, cada uno con un extremo en el punto \bar{x} y el otro en x_{α} . El denominador es la suma de cuadrados de los volúmenes de los paralelepípedos cuyos ejes principales tienen un extremo en μ_0 y el otro en x_{α} . Si la suma de los cuadrados de los volúmenes que involucran vectores que emanan de \bar{x} , el auto de x_{α} , es mucho menor que la de los volúmenes que involucran vectores que emanan de μ_0 , entonces rechazamos la hipótesis de que μ_0 es la media de la distribución.

Existe también una interpretación en el caso de una representación N -dimensional. Sea $y_i = (x_{i1}, \dots, x_{in})$ el vector i -ésimo. Entonces

$$\sqrt{N} \bar{x}_i = \sum_{\alpha=1}^N \frac{1}{\sqrt{N}} x_{i\alpha}$$

es la distancia desde el origen de la proyección de y_i en la línea equiangular (con senso directores $\frac{1}{\sqrt{N}}, \frac{1}{\sqrt{N}}, \dots, \frac{1}{\sqrt{N}}$). Las coordenadas de la proyección son $(\bar{x}_i, \dots, \bar{x}_i)$. Entonces $(x_{i1} - \bar{x}_i, \dots, x_{in} - \bar{x}_i)$ es la proyección de y_i sobre el plano, a través del origen, perpendicular a la línea equiangular. El numerador de χ^2/N es el cuadrado del volumen p -dimensional del paralelepípedo con ejes principales, los vectores $(x_{i1} - \bar{x}_i, \dots, x_{in} - \bar{x}_i)$. Un punto $(x_{i1} - \mu_{01}, \dots, x_{in} - \mu_{0n})$ se obtiene a partir de y_i una traslación paralela a la línea equiangular (a una distancia $\sqrt{N} \mu_{0i}$). El denominador es el cuadrado del volumen de los paralelepípedos con ejes principales dichos vectores. Entonces χ^2/N es la relación de los cuadrados de dichos volúmenes.

La distribución de T^2

Entendamos en esta sección la distribución de T^2 bajo condiciones generales, incluido el caso en que la hipótesis nula no es cierta. Sea $T^2 = Y'S^{-1}Y$, donde Y se distribuye $N(\mu, \Sigma)$ y NS se distribuye independientemente como $\sum_{\alpha=1}^n Z_{\alpha} Z_{\alpha}'$ con Z_{α} independientes entre sí, cada una con una distribución $N(0, \Sigma)$. El T^2 antes utilizado es un caso especial de esta definición con $Y = \sqrt{N}(\bar{x} - \mu_0)$ y $U = \sqrt{N}(\mu - \mu_0)$ y $n = N-1$. Sea D una matriz no singular tal que $D\Sigma D' = I$, y definimos

$$Y^* = DY$$

$$S^* = DSD'$$

$$U^* = DU$$

Entonces $T^2 = Y^{*'} S^{*-1} Y^*$, de acuerdo con el lema anterior, ~~donde~~ Y^* se distribuye $N(U^*, I)$ y NS^* se distribuye independientemente como $\sum_{\alpha=1}^n Z_{\alpha}^* Z_{\alpha}^{*'} = \sum_{\alpha=1}^n DZ_{\alpha} (DZ_{\alpha})'$ con $Z_{\alpha}^* = DZ_{\alpha}$ independientes entre sí, cada una con distribución $N(0, I)$. Observamos que $U^{*'} \Sigma^{-1} U = U^{*'} (I)^{-1} U^* = U^{*'} U^*$.

Sea la primera fila de una matriz ortogonal Q , de dimensión $p \times p$, definida mediante

$$q_{1i} = \frac{Y_i^*}{\sqrt{Y^{*'} Y^*}} \quad i=1, \dots, p$$

lo que es permisible, por cuanto se ha de verificar $\sum_i q_{1i}^2 = 1$. Las otras $p-1$ filas pueden definirse de forma arbitraria de acuerdo con una propiedad de estas matrices. Como Q depende de las Y^* se trata de una matriz aleatoria. Sean ahora

$$U = QY^*$$

$$B = QNS^*Q'$$

De la definición de Q se sigue:

$$U_1 = \sum q_{1i} Y_i^* = \sqrt{Y^{*'} Y^*}$$

$$U_j = \sum q_{ji} Y_i^* = \sqrt{Y^{*'} Y^*} \sum q_{ji} q_{1i} = 0, \quad j \neq 1.$$

Entonces

$$\frac{T^2}{n} = U' B^{-1} U = (U_1 \ 0 \dots 0) \begin{bmatrix} b^{11} & b^{12} & \dots & b^{1p} \\ \vdots & \vdots & \ddots & \vdots \\ b^{p1} & b^{p2} & \dots & b^{pp} \end{bmatrix} \begin{pmatrix} U_1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = U_1^2 \cdot b^{11}$$

donde $(b^{ij}) = B^{-1}$. Haciendo uso de una propiedad de las matrices particionadas tenemos,

$$1/b^{11} = b_{11} - b_{(1)} \cdot B_{22}^{-1} b'_{(1)} = b_{11.2-p}$$

donde

$$B = \begin{pmatrix} b_{11} & b_{(1)} \\ b'_{(1)} & B_{22} \end{pmatrix},$$

y $T^2 = n U_1^2 / b_{11.2-p} = n Y^{*'} Y^* / b_{11.2-p}$. La distribución condicional de B dado Q es la de $\sum_{\alpha=1}^n V_{\alpha} V'_{\alpha}$, donde condicionalmente los $V_{\alpha} = Q Z_{\alpha}^*$ son independientes entre sí, cada uno con una distribución $N(0, I)$. Por otra parte, aplicando las propiedades obtenidas en el capítulo anterior para la distribución de las matrices A_{n-2} , cuando estudiamos el coeficiente de correlación parcial, sabemos que $b_{11.2-p}$ se distribuye condicionalmente como $\sum_{\alpha=1}^{n-(p-1)} W_{\alpha}^2$ donde condicionalmente los W_{α} son independientes entre sí, cada una $N(0, 1)$; es decir, $b_{11.2-p}$ se distribuye condicionalmente como una χ^2 con $n-(p-1)$ gr. de libertad. Puesto que la distribución condicional de $b_{11.2-p}$ no depende de Q , ello significa que se distribuye incondicionalmente como χ^2 . La cantidad $Y^{*'} Y^*$ tiene una distribución χ^2 no centrada con p grados de libertad y parámetro de no centralización $V^{*'} V^* = V' \Sigma^{-1} V$. En definitiva T^2/n se distribuye como la razón de una χ^2 no centrada y una χ^2 independientes.

TEOREMA.- Sea $T^2 = Y' S^{-1} Y$, donde Y se distribuye $N(V, \Sigma)$ y NS independientemente como $\sum_{\alpha=1}^n Z_{\alpha} Z'_{\alpha}$ con Z_{α} independientes entre sí, cada una $N(0, \Sigma)$. Entonces $(T^2/n) [(n-p+1)/p]$ se distribuye como una F no centrada con p y $n-p+1$ grados de libertad y parámetro de no centralización $V' \Sigma^{-1} V$. Si $V=0$ la distribución es una F centrada.

Clamaremos a esta distribución una distribución T^2 con n grados de libertad.

COROLARIO.- Sea x_1, \dots, x_N una muestra de una $N(\mu, \Sigma)$ y sea $T^2 = N(\bar{x} - \mu_0)' S^{-1} (\bar{x} - \mu_0)$. La distribución de $[T^2/(N-1)] [(N-p)/p]$ es una F no centrada con p y $N-p$ grados de libertad y parámetro de no centralización $N(\mu - \mu_0)' \Sigma^{-1} (\mu - \mu_0)$. Si $\mu = \mu_0$, entonces la distribución es una F centrada.

La densidad y tablas de la F no centrada se estudiarán en un párrafo posterior.

3. USOS DEL ESTADISTICO T^2

Ya hemos visto que el test de la unión de verosimilitud de la hipótesis $\mu = \mu_0$ sobre la base de una muestra de tamaño N de una $N(\mu, \Sigma)$ es equivalente a

$$T^2 \geq T_0^2.$$

Si el nivel de significación es α , elegimos T_0 a partir de unas tablas F centradas de manera que

$$T_0^2 = \frac{(N-1)p}{N-p} F_{p, N-p}(\alpha),$$

siendo $F_{p, N-p}(\alpha)$ el punto que deja a su derecha un área α . La elección del nivel de significación depende o puede hacerse de la potencia del test, de ello nos ocuparemos más tarde.

El estadístico T^2 se calcula con facilidad a partir de la muestra. En efecto

$$A^{-1}(\bar{x} - \mu_0) = b$$

es la solución de la ecuación

$$Ab = (\bar{x} - \mu_0),$$

entonces

$$\frac{T^2}{N-1} = N(\bar{x} - \mu_0)' b.$$

Observase que $T^2/(N-1)$ es la raíz no nula de

4 ③

$$|N(\bar{x} - \mu_0)(\bar{x} - \mu_0)' - \lambda A| = 0.$$

LEMA.- Si v es un vector de p componentes y si B es una matriz no singular de orden p , entonces $v'B^{-1}v$ es la raíz no nula de

$$|vv' - \lambda B| = 0.$$

Demostración.- La raíz λ_1 , no nula de la ecuación anterior está asociada con un vector característico β que satisface

$$vv'\beta = \lambda_1 B\beta.$$

Como $\lambda_1 \neq 0$, $v'\beta \neq 0$. Multiplicando por la izquierda por $v'B^{-1}$, tenemos

$$(v'B^{-1}v)(v'\beta) = \lambda_1(v'\beta),$$

lo que demuestra el lema. En el caso anterior $v = \sqrt{N}(\bar{x} - \mu_0)$ y $B = A$.

Una región de confianza para el vector media

Si se es la media de $N(\mu, \Sigma)$ sabemos que la probabilidad de extraer una muestra de tamaño N con media \bar{x} y matriz de covarianzas S , tal que

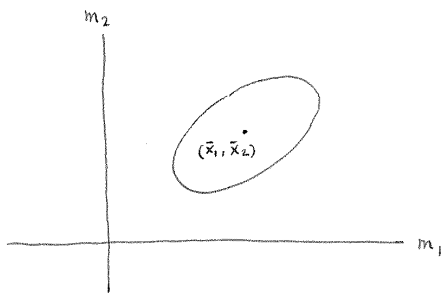
$$N(\bar{x} - \mu)' S^{-1} (\bar{x} - \mu) \leq T_0^2(\alpha)$$

es precisamente $1 - \alpha$.

Así pues, si calculamos la expresión anterior para una muestra particular, tenemos una confianza $1 - \alpha$ de que la citada expresión sea cierta en lo que a μ se refiere. La desigualdad

$$N(\bar{x} - \mu)' S^{-1} (\bar{x} - \mu) \leq T_0^2(\alpha)$$

es el interior y la frontera de un elipsoide p -dimensional con centro en \bar{x} y cuyo tamaño y forma dependen de S^{-1} y α . (Véase la figura para el caso bidimensional). Afirmamos que μ está en el elipsoide. lógicamente el elipsoide es aleatorio por cuanto depende de muestras aleatorias.



Problema de dos muestras

Otra situación en la cual se utiliza el estadístico T^2 es aquella en la que se pretende contrastar la hipótesis de la igualdad de medias para dos poblaciones normales multivariantes, cuyas matrices de covarianza se desconocen pero se suponen iguales. Supongamos $x_1^{(i)} \dots x_{N_1}^{(i)}$ es una muestra de $N(\mu^{(1)}, \Sigma)$, $i=1,2$. Querríamos contrastar la hipótesis $\mu^{(1)} = \mu^{(2)}$. $\bar{y}^{(i)}$ se distribuye $N(\mu^{(i)}, 1/N_i \Sigma)$. Consecuentemente $\sqrt{N_1 N_2 / (N_1 + N_2)} (\bar{y}^{(1)} - \bar{y}^{(2)})$ se distribuye como $N(\mu^{(1)} - \mu^{(2)}, \Sigma)$ y si la hipótesis nula $\mu^{(1)} = \mu^{(2)}$ es cierta esta distribución será $N(0, \Sigma)$. Si hacemos

$$S = \frac{1}{N_1 + N_2 - 2} \left\{ \sum_{\alpha=1}^{N_1} (y_\alpha^{(1)} - \bar{y}^{(1)}) (y_\alpha^{(1)} - \bar{y}^{(1)})' + \sum_{\alpha=1}^{N_2} (y_\alpha^{(2)} - \bar{y}^{(2)}) (y_\alpha^{(2)} - \bar{y}^{(2)})' \right\}$$

entonces $(N_1 + N_2 - 2)S$ se distribuye como $\sum_{\alpha=1}^{N_1 + N_2 - 2} Z_\alpha Z_\alpha'$ con $Z_\alpha \sim N(0, \Sigma)$. Así pues

$$T^2 = \frac{N_1 N_2}{N_1 + N_2} (\bar{y}^{(1)} - \bar{y}^{(2)})' S^{-1} (\bar{y}^{(1)} - \bar{y}^{(2)})$$

se distribuye como T^2 con $N_1 + N_2 - 2$ grados de libertad. La región crítica viene dada por

$$T^2 \geq \frac{(N_1+N_2-2)p}{N_1+N_2-p-1} F_{p, N_1+N_2-p-1}(\alpha) \quad \text{con nivel de significación } \alpha.$$

Un ejemplo de aplicación de esta teoría lo encontramos en Fisher (1936). Sea x_1 = longitud del sépalo, x_2 = anchura del sépalo, x_3 = longitud del pétalo y x_4 = anchura del pétalo. Se tienen 50 observaciones de la población Iris versicolor (1) y otras 50 de la población Iris setosa (2). Los datos resumidos, en centímetros, son los siguientes:

$$\bar{x}^{(1)} = \begin{pmatrix} 5.936 \\ 2.770 \\ 4.260 \\ 1.326 \end{pmatrix}, \quad \bar{x}^{(2)} = \begin{pmatrix} 5.006 \\ 3.428 \\ 1.462 \\ 0.246 \end{pmatrix}, \quad 98S = \begin{pmatrix} 19.1434 & 9.0356 & 4.7634 & 3.2394 \\ & 11.8658 & 4.6232 & 2.4746 \\ & & 12.2978 & 3.8744 \\ & & & 2.4604 \end{pmatrix}$$

El valor de $T^2/98 \approx 26.334$ y $T^2/98 \times 95/4 = 625.5$. Este valor es altamente significativo (comparado con el punto F para 4 y 95 grados de libertad), en consecuencia rechazamos la hipótesis de igualdad de medias en ambas poblaciones.

Un problema de q-muestras

Después de considerar el anterior ejemplo, Fisher obtiene una tercera muestra de una población a la que supone con igual matriz de covarianzas. Lleva a cabo 50 medidas, las mismas, sobre individuos Iris virginica. Existe una razón teórica para creer que la estructura genética de estas tres poblaciones es tal que sus vectores medios están relacionados como sigue:

$$3\mu^{(1)} = \mu^{(2)} + 2\mu^{(3)}$$

siendo $\mu^{(3)}$ el vector medio de la tercera población.

Se trata de un caso particular del siguiente problema general. Sean $\{x_\alpha^{(i)}\}$ $\alpha=1, \dots, N_i$, $i=1, \dots, q$ muestras extraídas de $N(\mu^{(i)}, \Sigma)$ $i=1, \dots, q$, respectivamente. Hagamos la siguiente hipótesis H :

$$\sum_{i=1}^q \beta_i \mu^{(i)} = \mu,$$

donde β_1, \dots, β_q son escalares dados y μ es un vector dado. El criterio a aplicar en esta situación es

$$T^2 = c \left(\sum_{i=1}^q \beta_i \bar{x}^{(i)} - \mu \right)' S^{-1} \left(\sum_{i=1}^q \beta_i \bar{x}^{(i)} - \mu \right),$$

$$\text{donde } \bar{x}^{(i)} = \frac{1}{N_i} \sum_{\alpha=1}^{N_i} x_\alpha^{(i)}, \quad \left(\sum_{i=1}^q N_i - q \right) S = \sum_{i=1}^q \sum_{\alpha=1}^{N_i} (x_\alpha^{(i)} - \bar{x}^{(i)}) (x_\alpha^{(i)} - \bar{x}^{(i)})', \quad 1/c = \sum_{i=1}^q \frac{\beta_i^2}{N_i}.$$

Esta T^2 tiene una distribución T^2 de Hotelling con $\sum_{i=1}^q N_i - q$ grados de libertad. Fisher, en realidad, supuso que las matrices de covarianzas podían ser diferentes y resolvió el problema aplicando una técnica de la que nos ocuparemos más tarde.

Un problema de simetría

Consideremos el siguiente contraste de hipótesis $H: \mu_1 = \mu_2 = \dots = \mu_p$ sobre la base de una muestra x_1, \dots, x_N extraída de una población $N(\mu, \Sigma)$ donde $\mu' = (\mu_1, \dots, \mu_p)$. Sea G cualquier matriz $(p-1) \times p$ de rango $p-1$ tal que

$$G\Sigma = 0$$

donde $\Sigma' = (1, \dots, 1)$. Entonces

$$y_\alpha = Gx_\alpha, \quad \alpha=1, \dots, N$$

tiene media $G\mu$ y matriz de covarianzas $G\Sigma G'$. La hipótesis H impone ahora $G\mu = 0$, en efecto de $G\Sigma = 0$ obtenemos

$$\sum_{j=1}^p c_{ij} \mu_j = 0, \quad i=1, \dots, p-1$$

pero $G\mu = \begin{bmatrix} \sum_{j=1}^p c_{1j} \mu_j \\ \vdots \\ \sum_{j=1}^p c_{p-1,j} \mu_j \end{bmatrix}$ si H se acepta $\sum_{j=1}^p c_{ij} \mu_j = \mu \left(\sum_{j=1}^p c_{ij} \right) = 0$, donde $\mu = \mu_1 = \mu_2 = \dots = \mu_p$.

$$T^2 = N \bar{y}' S^{-1} \bar{y}$$

$$\text{con } \bar{y} = \frac{1}{N} \sum_{\alpha=1}^N y_{\alpha} = G \bar{x}, \quad S = \frac{1}{N-1} \sum_{\alpha} (y_{\alpha} - \bar{y})(y_{\alpha} - \bar{y})' = \frac{1}{N-1} G \left[\sum_{\alpha} (x_{\alpha} - \bar{x})(x_{\alpha} - \bar{x})' \right] G'$$

Este estadístico tiene una distribución T^2 con $N-1$ grados de libertad para una distribución $(p-1)$ -dimensional. Este estadístico T^2 es invariante bajo cualquier transformación lineal en las $p-1$ dimensiones que sea ortogonal a E . En consecuencia el estadístico es independiente de la elección que se haga para G .

Un ejemplo que ilustra este problema viene dado por Rao (1948). Sea N la cantidad de arboles extraídos, de un árbol productor de urcho, en la ~~zona~~ ^{misma} ~~parte~~ ^{zona} Norte del ~~misma~~ ^{misma} ~~parte~~ ^{zona}, sean E, S y W definidas de forma similar para las ~~caras~~ ^{caras} Este, Sur y Oeste respectivamente. El conjunto de estas cuatro cantidades para un mismo árbol se considera una observación de una normal 4-variante. La cuestión a dilucidar es si los árboles tienen la misma cantidad de urcho en cada una de sus caras. Hacemos la siguiente transformación,

$$y_1 = N - E - W + S$$

$$y_2 = S - W$$

$$y_3 = N - S.$$

El número de observaciones es de 28.

El vector de medias es

$$\bar{y} = \begin{bmatrix} 8.86 \\ 4.50 \\ 0.86 \end{bmatrix}$$

La matriz de varianzas para y es

$$S = \begin{bmatrix} 128.72 & 61.41 & -21.02 \\ & 56.93 & -28.30 \\ & & 63.53 \end{bmatrix}$$

El valor de $T^2/(N-1)$ es 0.768. El estadístico $0.768 \cdot 25/3 = 6.402$ es una F con 3, 25 grados de libertad. Para $\alpha = 0.01$ $F_{3,25} = 4.68$ y rechazamos la hipótesis.

4. LA DISTRIBUCIÓN DE T^2 BAJO HIPÓTESIS ALTERNATIVAS: LA FUNCIÓN POTENCIA.

Hemos visto anteriormente que $(T^2/n)/(N-p)/p$ tiene una distribución F no centrada. Se trata ahora de estudiar las distribuciones χ^2 y F no centradas, la tabulación de estas últimas y su aplicación a procedimientos basados en T^2 .

La distribución χ^2 es la distribución de la suma de cuadrados de variables aleatorias (scalars) independientes normales con media 0 y varianzas 1; la χ^2 no centrada es la generalización de la anterior cuando las medias pueden ser distintas de cero. Sea Y un vector de p componentes distribuido $N(\nu, I)$. Sea Q una matriz ortogonal cuya primera fila viene dada por

$$q_{1i} = \frac{\nu_i}{\sqrt{\nu' \nu}}.$$

Entonces $Z = QY$ se distribuye $N(\lambda, I)$, donde $\lambda' = [2, 0, \dots, 0]$ y $Z = \sqrt{\nu' \nu}$, como ya vimos anteriormente (apartado 2 de este capítulo). Sea $V = Y'Y = Z'Z = \sum_{i=1}^p Z_i^2$, entonces $W = \sum_{i=2}^p Z_i^2$ tiene una distribución χ^2 con $p-1$ grados de libertad y Z_1 y W tienen una densidad conjunta

$$\begin{aligned} & \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(z_1 - z)^2} \cdot \frac{1}{2^{\frac{1}{2}(p-1)} \Gamma(\frac{1}{2}(p-1))} w^{\frac{1}{2}(p-1)-1} e^{-\frac{1}{2}w} = c e^{-\frac{1}{2}(z_1^2 + z_1^2 + w)} w^{\frac{1}{2}(p-3)} e^{-\frac{1}{2}w} = \\ & = c e^{-\frac{1}{2}(z_1^2 + z_1^2 + w)} w^{\frac{1}{2}(p-3)} \sum_{\alpha=0}^{\infty} \frac{z_1^{\alpha} z_1^{\alpha}}{\alpha!} \end{aligned}$$

donde $c' = 2^{\frac{1}{2}p} \sqrt{\pi} \Gamma(\frac{1}{2}(p-1))$. La densidad conjunta de $V = W + Z_1^2$ y Z_1 se obtiene substituyendo $w = V - Z_1^2$ (el jacobiano

de la transformación siendo la 1,

$$C e^{-\frac{1}{2}(z^2+u)} (v-z_1^2)^{\frac{1}{2}(p-3)} \sum_{\alpha=0}^{\infty} \frac{z_1^\alpha z_1^\alpha}{\alpha!}.$$

la densidad conjunta de V y $U = Z_1/\sqrt{V}$ es ($dz_1 = \sqrt{v} du$)

$$C e^{-\frac{1}{2}(z^2+u)} v^{\frac{1}{2}(p-2)} (1-u^2)^{\frac{1}{2}(p-3)} \sum_{\alpha=0}^{\infty} \frac{z^\alpha u^\alpha v^{\frac{1}{2}\alpha}}{\alpha!}$$

respecto al campo de variación de z_1 dado v , $[-\sqrt{v}, \sqrt{v}]$, el de u , $[-1, 1]$. Cuando integramos en esta última expresión respecto de u , término a término, los términos para α impar dan una integral nula por ser los términos a la vez de una función impar de u . En las otras integraciones haciendo el cambio $u = \sqrt{s}$ ($du = \frac{1}{2} ds/\sqrt{s}$), obtenemos

$$\int_{-1}^1 (1-u^2)^{\frac{1}{2}(p-3)} u^{2\beta} du = 2 \int_0^1 (1-u^2)^{\frac{1}{2}(p-3)} u^{2\beta} du = \int_0^1 (1-s)^{\frac{1}{2}(p-3)} s^{\beta-\frac{1}{2}} ds = B[\frac{1}{2}(p-1), \beta+\frac{1}{2}] =$$

$$= \frac{\Gamma[\frac{1}{2}(p-1)] \cdot \Gamma(\beta+\frac{1}{2})}{\Gamma(\frac{1}{2}p+\beta)}$$

de acuerdo con la relación existente entre las funciones beta y gamma.

la densidad de V viene pues dada por:

$$\frac{1}{2^{\frac{1}{2}p} \sqrt{\pi}} e^{-\frac{1}{2}(z^2+v)} v^{\frac{1}{2}p-1} \sum_{\beta=0}^{\infty} \frac{(z^2)^{\beta} v^{\beta}}{(2\beta)!} \cdot \frac{\Gamma(\beta+\frac{1}{2})}{\Gamma(\frac{1}{2}p+\beta)}.$$

Esta es la densidad de una χ^2 no centrada con p grados de libertad y parámetro de no centralización z^2 .

TEOREMA - Si el vector Y de p componentes se distribuye $N(0, I)$, entonces $V = Y'Y$ tiene una densidad χ^2 no centrada con p grados de libertad, siendo $z^2 = v'/v$ el parámetro de no centralización.

Sea ahora V una χ^2 no centrada con p grados de libertad y parámetro de no centralización z^2 , sea W , independiente de la anterior, una χ^2 con m grados de libertad. Vamos a encontrar la densidad de $F = (V/p)/(W/m)$, que es una F no centrada con parámetro z^2 . La densidad conjunta de V y W es la densidad de V multiplicada por la de W , que es $2^{-\frac{1}{2}m} \pi^{-\frac{1}{2}} (\frac{1}{2}m)^{\frac{1}{2}m-1} e^{-\frac{1}{2}w}$. La densidad conjunta de F y W ($dw = pm df/m$) es

$$\frac{e^{-\frac{1}{2}z^2}}{2^{\frac{1}{2}(p+m)} \sqrt{\pi} \Gamma(\frac{1}{2}m)} e^{-\frac{1}{2}w(1+pf/m)} \cdot \frac{p}{m} \sum_{\beta=0}^{\infty} \frac{(z^2)^{\beta} \Gamma(\beta+\frac{1}{2})}{(2\beta)! \Gamma(\frac{1}{2}p+\beta)} \cdot \left(\frac{pf}{m}\right)^{\frac{1}{2}p+\beta-1} w^{\frac{1}{2}(p+m)+\beta-1}.$$

la densidad marginal se obtiene integrando para w entre 0 y $+\infty$,

$$\frac{p e^{-\frac{1}{2}z^2}}{m \sqrt{\pi} \Gamma(\frac{1}{2}m)} \sum_{\beta=0}^{\infty} \frac{(2e^{\frac{1}{2}z^2})^{\beta} \Gamma(\beta+\frac{1}{2})}{(2\beta)! \Gamma(\frac{1}{2}p+\beta)} \left(\frac{pf}{m}\right)^{\frac{1}{2}p+\beta-1} \frac{\Gamma(\frac{1}{2}(p+m)+\beta)}{(1+pf/m)^{\frac{1}{2}(p+m)+\beta}}.$$

Si utilizamos la fórmula de la duplicación de la función gamma

$$(2\beta)! = \Gamma(2\beta+1) = \Gamma(\beta+\frac{1}{2}) \Gamma(\beta+\frac{1}{2}) 2^{2\beta}/\sqrt{\pi}$$

la densidad queda finalmente

$$\frac{p e^{-\frac{1}{2}z^2}}{m \Gamma(\frac{1}{2}m)} \sum_{\beta=0}^{\infty} \frac{(z^2/2)^{\beta} (pf/m)^{\frac{1}{2}p+\beta-1} \Gamma(\frac{1}{2}(p+m)+\beta)}{\beta! \Gamma(\frac{1}{2}p+\beta) (1+pf/m)^{\frac{1}{2}(p+m)+\beta}}$$

que es la densidad de una F no centrada con p y m grados de libertad respectivamente y parámetro de no centralización z^2 .

Si $T^2 = N(\bar{x} - \mu_0)' S^{-1} (\bar{x} - \mu_0)$ se basa en una muestra de tamaño N de una $N(\mu, \Sigma)$, entonces $(T^2/n) \left(\frac{N-p}{p} \right)$ 4 (5) tiene una distribución F no centrada con p y $N-p$ grados de libertad y parámetro de no centralización $N(\mu - \mu_0)' \Sigma^{-1} (\mu - \mu_0) = Z^2$. La expresión de la densidad de T^2 será

$$\frac{e^{-\frac{1}{2}Z^2}}{(N-1)^p \left(\frac{1}{2}(N-p) \right)} \sum_{\beta=0}^{\infty} \frac{(Z^2/2)^{\beta} \left[\frac{1}{2}(N-1) \right]^{\frac{1}{2}p+\beta-1} \Gamma\left(\frac{1}{2}N+\beta\right)}{\beta! \Gamma\left(\frac{1}{2}p+\beta\right) \left[1 + \frac{1}{2}(N-1) \right]^{\frac{1}{2}N+\beta}}$$

Zang (1938) ha obtenido tablas de la probabilidad de aceptar la hipótesis nula (es decir, la probabilidad del error del tipo II) para varios valores de Z^2 y para niveles de significación 0.05 y 0.01. ~~El~~ número de grados de libertad f_1 que allí aparece es el p aquí considerado $[1, (1) 8]$, y f_2 es $n-p+1$ $[2, 4 (1) 30, 60, \infty]$, el parámetro de no centralización ϕ está relacionado con Z^2 mediante la expresión

$$\phi = \frac{Z^2}{p+1} \quad [1 (1) 3 (1) 8].$$

A modo de ejemplo, supongamos que $p=4$, $n-p+1=20$ y que queremos contrastar la hipótesis $\mu=0$ a un nivel 1%. Desearíamos conocer la probabilidad de aceptar la hipótesis nula cuando $\phi=2.5$ ($Z^2=31.25$). Dicha probabilidad es 0.227. Si pensamos que la desventaja de aceptar la hipótesis nula cuando N, μ, Σ son tales que $Z^2=31.25$ es menor que la desventaja de rechazarla cuando es cierta, entonces podemos encontrar razonable aceptar el test en estas condiciones. Si embargo, si la desventaja del error del tipo I es aproximadamente igual que la del tipo II, parece razonable desear disminuir la probabilidad del error del tipo II. Para ello, adoptando un nivel $\alpha=5\%$, la probabilidad de un error del tipo II (para $\phi=2.5$) es solo 0.043.

Existen también tablas (Emma Lehmer, 1944) que proporcionan ϕ para α dado y para una probabilidad de error del tipo II dada. Estas tablas son útiles para ver que valor de Z^2 es necesario para hacer la probabilidad de aceptar la hipótesis nula suficientemente baja cuando $\mu \neq 0$. Por ejemplo, si queremos ser capaces de rechazar la hipótesis $\mu=0$ sobre la base de una muestra para unos μ y Σ dados, podemos ser capaces de elegir N de manera que $N\mu'\Sigma^{-1}\mu = Z^2$ sea suficientemente grande. Otra bien, la dificultad con estas consideraciones es que normalmente no conocemos los valores de μ y Σ (y por tanto, Z^2) para los cuales queremos que la probabilidad de rechazo tenga un cierto valor.

J. ALGUNAS PROPIEDADES OPTIMAS DEL TEST T^2

Supongamos que queremos contrastar la hipótesis $\mu=0$ sobre la base de N observaciones x_1, \dots, x_N de una $N(\mu, \Sigma)$. Consideremos en primer lugar la clase de los tests basados en los estadísticos $A = [(\bar{x} - \bar{x})(\bar{x} - \bar{x})']$ y \bar{x} que son invariantes con respecto a transformaciones del tipo $A^* = CAC'$ y $\bar{x}^* = C\bar{x}$, donde C es no singular. La transformación $x_i^* = Cx_i$ deja el problema invariante; esto es, estamos de x_i^* contrastando la hipótesis $E(x_i^*)=0$ dado que x_1^*, \dots, x_N^* son N observaciones procedentes de una población normal multivariante. Parece razonable buscar una solución que sea también invariante con respecto a estas transformaciones; es decir, buscaremos una región crítica que no se altere mediante una transformación lineal no singular (la definición de la región es la misma en diferentes sistemas de coordenadas).

TEOREMA. - Dadas las observaciones x_1, \dots, x_N procedentes de una $N(\mu, \Sigma)$, de todos los tests de $\mu=0$ basados en \bar{x} , $A = [(\bar{x} - \bar{x})(\bar{x} - \bar{x})']$ que son invariantes con respecto a transformaciones del tipo $\bar{x}^* = C\bar{x}$, $A^* = CAC'$ (C no singular), el test T^2 es uniformemente más potente.

Demostración. - Como vimos anteriormente cualquier test basado en T^2 es invariante. Por otra parte, esta función es la única invariante para la que si $f(\bar{x}, A)$ es invariante, entonces $f(\bar{x}^*, A^*) = f(\bar{x}^*, I)$, donde la única coordenada de \bar{x}^* diferente de cero es la primera y su valor es $\sqrt{\bar{x}' A^* \bar{x}}$ (existe una matriz G , no singular, tal que $C\bar{x} = \bar{x}^* \sqrt{CAC'} = I$). Así $f(\bar{x}, A)$ depende solo de $\bar{x}' A^* \bar{x}$. Por tanto un test invariante debe de estar basado en $\bar{x}' A^* \bar{x}$. Finalmente, podemos aplicar el lema fundamental de Neyman-Pearson a la distribución de T^2 para encontrar el test uniformemente más potente basado en T^2 , frente a la alternativa simple $Z^2 = N\mu'\Sigma^{-1}\mu$. El test uniformemente más potente ~~de~~ de $Z^2=0$ está basado en el cociente de las distribuciones de T^2 cuando $Z^2 = N\mu'\Sigma^{-1}\mu$ y cuando $Z^2=0$. Dicho test es

Si $T^2 = N(\bar{x} - \mu_0)' S^{-1} (\bar{x} - \mu_0)$ se basa en una muestra de tamaño N de una $N(\mu, \Sigma)$, entonces $(T^2/n) \left(\frac{N-p}{p} \right)$ 4 (5) tiene una distribución F no centrada con p y $N-p$ grados de libertad y parámetro de no centralización $N(\mu - \mu_0)' \Sigma^{-1} (\mu - \mu_0) = Z^2$. La expresión de la densidad de T^2 será

$$\frac{e^{-\frac{1}{2}Z^2}}{(N-1)^p \left(\frac{1}{2}(N-p) \right)} \sum_{\beta=0}^{\infty} \frac{(Z^2/2)^{\beta} [t^2/(N-1)]^{\frac{1}{2}p+\beta-1} \Gamma\left(\frac{1}{2}N+p\right)}{\beta! \Gamma\left(\frac{1}{2}p+\beta\right) [1+t^2/(N-1)]^{\frac{1}{2}N+p}}$$

Zang (1938) ha obtenido tablas de la probabilidad de aceptar la hipótesis nula (es decir, la probabilidad del error del tipo II) para varios valores del Z^2 y para niveles de significación 0.05 y 0.01. ~~El~~ número de grados de libertad f_1 que allí aparece es el p aquí considerado $[1(1)8]$, y f_2 es $n-p+1$ $[2, 4(1)30, 60, \infty]$, el parámetro de no centralización ϕ está relacionado con Z^2 mediante la expresión

$$\phi = \frac{Z^2}{p+1} \quad [1(1)3(1)8].$$

A modo de ejemplo, supongamos que $p=4$, $n-p+1=20$ y que queremos contestar la hipótesis $\mu=0$ a un nivel 1%. Desearíamos conocer la probabilidad de aceptar la hipótesis nula cuando $\phi=2.5$ ($Z^2=31.25$). Dicha probabilidad es 0.227. Si pensamos que la desventaja de aceptar la hipótesis nula cuando N, μ, Σ cambian que $Z^2=31.25$ es menor que la desventaja de rechazarla cuando es cierta, entonces podemos considerar razonable aceptar el test en estas condiciones. Si embargo, si la desventaja del error del tipo I es aproximadamente igual que la del tipo II, parece razonable desear disminuir la probabilidad del error del tipo II. Para ello, adoptando un nivel $\alpha=5\%$, la probabilidad de un error del tipo II (para $\phi=2.5$) es sólo 0.043.

Existen también tablas (Emma Lehmer, 1944) que proporcionan ϕ para α dado y para una probabilidad de error del tipo II dada. Estas tablas son útiles para ver que valor del Z^2 es necesario para hacer la probabilidad de aceptar la hipótesis nula suficientemente baja cuando $\mu \neq 0$. Por ejemplo, si queremos ser capaces de rechazar la hipótesis $\mu=0$ sobre la base de una muestra para unos μ y Σ dados, podemos ser capaces de elegir N de manera que $N\mu'\Sigma^{-1}\mu = Z^2$ sea suficientemente grande. Ahora bien, la dificultad en estas consideraciones es que normalmente no conocemos los valores de μ y Σ (y por tanto, Z^2) para los cuales queremos que la probabilidad de rechazo tenga un cierto valor.

J. ALGUNAS PROPIEDADES OPTIMAS DEL TEST T^2

Supongamos que queremos contestar la hipótesis $\mu=0$ sobre la base de N observaciones x_1, \dots, x_N de una $N(\mu, \Sigma)$. Consideremos en primer lugar la clase de los tests basados en los estadísticos $A = [(\bar{x} - \bar{x})(\bar{x} - \bar{x})']$ y \bar{x} que son invariantes con respecto a transformaciones del tipo $A^* = CAC'$ y $\bar{x}^* = C\bar{x}$, donde C es no singular. La transformación $x_i^* = Cx_i$ deja el problema invariante; esto es, estamos de x_i^* contestando la hipótesis $E(x_i^*) = 0$ dado que x_1^*, \dots, x_N^* son N observaciones procedentes de una población normal multivariante. Parece razonable buscar una solución que sea también invariante con respecto a estas transformaciones; es decir, buscaremos una región crítica que no se altere mediante una transformación lineal no singular (la definición de la región es la misma en diferentes sistemas de coordenadas).

TEOREMA. - Dadas las observaciones x_1, \dots, x_N procedentes de una $N(\mu, \Sigma)$, de todos los tests de $\mu=0$, basados en \bar{x} , $A = [(\bar{x} - \bar{x})(\bar{x} - \bar{x})']$ que son invariantes con respecto a transformaciones del tipo $\bar{x}^* = C\bar{x}$, $A^* = CAC'$ (C no singular), el test T^2 es uniformemente más potente.

Demostración. - Como vimos anteriormente cualquier test basado en T^2 es invariante. Por otra parte, esta función es la única invariante para la que si $f(\bar{x}, A)$ es invariante, entonces $f(\bar{x}, A) = f(\bar{x}^*, I)$, donde la única coordenada de \bar{x}^* diferente de cero es la primera y su valor es $\sqrt{\bar{x}' A^{-1} \bar{x}}$ (Existe una matriz G , no singular, tal que $C\bar{x} = \bar{x}^* \sqrt{CAC' = I}$). Así $f(\bar{x}, A)$ depende sólo de $\bar{x}' A^{-1} \bar{x}$. Por tanto un test invariante debe de estar basado en $\bar{x}' A^{-1} \bar{x}$. Finalmente, podemos aplicar el lema fundamental de Neyman-Pearson a la distribución de T^2 para encontrar el test uniformemente más potente basado en T^2 , frente a la alternativa simple $Z^2 = N\mu'\Sigma^{-1}\mu$. El test uniformemente más potente ~~de~~ de $Z^2=0$ está basado en el cociente de las distribuciones de T^2 cuando $Z^2 = N\mu'\Sigma^{-1}\mu$ y cuando $Z^2=0$. Dichos test es

$$C < e^{-\frac{1}{2}t^2} \sum_{\alpha=0}^{\infty} \frac{(t^2/2)^{\alpha} (t^2/n)^{\frac{1}{2}p+\alpha-1} (1+t^2/n)^{-\frac{1}{2}(n+1)+\alpha} \Gamma(\frac{1}{2}(n+1)+\alpha)}{\alpha! \Gamma(\frac{1}{2}p+\alpha)} \bigg/ \frac{(t^2/n)^{\frac{1}{2}p-1} (1+t^2/n)^{-\frac{1}{2}(n+1)} \Gamma(\frac{1}{2}(n+1))}{\Gamma(\frac{1}{2}p)} =$$

$$= \frac{\Gamma(\frac{1}{2}p)}{\Gamma(\frac{1}{2}(n+1))} e^{-\frac{1}{2}t^2} \sum_{\alpha=0}^{\infty} \frac{(t^2/2)^{\alpha} \Gamma(\frac{1}{2}(n+1)+\alpha)}{\alpha! \Gamma(\frac{1}{2}p+\alpha)} \left(\frac{t^2/n}{1+t^2/n} \right)^{\alpha}.$$

La parte derecha de la desigualdad es una función estrictamente creciente de $\frac{t^2/n}{1+t^2/n}$ y portanto de t^2 . Así pues la desigualdad se reduce a una de la forma $t^2 > k$ para k elegido adecuadamente. Como no depende de la alternativa Z^2 , el test es uniformemente más potente.

DEFINICIÓN.- Una función crítica $\psi(\bar{x}, A)$ es una función con valores entre 0 y 1 (ambos incluidos) tal que $E(\psi(\bar{x}, A)) = \varepsilon$, al nivel de significación, cuando $\mu = 0$.

Un test aleatorizado existe en seleccionar la hipótesis con probabilidad $\psi(\bar{x}, B)$, cuando $\bar{x} = \bar{x}$, $A = B$. Un test no aleatorizado se define cuando $\psi(\bar{x}, A)$ toma únicamente los valores 0 y 1. Utilizando de manera apropiada el lema de Neyman-Pearson para funciones críticas obtenemos el siguiente resultado.

LEMMA.- Sobre la base de las observaciones x_1, \dots, x_n de $N(\mu, \Sigma)$, de todos los test aleatorizados basados en \bar{x} y A que son invariantes te con respecto a transformaciones $\bar{x}^* = C\bar{x}$ y $A^* = CA' (C \text{ no singular})$, el test T^2 es uniformemente más potente.

TEOREMA.- Sobre la base de N observaciones de $N(\mu, \Sigma)$, de todos los test de $\mu = 0$ que son invariantes con respecto a transformaciones $x_i^* = Cx_i$ (C no singular), el test T^2 es un test uniformemente más potente; y de más, el test T^2 es al menos tan potente como cualquier otro test invariante.

Demostración.- Sea $\psi(x_1, \dots, x_n)$ una función crítica de un test invariante. Entonces

$$E(\psi(x_1, \dots, x_n)) = E_{\bar{x}, A} \{ E(\psi(x_1, \dots, x_n) | \bar{x}, A) \}$$

Puesto que \bar{x}, A son estadísticas suficientes para μ, Σ , $E(\psi(x_1, \dots, x_n) | \bar{x}, A)$ depende sólo de \bar{x}, A . Es invariante, tiene la misma potencia que $\psi(x_1, \dots, x_n)$. Resulta así que cada test en la clase mayor puede ser reemplazado por un test en la clase menor (los que dependen sólo de \bar{x}, A) que tiene la misma potencia. Aplicando el lema anterior se completa la demostración.

TEOREMA.- Dadas las observaciones x_1, \dots, x_n de una $N(\mu, \Sigma)$, de todos los test de $\mu = 0$ basados en \bar{x} y A con potencia dependiendo sólo de $Z^2 = N\mu'Z^{-1}\mu$, el test T^2 es uniformemente más potente.

Demostración.- Puede verse en Anderson pags. 117, 118.

6. EL PROBLEMA DE BEHRENS-FISHER EN EL CASO MULTIVARIANTE

Vamos a dar ahora una solución análoga a la solución de Scheffé (1943) para el problema de Behrens-Fisher. Sean $\{x_{\alpha}^{(i)}\}$, $\alpha = 1, \dots, N_i$, $i = 1, 2$ muestras procedente de $N(\mu^{(i)}, \Sigma_i)$, $i = 1, 2$. Queramos contrastar la hipótesis $\mu^{(1)} = \mu^{(2)}$. La media $\bar{x}^{(1)}$ de la primera muestra se distribuye normalmente con esperanza

$$E(\bar{x}^{(1)}) = \mu^{(1)}$$

y matriz de varianzas

$$E[(\bar{x}^{(1)} - \mu^{(1)})(\bar{x}^{(1)} - \mu^{(1)})'] = \frac{1}{N_1} \Sigma_1.$$

Análogamente, para la media de la segunda muestra tenemos

$$E(\bar{x}^{(2)}) = \mu^{(2)}$$

y matriz de varianzas

$$E[(\bar{x}^{(2)} - \mu^{(2)})(\bar{x}^{(2)} - \mu^{(2)})'] = \frac{1}{N_2} \Sigma_2.$$

Así $\bar{x}^{(1)} - \bar{x}^{(2)}$ tiene media $\mu^{(1)} - \mu^{(2)}$ y matriz de covarianzas $(1/N_1) \Sigma_1 + (1/N_2) \Sigma_2$. Ahora bien, no podemos utilizar la técnica desarrollada anteriormente puesto que

$$\sum_{\alpha=1}^{N_1} (x_{\alpha}^{(1)} - \bar{x}^{(1)})(x_{\alpha}^{(1)} - \bar{x}^{(1)})' + \sum_{\alpha=1}^{N_2} (x_{\alpha}^{(2)} - \bar{x}^{(2)})(x_{\alpha}^{(2)} - \bar{x}^{(2)})'$$

no tiene la distribución de Wishart con matriz de covarianzas múltiple de $(1/N_1) \Sigma_1 + (1/N_2) \Sigma_2$.

Si $N_1 = N_2 = N$, podemos utilizar el t^2 de una manera obvia. Veamoslo. Sea $y_{\alpha} = x_{\alpha}^{(1)} - x_{\alpha}^{(2)}$ (suponiendo que la numeración de las observaciones de las dos muestras es independiente de las mismas observaciones). Entonces y_{α} se distribuye normalmente con media $\mu^{(1)} - \mu^{(2)}$ y matriz de covarianzas $\Sigma_1 + \Sigma_2$, e independientemente de y_{β} ($\beta \neq \alpha$). Sea $\bar{y} = \frac{1}{N} \sum_{\alpha} y_{\alpha} = \bar{x}^{(1)} - \bar{x}^{(2)}$ y definamos S mediante

$$(N-1)S = \sum_{\alpha=1}^N (y_{\alpha} - \bar{y})(y_{\alpha} - \bar{y})' = \sum_{\alpha=1}^N (x_{\alpha}^{(1)} - x_{\alpha}^{(2)} - \bar{x}^{(1)} + \bar{x}^{(2)})(x_{\alpha}^{(1)} - x_{\alpha}^{(2)} - \bar{x}^{(1)} + \bar{x}^{(2)})'$$

Entonces

$$T^2 = N \bar{y}' S^{-1} \bar{y}$$

es adecuado para contrastar la hipótesis $\mu_y = \mu^{(1)} - \mu^{(2)} = 0$ y T^2 tiene la distribución T^2 con $N-1$ grados de libertad. Hay que señalar en este punto que si hubiéramos sabido que $\Sigma_1 = \Sigma_2$ habríamos utilizado un estadístico T^2 con $2N-2$ grados de libertad; hemos pues perdido $N-1$ grados de libertad en la construcción de un test que es independiente de las dos matrices de covarianzas.

Volvamos nuevamente al caso en que $N_1 \neq N_2$. Por conveniencia supongamos $N_1 < N_2$. Definimos entonces

$$y_{\alpha} = x_{\alpha}^{(1)} - \sqrt{\frac{N_1}{N_2}} x_{\alpha}^{(2)} + \frac{1}{\sqrt{N_1 N_2}} \sum_{\beta=1}^{N_1} x_{\beta}^{(2)} - \frac{1}{N_2} \sum_{\gamma=1}^{N_2} x_{\gamma}^{(2)} \quad \alpha = 1, \dots, N_1$$

El valor esperado de y_{α} será

$$E(y_{\alpha}) = \mu^{(1)} - \sqrt{\frac{N_1}{N_2}} \mu^{(2)} + \frac{N_1}{\sqrt{N_1 N_2}} \mu^{(2)} - \frac{N_2}{N_2} \mu^{(2)} = \mu^{(1)} - \mu^{(2)}$$

La matriz de covarianzas de $y_{\alpha} \in y_{\beta}$ es

$$\begin{aligned} E(y_{\alpha} - E(y_{\alpha}))(y_{\beta} - E(y_{\beta}))' &= E\left[(x_{\alpha}^{(1)} - \mu^{(1)}) - \sqrt{\frac{N_1}{N_2}}(x_{\alpha}^{(2)} - \mu^{(2)}) + \frac{1}{\sqrt{N_1 N_2}} \sum_{\gamma=1}^{N_1} (x_{\gamma}^{(2)} - \mu^{(2)}) - \frac{1}{N_2} \sum_{\gamma=1}^{N_2} (x_{\gamma}^{(2)} - \mu^{(2)})\right] \\ &\quad \left[(x_{\beta}^{(1)} - \mu^{(1)}) - \sqrt{\frac{N_1}{N_2}}(x_{\beta}^{(2)} - \mu^{(2)}) + \frac{1}{\sqrt{N_1 N_2}} \sum_{\gamma=1}^{N_1} (x_{\gamma}^{(2)} - \mu^{(2)}) - \frac{1}{N_2} \sum_{\gamma=1}^{N_2} (x_{\gamma}^{(2)} - \mu^{(2)})\right]' \\ &= \delta_{\alpha\beta} \Sigma_1 + \frac{N_1}{N_2} \delta_{\alpha\beta} \Sigma_2 + \Sigma_2 \left(-2 \frac{1}{N_2} + \frac{2}{N_2} \sqrt{\frac{N_1}{N_2}} + \frac{N_1}{N_1 N_2} - 2 \frac{N_1}{\sqrt{N_1 N_2} N_2} + \frac{N_2}{N_2^2}\right) \\ &= \delta_{\alpha\beta} \left(\Sigma_1 + \frac{N_1}{N_2} \Sigma_2\right) \end{aligned}$$

Así un estadístico apropiado para contrastar la hipótesis $\mu^{(1)} - \mu^{(2)} = 0$, que tiene una distribución T^2 con $N_1 - 1$ grados de libertad, es

$$T^2 = N_1 \bar{y}' S^{-1} \bar{y}$$

donde

$$\bar{y} = \frac{1}{N_1} \sum_{\alpha=1}^{N_1} y_{\alpha} = \bar{x}^{(1)} - \bar{x}^{(2)}$$

y

$$(N_1 - 1)S = \sum_{\alpha=1}^{N_1} (y_{\alpha} - \bar{y})(y_{\alpha} - \bar{y})' = \sum_{\alpha=1}^{N_1} \left(x_{\alpha}^{(1)} - \bar{x}^{(1)} - \sqrt{\frac{N_1}{N_2}} \left(x_{\alpha}^{(2)} - \frac{1}{N_1} \sum_{\beta=1}^{N_1} x_{\beta}^{(2)}\right)\right) \left(x_{\alpha}^{(1)} - \bar{x}^{(1)} - \sqrt{\frac{N_1}{N_2}} \left(x_{\alpha}^{(2)} - \frac{1}{N_1} \sum_{\beta=1}^{N_1} x_{\beta}^{(2)}\right)\right)'$$

lo que en términos de $u_{\alpha} = x_{\alpha}^{(1)} - \sqrt{N_1/N_2} x_{\alpha}^{(2)}$, $\alpha = 1, \dots, N_1$, puede escribirse como

$$(N_1 - 1) S = \sum_{\alpha=1}^{N_1} (u_\alpha - \bar{u})(u_\alpha - \bar{u})'$$

$$\text{donde } \bar{u} = \frac{1}{N_1} \sum_{\alpha=1}^{N_1} u_\alpha.$$

Este procedimiento ha sido sugerido por Scheffé (1943), que volvió de esta manera el caso univariante. Demostró Scheffé que en dicho caso, esta técnica da los intervalos de confianza más pequeños utilizando la t de Student. La ventaja del método estriba en el hecho de utilizar $\bar{x}^{(1)} - \bar{x}^{(2)}$, que es el estadístico más importante para contrastar hipótesis acerca de $\mu^{(1)} - \mu^{(2)}$. El sacrificio de observaciones a la hora de estimar la matriz de covarianzas no es demasiado importante. La extensión del método de Scheffé al caso multivariante se debe a Bennett (1951).

El método es generalizable a situaciones con más de dos poblaciones. Sean $\{x_\alpha^{(i)}\}$, $\alpha=1, \dots, N_i$, $i=1, \dots, q$ muestras provenientes de $N(\mu^{(i)}, \Sigma_i)$ $i=1, \dots, q$, respectivamente. Consideremos la siguiente hipótesis

$$\sum_{i=1}^q \beta_i \mu^{(i)} = \mu.$$

donde los β_i son scalares dados, al igual que el vector μ . Si las N_i son distintas, sea N_1 la menor de ellas.

Definamos

$$y_\alpha = \beta_1 x_\alpha^{(1)} + \sum_{i=2}^q \beta_i \sqrt{\frac{N_1}{N_i}} \left(x_\alpha^{(i)} - \frac{1}{N_i} \sum_{\beta=1}^{N_i} x_\beta^{(i)} + \frac{1}{\sqrt{N_1 N_i}} \sum_{\beta=1}^{N_i} x_\beta^{(i)} \right).$$

Entonces

$$E(y_\alpha) = \beta_1 \mu^{(1)} + \sum_{i=2}^q \beta_i \sqrt{\frac{N_1}{N_i}} \left(\mu^{(i)} - \frac{1}{N_i} \sum_{\beta=1}^{N_i} \mu^{(i)} + \frac{N_i}{\sqrt{N_1 N_i}} \mu^{(i)} \right) = \sum_{i=1}^q \beta_i \mu^{(i)}.$$

y

$$E(y_\alpha - E(y_\alpha))(y_\beta - E(y_\beta))' = \delta_{\alpha\beta} \left(\sum_{i=1}^q \frac{\beta_i^2 N_1}{N_i} \Sigma_i \right).$$

Sean \bar{y} y S definidas mediante

$$\bar{y} = \frac{1}{N_1} \sum_{\alpha=1}^{N_1} y_\alpha = \sum_{i=1}^q \beta_i \bar{x}^{(i)}, \quad \text{con } \bar{x}^{(i)} = \frac{1}{N_i} \sum_{\beta=1}^{N_i} x_\beta^{(i)},$$

$$(N_1 - 1) S = \sum_{\alpha=1}^{N_1} (y_\alpha - \bar{y})(y_\alpha - \bar{y})'.$$

Entonces

$$T^2 = N_1 (\bar{y} - \mu)' S^{-1} (\bar{y} - \mu)$$

es un estadístico apropiado para contrastar la hipótesis propuesta, cuando la hipótesis es cierta este estadístico tiene una distribución T^2 con $N_1 - 1$ grados de libertad. mediante un cambio del tipo

$$u_\alpha = \sum_{i=1}^q \beta_i \sqrt{\frac{N_1}{N_i}} x_\alpha^{(i)} \quad \alpha=1, \dots, N_1$$

S puede expresarse de la forma

$$(N_1 - 1) S = \sum_{\alpha=1}^{N_1} (u_\alpha - \bar{u})(u_\alpha - \bar{u})'.$$

Para terminar un último problema que puede resolverse mediante este procedimiento es el de contrastar la hipótesis de que dos subvectores de un determinado vector normal multivariante tienen igual media. Sea

$$x = \begin{pmatrix} x^{(1)} \\ x^{(2)} \end{pmatrix}$$

distribuido normalmente con media

$$\mu = \begin{pmatrix} \mu^{(1)} \\ \mu^{(2)} \end{pmatrix}$$

y matriz de varianzas

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}.$$

Supongamos que $x^{(1)}$ y $x^{(2)}$ tienen cada uno p componentes. Entonces $x^{(1)} - x^{(2)}$ se distribuye normalmente con media $\mu^{(1)} - \mu^{(2)}$ y matriz de varianzas

$$E \left((x^{(1)} - \mu^{(1)}) - (x^{(2)} - \mu^{(2)}) \right) \left((x^{(1)} - \mu^{(1)}) - (x^{(2)} - \mu^{(2)}) \right)' = \Sigma_{11} - \Sigma_{21} - \Sigma_{12} + \Sigma_{22}.$$

Para contrastar la hipótesis $\mu^{(1)} = \mu^{(2)}$ utilizaremos el estadístico T^2 definido mediante:

$$T^2 = N (\bar{x}^{(1)} - \bar{x}^{(2)})' (S_{11} - S_{21} - S_{12} + S_{22})^{-1} (\bar{x}^{(1)} - \bar{x}^{(2)})$$

donde la media muestral y la matriz de varianzas muestral han sido fraccionadas de forma similar a como lo han sido μ y Σ .

EJEMPLOS

- ① En una investigación de las respuestas al test de Wechsler para la inteligencia de adultos aplicado a hombres y mujeres de edad avanzada se obtuvieron las siguientes medias para las respuestas verbal y de computarizado

$$\begin{bmatrix} \bar{x}_v \\ \bar{x}_c \end{bmatrix} = \begin{bmatrix} 55.24 \\ 34.97 \end{bmatrix}$$

siendo la muestra de un tamaño $N=101$ y los sujetos de edad comprendidos entre los 60 y 64 años. La matriz de varianzas muestral de las respuestas fue

$$S = \begin{bmatrix} 210.54 & 126.99 \\ 126.99 & 119.68 \end{bmatrix}$$

Desamos contrastar, a un nivel $\alpha = 0.01$ la hipótesis nula de que las observaciones provienen de una población con vector media

$$\mu_0 = \begin{bmatrix} 60 \\ 50 \end{bmatrix}.$$

El test está basado en la aplicación del estadístico T^2 , cuyo valor vendría dado por

$$T^2 = N (\bar{x} - \mu_0)' S^{-1} (\bar{x} - \mu_0) = 101 \cdot [55.24 - 60, 34.97 - 50] \begin{bmatrix} 0.01319 & -0.01400 \\ -0.01400 & 0.02321 \end{bmatrix} \begin{bmatrix} 55.24 - 60 \\ 34.97 - 50 \end{bmatrix} = 357.43$$

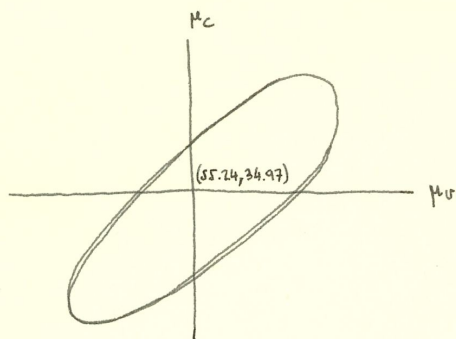
el valor $[T^2/(N-1)] [(N-p)/p]$ es una F de Snedecor con p y $N-p$ grados de libertad, por tanto

$$F = \frac{357.43}{100} \cdot \frac{99}{2} = 176.93$$

como $F_{0.01, 2, 99} \approx 4.98$, y $F \gg F_{0.01, 2, 99}$ rechazamos la hipótesis formulada. A la vista de los resultados siendo cabe pensar que la población origen de estos datos debe poseer un vector media menor que el postulado en H_0 .

Podemos también hacer uso de T^2 para obtener una región al 99% de confianza para el vector media de la población. Necesitaremos el valor de $F_{0.01, 2, 99}$ que podemos obtener por interpolación entre los de $F_{0.01, 2, 60}$ y $F_{0.01, 2, 120}$, lo que da para $F_{0.01, 2, 99}$ un valor 4.83.

la elipse de confianza vendrá dada por la ecuación



$$1.33 (\mu_v - 55.24)^2 - 2.83 (\mu_v - 55.24)(\mu_c - 34.97) + 2.34 (\mu_c - 34.97)^2 \leq 9.76$$

- ② Cuarenta ~~niños~~ ancianos que participaron en un estudio interdisciplinario del envejecimiento humano fueron clasificados en dos categorías de diagnóstico, "factor senil presente" y "factor senil ausente", sobre la base de un intenso examen psiquiátrico. La escala de inteligencia de adultos de Wechsler fue administrada a todos los sujetos por un investigador independiente, y ciertos subtests mostraron grandes diferencias entre los dos grupos. A la vista de estos resultados queremos intentar la hipótesis de que los grupos provienen de una misma población.

Los datos muestrales se recogen en la siguiente tabla (datos medios)

Subtest	Grupo	
	Factor senil ausente $N_1 = 37$	Factor senil $N_2 = 12$
Información	12.57	8.75
Similitudes	9.57	5.33
Aritmética	11.49	8.50
Figuras complejas	7.97	4.75

La matriz de covarianzas empírica viene dada por

$$S = \begin{bmatrix} 11.2583 & 9.4042 & 7.1489 & 3.3830 \\ & 13.5318 & 7.3830 & 2.5532 \\ & & 11.5744 & 2.6170 \\ & & & 5.8085 \end{bmatrix}$$

El valor de T^2 para dos muestras, en este caso vale

$$T^2 = \frac{N_1 N_2}{N_1 + N_2} (\bar{X}_1 - \bar{X}_2)' S^{-1} (\bar{X}_1 - \bar{X}_2) = 22.05$$

y el correspondiente valor de F , es

$$F = \frac{N_1 + N_2 - p - 1}{(N_1 + N_2 - 2)p} T^2 = \frac{44}{47.4} 22.05 = 5.16$$

Bajo la hipótesis de igualdad de medias la probabilidad de exceder semejante valor de F sería menor que 0.005, por tanto rechazamos la hipótesis formulada para niveles de significación convencionales del 5% o del 1%.

- ③ Se estudian los cambios de concentración en el plasma de ácidos grasos libres (FFA) y glucosa sanguínea (G) en 12 pacientes esquizofrénicos y 13 voluntarios normales, después de inyecciones intramusculares de insulina. Los cambios resumen en la siguiente tabla

	Cambio medio	
	Esquizofrénicos	Normales
G, mg, %	-25.6	-31.1
FFA, mcg/litro	-0.06	-0.15
Matrices de suma de cuadrados y productos cruzados	$\begin{bmatrix} 3455 & 9.9492 \\ 9.9492 & 0.1105 \end{bmatrix}$	$\begin{bmatrix} 3509 & -3.2408 \\ -3.2408 & 0.0865 \end{bmatrix}$
Matriz de covarianzas "pooled"	$S = \begin{bmatrix} 302 & 0.292 \\ 0.292 & 0.00856 \end{bmatrix}$	

Bajo el supuesto de que los valores de FFA, G tienen una distribución normal bivariante en una matriz de variables que no se afecta por diagnósticos psiquiátricos, debemos contrastar la hipótesis de que los vectores medios en los sujetos normales y en los esquizofrénicos son equivalentes. Observamos que la hipótesis de la igualdad de la matriz de covarianzas aparece válida respecto de los elementos de la diagonal principal de la matriz de varianzas muestral, aunque no tanto para los términos fuera de ésta, es decir la covarianza entre ambas variables.

El valor de T^2 para esta hipótesis

$$H_0: \mu_1 - \mu_2 = 0$$

viene dado por $T^2 = 6.03$ y el valor de F asociado es $F = 2.88$. Observamos que este valor está comprendido entre los de $F_{0.1, 2, 22} = 2.56$ y $F_{0.05, 2, 22} = 3.44$. Ello significa que a un nivel $\alpha = 0.05$ no rechazaremos la hipótesis nula, mientras que la aceptaremos para un nivel superior $\alpha = 0.1$.

- ④ Los niveles de FFA en sangre se midieron en 15 sujetos normales que se sometieron voluntariamente a series de hipnosis. Durante estas series se pidió que experimentaran miedo, angustia y alegría y se le midió la cantidad (o mejor el cambio en la concentración) de FFA después de cada una de estas situaciones ficticias. Bajo el supuesto de que cada una de estas situaciones produce el mismo grado de cambio, los investigadores desearon comprobar que no había cambios apreciables entre las diferentes series de stress a las que se sometió a los individuos. Los cambios medios de FFA fueron

$$\bar{x}_1 = 2.699 \quad \bar{x}_2 = 2.178 \quad \bar{x}_3 = 2.558$$

El problema puede resolverse como un caso particular de los problemas de simetría antes presentados. Podemos efectuar el siguiente cambio

$$y_{i1} = x_{i1} - x_{i2} \quad y_{i2} = x_{i1} - x_{i3} \quad , \quad i = 1, \dots, 15$$

es decir

$$Y = \begin{bmatrix} 1 & -1 & 0 \\ 1 & 0 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

La hipótesis $H_0: \mu_1 = \mu_2 = \mu_3$ es ahora equivalente a $H'_0: \mu_Y = 0$, donde μ_Y es la media de Y y μ_X es la media de X . La matriz de varianzas muestral para Y es

$$S = \begin{bmatrix} 1.7343 & 1.1666 \\ 1.1666 & 2.7733 \end{bmatrix} \quad , \quad S^{-1} = \begin{bmatrix} 0.8041 & -0.3382 \\ -0.3382 & 0.5029 \end{bmatrix}$$

El valor de $T^2 = N \bar{y}' S^{-1} \bar{y}$ viene dado por $T^2 = 2.64$ y el correspondiente valor de $F = 1.24$ y como $F_{0.5, 2, 13} = 3.80$, aceptaremos la hipótesis formulada, es decir, que los cambios de FFA en las diferentes situaciones de stress son iguales.

Capítulo 5

Clasificación de observaciones

1. EL PROBLEMA DE LA CLASIFICACION

El problema de la clasificación surge cuando un investigador realiza varias medidas sobre un individuo y desea clasificarlo en una de varias categorías sobre la base de estas medidas. El investigador no puede identificar al individuo directamente con una categoría y debe por tanto reunir a las medidas en un todo. En muchos casos puede suponerse que hay un número finito de categorías o poblaciones a las que el individuo puede pertenecer y cada una de las poblaciones se caracteriza por una distribución de probabilidad de las medidas. Así, un individuo es considerado como una observación aleatoria de una población. La cuestión es: Dado un individuo con ciertas medidas, de qué población proviene?

El problema de la clasificación puede ser considerado como un problema de funciones estadísticas de decisión. Tenemos un cierto número de hipótesis: cada una de las hipótesis es que la distribución de la observación es una dada. Debemos aceptar una de estas hipótesis y rechazar las otras. Si solo existen dos poblaciones, nos enfrentamos ante un problema elemental de contrastar una hipótesis de una distribución específica frente a otra.

En algunos casos, las categorías están especificadas de antemano en el sentido de que la distribución de probabilidad de las medidas es completamente conocida. En otros casos, la forma de cada una de las distribuciones puede ser conocida, pero los parámetros de la distribución deben estimarse a partir de una muestra de la población.

Consideremos un ejemplo de clasificación. Los alumnos que desean ingresar en la universidad son sometidos a una batería de tests; el vector de los resultados es un conjunto de medidas x . El estudiante puede ser miembro de una población que consiste en aquellos estudiantes que están capacitados para sacar adelante con éxito los estudios, potencialmente clasificados como 'sí', o bien pertenecer a la población de los que no completa con éxito dichos estudios. El problema es clasificar un estudiante sobre la base de sus resultados en el examen de admisión.

En este capítulo desarrollaremos la teoría de clasificación en términos generales para aplicarla posteriormente a aquellos casos en que se involucran poblaciones con distribución normal.

2. NORMAS PARA UNA BUENA CLASIFICACION

Consideraciones previas. - En la intención de un proceso de clasificación se desea minimizar la probabilidad de mala clasificación; o más específicamente, se desea minimizar, por término medio, los malos efectos de clasificaciones erróneas. Leamos de preürar más este concepto. Por conveniencia consideraremos el caso de dos categorías. Luego trataremos el caso general.

Supongamos que un individuo es una observación de una de estas dos poblaciones, π_1 o π_2 . La clasificación de una observación depende del vector de medidas $x' = (x_1, \dots, x_p)$ sobre este individuo. Construiremos una regla de manera que si el individuo está caracterizado mediante ciertos conjuntos de valores de x_1, \dots, x_p sea clasificado como perteneciente a π_1 ; si tiene otros valores sea clasificado como perteneciente a π_2 .

Podemos considerar una observación como un punto en un espacio p -dimensional. Dividimos al espacio en dos regiones. Si la observación está en R_1 , clasificamos al individuo como de π_1 ; si la observación cae en R_2 lo clasificamos como de π_2 .

Siguiendo un procedimiento de este tipo, el investigador puede cometer dos tipos de errores. Clasificar al individuo como de π_2 cuando realmente procede de π_1 y viceversa. Acostumamos conocer la indeseabilidad relativa de estos dos clases de mala clasificación. Denotaremos por $C(2/1) > 0$ el coste del primero de los errores y por $C(1/2) > 0$ el otro. Estos costes pueden venir medidos en cualquier tipo de unidades. Como más tarde veremos, lo verdaderamente importante es la relación entre los dos costes.

En la tabla que presentamos a continuación aparecen descritas las anteriores situaciones. Obviamente, un buen procedimiento

		Decisiones del estadístico	
		π_1	π_2
Población de origen	π_1	0	$C(2/1)$
	π_2	$C(1/2)$	0

es aquel que minimiza, en algún sentido, el coste de las clasificaciones erróneas.

Caso de dos poblaciones

Veamos a continuación dos formas de definir el criterio mínimo. En un caso suponemos que tenemos probabilidades a priori de ambas poblaciones. Sean q_1, q_2 las probabilidades de que una observación provenga de las poblaciones π_1 y π_2 respectivamente. Las propiedades probabilísticas de π_1 están especificadas mediante una función de densidad. Por comodidad consideremos el caso en que la distribución de probabilidad es absolutamente continua y existe la función de densidad, aunque el caso discreto requiere intuitivamente muy similares. Sean $p_i(x)$, $i=1,2$ las densidades de probabilidad de las poblaciones π_i , $i=1,2$. Si tenemos una región R_1 de clasificación en la que asignamos la observación a π_1 , la probabilidad de clasificar correctamente una observación que realmente proviene de π_1 es

$$P(1/1, R) = \int_{R_1} p_1(x) dx,$$

y la probabilidad de mala clasificación de una observación que procede de π_2 es

$$P(2/1, R) = \int_{R_1} p_2(x) dx.$$

Análogamente tendremos

$$P(2/2, R) = \int_{R_2} p_2(x) dx$$

y

$$P(1/2, R) = \int_{R_1} p_2(x) dx.$$

Puesto que la probabilidad de extraer una observación de π_1 es q_1 , la probabilidad de extraer una observación de π_1 y clasificarla correctamente es $q_1 \cdot P(1/1, R)$. Análogamente obtendríamos las probabilidades de los tres restantes resultados posibles.

¿Cuál será la pérdida esperada o el coste medio de mala clasificación? Es la suma de los productos de los costes de cada una de las clasificaciones erróneas por la probabilidad de que ocurra; a saber

$$C(2/1) \cdot P(2/1, R) \cdot q_1 + C(1/2) \cdot P(1/2, R) \cdot q_2$$

Es este coste medio el que se desea minimizar. El derivar, se quiere dividir el espacio en dos regiones R_1 y R_2 de manera que la pérdida esperada sea lo menos posible. Un procedimiento que minimiza la anterior expresión para unas q_1, q_2 dadas es conocido como un procedimiento Bayes.

El otro caso a considerar es aquel en el que no se conocen probabilidades a priori. En este caso la pérdida esperada si la observación proviene de π_2 es

$$C(2/1) \cdot P(2/1, R) = r(1, R);$$

y para observaciones procedentes de π_1 ,

$$C(1/2) \cdot P(1/2, R) = r(2, R).$$

No sabemos de donde proviene nuestra observación y no conocemos las probabilidades de que pertenezca a una u otra población.

Un procedimiento R es al menos tan bueno como un procedimiento R^* si $r(1, R) \leq r(1, R^*)$ y $r(2, R) \leq r(2, R^*)$. R es mejor que R^* si al menos una de estas desigualdades es estricta. Normalmente no hay procedimientos que sean mejor que cualquier otro o al menos tan buenos como cualquier otro. Un procedimiento R recibe que es admisiblemente si no existe ningún otro que sea mejor; estamos interesados en la clase de los procedimientos admisibles. Demostremos que bajo ciertas condiciones esta clase es la misma que la clase de los procedimientos Bayes. Una clase de procedimientos es completa si para cualquier procedimiento exterior existe uno en la clase que es mejor; una clase es esencialmente completa si para cualquier procedimiento exterior existe uno en la clase que al menos tan bueno como aquel. Una clase mínima completa (si existe) es una clase completa tal que ningún subconjunto propio es una clase completa; una definición similar puede darse para las clases mínimas esencialmente completas. Bajo ciertas condiciones demostraremos que la clase

admisibles y seriamente completa. Para multiplicar la discusión consideramos que dos procedimientos son iguales si difieren solo en conjuntos de probabilidad nula. De hecho, a través de cuanto requiriríamos afirmaciones que no vistas "excepción hecha de conjuntos de probabilidad nula" aunque no lo hagamos constar explícitamente.

El principio que usualmente conduce a un procedimiento único es el principio del minimax. Un procedimiento es minimax si la máxima pérdida esperada, $r(i, R)$ es un mínimo. Desde un punto de vista conservador, este puede ser considerado como un procedimiento óptimo. Un estudio detallado de estos conceptos puede encontrarse en los textos de Wald (1950), Blackwell y Goshikie (1954).

3. PROCEDIMIENTOS DE CLASIFICACIÓN EN UNA DE DOS POBLACIONES CUANDO SE CONocen LAS DISTRIBUCIONES DE PROBABILIDAD

Probabilidades a priori conocidas

Puesto que conocemos las probabilidades a priori podemos definir la distribución conjunta de las poblaciones y del conjunto de variables observadas. La probabilidad de que una observación provenga de Π_1 y de que cada una de sus variables sea menor que la correspondiente componente de y viene dada por

$$\int_{-\infty}^{y_1} \dots \int_{-\infty}^{y_k} q_1 p_1(x) dx_1 \dots dx_k$$

Podemos también definir la probabilidad condicional de que una observación provenga de una cierta población dados los valores de las variables observadas. Por ejemplo, la probabilidad condicional de venir de la población Π_1 , dada una observación x ,

$$\frac{q_1 p_1(x)}{q_1 p_1(x) + q_2 p_2(x)}$$

Supongamos por un momento que $C(1/2) = C(2/1) = 1$. Entonces la pérdida esperada es

$$q_1 \int_{R_2} p_1(x) dx + q_2 \int_{R_1} p_2(x) dx.$$

Esto también es la probabilidad de una mala clasificación; por tanto, deseamos minimizar la probabilidad de una clasificación errónea.

Para una observación dada x minimizaremos la probabilidad de una mala clasificación asignando la población que tiene una probabilidad condicional mayor. Si

$$\frac{q_1 p_1(x)}{q_1 p_1(x) + q_2 p_2(x)} \geq \frac{q_2 p_2(x)}{q_1 p_1(x) + q_2 p_2(x)}$$

elegiremos la población Π_1 . En cualquier otro caso elegiremos Π_2 . Puesto que minimizaremos la probabilidad de mala clasificación en cada punto, la minimizaremos sobre el espacio entero. Así, la regla es

$$\begin{aligned} R_1: & q_1 p_1(x) \geq q_2 p_2(x) \\ R_2: & q_1 p_1(x) < q_2 p_2(x) \end{aligned}$$

En caso de igualdad la observación puede ser asignada indistintamente a una u otra población. Si $q_1 p_1(x) + q_2 p_2(x) = 0$ para un x dado, también en este caso el punto puede asignarse a cualquiera de las regiones.

Proponemos ahora formalmente que esta regla es el mejor procedimiento. Para cualquier procedimiento $R^* = (R_1^*, R_2^*)$ la probabilidad de mala clasificación es

$$q_1 \int_{R_2^*} p_1(x) dx + q_2 \int_{R_1^*} p_2(x) dx = \int_{R_2^*} [q_1 p_1(x) - q_2 p_2(x)] dx + q_2 \int p_2(x) dx.$$

El segundo sumando del segundo miembro es una constante y en cuanto al primer sumando se minimiza si R_2^* incluye los puntos x tales que $q_1 p_1(x) - q_2 p_2(x) < 0$ y excluye los puntos para los que $q_1 p_1(x) - q_2 p_2(x) > 0$. Si suponemos que

$$P \left\{ \frac{p_1(x)}{p_2(x)} = \frac{q_2}{q_1} \mid \Pi_i \right\} = 0 \quad i=1,2$$

entonces el procedimiento Bayes es único excepto para conjuntos de probabilidad nula.

Hemos con esto resuelto el problema cuando los q_i son iguales e iguales a uno. Supongamos ahora que no ocurre así. Se trata en este caso de minimizar

$$G(2/1) \cdot q_1 \int_{R_2} p_1(x) dx + G(1/2) \cdot q_2 \int_{R_1} p_2(x) dx$$

elijamos entonces R_1 y R_2 de manera que

$$\begin{aligned} R_1 &: [G(2/1) \cdot q_1] p_1(x) \geq [G(1/2) \cdot q_2] p_2(x) \\ R_2 &: [G(2/1) \cdot q_1] p_1(x) < [G(1/2) \cdot q_2] p_2(x) \end{aligned}$$

puesto que $G(2/1)$ y $G(1/2)$ son constantes no negativas. Otra forma de definir R_1 y R_2 sería

$$\begin{aligned} R_1 &: \frac{p_1(x)}{p_2(x)} \geq \frac{G(1/2) \cdot q_2}{G(2/1) \cdot q_1} \\ R_2 &: \frac{p_1(x)}{p_2(x)} < \frac{G(1/2) \cdot q_2}{G(2/1) \cdot q_1} \end{aligned}$$

TEOREMA.- Si q_1 y q_2 son las probabilidades a priori de extraer una observación de la población Π_1 con probabilidad (densidad) $p_1(x)$ y de la población Π_2 con densidad de probabilidad $p_2(x)$, respectivamente, y si δ es una regla de mala clasificación de una observación de Π_1 como procedente de Π_2 es $G(2/1)$ y de una observación de Π_2 como procedente de Π_1 es $G(1/2)$, entonces las reglas de clasificación R_1 y R_2 , definidas mediante las anteriores expresiones, minimizan el coste esperado. Si además

$$P \left\{ \frac{p_1(x)}{p_2(x)} = \frac{q_2 G(1/2)}{q_1 G(2/1)} \mid \Pi_i \right\} = 0, \quad i=1,2$$

entonces el procedimiento es único casi por todas partes.

Cuando no se conocen probabilidades a priori

En muchos ejemplos de clasificación el estadístico no puede asignar probabilidades a priori a las dos poblaciones. En este caso trataremos de encontrar la clase de los procedimientos admisibles, es decir, el conjunto de procedimientos que no puede ser mejorado.

Probaremos primero que un procedimiento Bayes es admisible. Sea $R = (R_1, R_2)$ un procedimiento Bayes para q_1 y q_2 dados; existirá un procedimiento $R^* = (R_1^*, R_2^*)$ tal que $P(1/2, R^*) \leq P(1/2, R)$ y $P(2/1, R^*) \leq P(2/1, R)$ con al menos una desigualdad estricta. Puesto que R es Bayes

$$q_1 \cdot P(2/1, R) + q_2 \cdot P(1/2, R) \leq q_1 \cdot P(2/1, R^*) + q_2 \cdot P(1/2, R^*).$$

Esta desigualdad puede escribirse

$$q_1 [P(2/1, R) - P(2/1, R^*)] \leq q_2 [P(1/2, R^*) - P(1/2, R)].$$

Supongamos que $q_2 > 0$. Entonces si $P(1/2, R^*) \leq P(1/2, R)$, la parte derecha de la desigualdad será menor o igual a 0, por tanto $P(2/1, R) \leq P(2/1, R^*)$. Si $q_2 = 0$, razonando de forma análoga llegamos a que $P(1/2, R) \leq P(1/2, R^*)$. Así R^* no es mejor que R y R es admisible. Si $q_1 = 0$, entonces de la desigualdad obtenemos $0 \leq q_2 [P(1/2, R^*) - P(1/2, R)]$. Para un procedimiento Bayes esto impone que R_1 contiene aquellos x tales que $p_2(x) = 0$. Por tanto $P(1/2, R) = 0$ y si R^* ha de ser mejor que R , deberá ocurrir que $P(1/2, R^*) = 0$. En cuanto a $P(2/1, R)$ si $P(p_2(x) = 0 / \Pi_1) = 0$, entonces $P(2/1, R) = P(p_2(x) > 0 / \Pi_1) = 1$. Ahora bien si $P(1/2, R^*) = 0$, entonces R_1^* contiene sólo puntos para los que $p_2(x) = 0$, pero $P(2/1, R^*) = P(R_1^* / \Pi_1) = P(p_2(x) > 0 / \Pi_1) = 1$, y por tanto R^* no es mejor que R .

TEOREMA.- Si $P(p_2(x) = 0 / \Pi_1) = 0$ y $P(p_1(x) = 0 / \Pi_2) = 0$ entonces todo procedimiento Bayes es admisible.

Problema ahora el inverso, es decir, que todo procedimiento admisible es Bayes. Supondremos

5(3)

$$P \left\{ \frac{p_1(x)}{p_2(x)} = k/\pi_i \right\} = 0, \quad i=1,2; \quad 0 \leq k \leq \infty \quad \left[\frac{p_1(x)}{p_2(x)} = \infty \text{ significa que } p_2(x) = 0 \right].$$

Entonces para cualquier q_2 el procedimiento Bayes es único. Además la función de distribución de $p_1(x)/p_2(x)$ para π_1, π_2 es continua.

Sea R un procedimiento admisible. Entonces existe un k tal que

$$P(2/1, R) = P \left\{ \frac{p_1(x)}{p_2(x)} \leq k/\pi_1 \right\} = P(2/1, R^*)$$

donde R^* es el procedimiento Bayes independiente a $q_2/q_1 = k$ [es decir, $q_2 = 1/(1+k)$]. Puesto que R es admisible, tendremos $P(1/2, R) \leq P(1/2, R^*)$. Sin embargo el teorema anterior afirma que R^* es admisible, por tanto $P(1/2, R) \geq P(1/2, R^*)$ y de aquí la igualdad de ambos. En definitiva R es un procedimiento Bayes y por la unicidad de este procedimiento debe de ser precisamente R^* .

TEOREMA.- Si $P \left\{ p_1(x)/p_2(x) = k/\pi_i \right\} = 0, \quad i=1,2; \quad 0 \leq k \leq \infty$, entonces cualquier procedimiento admisible es un procedimiento Bayes.

De acuerdo con lo visto en la demostración de este último teorema la clase de los procedimientos Bayes es completa. En realidad, y desde las consideraciones relativas de una clase minimal completa justo que incluye en la clase de los procedimientos admisibles.

Consideremos finalmente el procedimiento minimax. Sea $P(i/j, q_2) = P(i/j, R)$, donde R es el procedimiento Bayes correspondiente a q_2 . $P(i/j, q_2)$ es una función continua de q_2 . Así, $P(2/1, q_2)$ varía de 1 a 0 cuando q_2 lo hace de 0 a 1, mientras que $P(1/2, q_2)$ lo hace, al mismo tiempo, de 0 a 1. Existirá pues un valor de q_2 , q_2^* , tal que $P(2/1, q_2^*) = P(1/2, q_2^*)$. Esta es la solución minimax y es óptima, por cuanto no existiera otro procedimiento R^* tal que $\max \{P(2/1, R^*), P(1/2, R^*)\} \leq P(2/1, q_2^*) = P(1/2, q_2^*)$ esta contradeciría el hecho de ser toda solución Bayes admisible.

4. CLASIFICACIÓN EN UNA DE DOS POBLACIONES NORMALES MULTIVARIANTES CONOCIDAS

Utilizaremos el procedimiento general descrito anteriormente para el caso de dos poblaciones normales multivariantes con igual matriz de covarianza, a saber, $N(\mu^{(1)}, \Sigma)$ y $N(\mu^{(2)}, \Sigma)$, donde $\mu^{(i)} = (\mu_1^{(i)}, \dots, \mu_p^{(i)})$ es el vector de medias de la i -ésima población ($i=1,2$) y Σ la matriz de covarianzas común a ambas poblaciones. La densidad de la población i es

$$p_i(x) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \cdot \exp \left[-\frac{1}{2} (x - \mu^{(i)})' \Sigma^{-1} (x - \mu^{(i)}) \right].$$

La función de densidad viene dada por

$$\frac{p_1(x)}{p_2(x)} = \frac{\exp \left[-\frac{1}{2} (x - \mu^{(1)})' \Sigma^{-1} (x - \mu^{(1)}) \right]}{\exp \left[-\frac{1}{2} (x - \mu^{(2)})' \Sigma^{-1} (x - \mu^{(2)}) \right]} = \exp \left[-\frac{1}{2} \{ (x - \mu^{(1)})' \Sigma^{-1} (x - \mu^{(1)}) - (x - \mu^{(2)})' \Sigma^{-1} (x - \mu^{(2)}) \} \right].$$

La región de clasificación en π_1, π_2 es el conjunto de x 's para los que este cociente es $\geq k$ (k elegido adecuadamente). Puesto que la función logarítmica es una función monótona creciente, la desigualdad puede ser reescrita en términos de sus logaritmos como

$$-\frac{1}{2} \{ (x - \mu^{(1)})' \Sigma^{-1} (x - \mu^{(1)}) - (x - \mu^{(2)})' \Sigma^{-1} (x - \mu^{(2)}) \} \geq \log k.$$

Lo que desarrollado conduce a

$$-\frac{1}{2} [x' \Sigma^{-1} x - x' \Sigma^{-1} \mu^{(1)} - \mu^{(1)'} \Sigma^{-1} x + \mu^{(1)'} \Sigma^{-1} \mu^{(1)} - x' \Sigma^{-1} x + x' \Sigma^{-1} \mu^{(2)} + \mu^{(2)'} \Sigma^{-1} x - \mu^{(2)'} \Sigma^{-1} \mu^{(2)}] \geq \log k$$

y anulado adecuadamente, llegamos a

$$x' \Sigma^{-1} (\mu^{(2)} - \mu^{(1)}) - \frac{1}{2} (\mu^{(1)} + \mu^{(2)})' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)}) \geq \log k.$$

El primer sumando de la izquierda es conocido como la función discriminante, y es una función lineal de las componentes del vector de las observaciones.

Aplicando los resultados antes obtenidos llegamos al siguiente teorema.

TEOREMA - Si Π_i tiene una densidad normal multivariante $p_i(x)$, $i=1,2$, las mejores reglas de clasificación vienen dadas por

$$R_1: x' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)}) - \frac{1}{2} (\mu^{(1)} + \mu^{(2)})' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)}) \geq \log k$$

$$R_2: x' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)}) - \frac{1}{2} (\mu^{(1)} + \mu^{(2)})' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)}) < \log k.$$

Si las probabilidades a priori q_1 y q_2 son iguales, entonces k viene dado por

$$k = \frac{q_2 \sigma(1/2)}{q_1 \sigma(1/1)}.$$

El caso particular en que las dos poblaciones son igualmente probables y los costes son iguales, $k=1$ y $\log k=0$. Entonces la regla de clasificación en Π_1 es

$$R_1: x' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)}) \geq \frac{1}{2} (\mu^{(1)} + \mu^{(2)})' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)}).$$

Si no sabemos nada acerca de q_1 y q_2 podemos seleccionar $\log k = c$, sobre la base de igualar las pérdidas esperadas debidas a malclasificación.

Sea x una observación aleatoria. Queremos encontrar la distribución de

$$U = x' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)}) - \frac{1}{2} (\mu^{(1)} + \mu^{(2)})' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)}),$$

sobre el supuesto de que x se distribuye $N(\mu^{(1)}, \Sigma)$ y luego sobre el supuesto de que x lo hace $N(\mu^{(2)}, \Sigma)$. Cuando $x \sim N(\mu^{(1)}, \Sigma)$, U es normal con media

$$E_1(U) = \mu^{(1)'} \Sigma^{-1} (\mu^{(1)} - \mu^{(2)}) - \frac{1}{2} (\mu^{(1)} + \mu^{(2)})' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)}) = \frac{1}{2} (\mu^{(1)} - \mu^{(2)})' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)}),$$

y varianza

$$\text{var}_1(U) = E_1[(\mu^{(1)} - \mu^{(2)})' \Sigma^{-1} (x - \mu^{(1)}) (x - \mu^{(1)})' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)})] = (\mu^{(1)} - \mu^{(2)})' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)}).$$

Sea α la distancia entre las dos poblaciones, a saber

$$\alpha = (\mu^{(1)} - \mu^{(2)})' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)}).$$

Entonces U se distribuye $N(\frac{1}{2}\alpha, \alpha)$ cuando $x \sim N(\mu^{(1)}, \Sigma)$. Si $x \sim N(\mu^{(2)}, \Sigma)$ entonces

$$E_2(U) = \frac{1}{2} (\mu^{(2)} - \mu^{(1)})' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)}) = -\frac{1}{2}\alpha$$

viendo la varianza la misma, puesto que esto depende de los momentos de segundo orden de x y éstos son iguales en ambos casos. Así, U se distribuye $N(-\frac{1}{2}\alpha, \alpha)$ en este segundo caso.

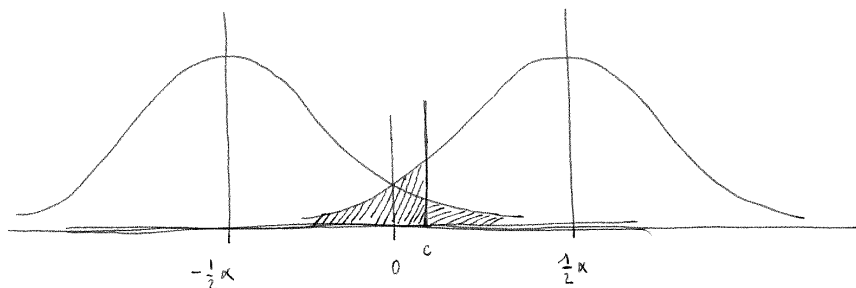
La probabilidad de malclasificación, si la observación proviene de Π_1 es

$$P(2/1) = \int_{-\infty}^c \frac{1}{\sqrt{2\pi\alpha}} e^{-\frac{1}{2}(z - \frac{1}{2}\alpha)^2/\alpha} dz = \int_{-\infty}^{(c - \frac{1}{2}\alpha)/\sqrt{\alpha}} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}y^2} dy$$

y la probabilidad de mala clasificación cuando la observación proviene de Π_2 es

$$P(1/2) = \int_c^{\infty} \frac{1}{\sqrt{2\pi\alpha}} e^{-\frac{1}{2}(z + \frac{1}{2}\alpha)^2/\alpha} dz = \int_{(c + \frac{1}{2}\alpha)/\sqrt{\alpha}}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}y^2} dy$$

ambas probabilidades aparecen representadas, en la figura de la página siguiente, como las áreas sombreadas de cada una de las colas adyacentes.



Para la solución minimax elegiremos c de manera que

$$G(1/2) \int_{(c-1/2)/\sqrt{x}}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}y^2} dy = G(1/2) \int_{-\infty}^{(c-1/2)/\sqrt{x}} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}y^2} dy.$$

Teorema. Si las poblaciones π_i tienen densidades multivariantes $p_i(x)$, $i=1,2$, la solución minimax viene dada por aquellas regiones que satisfacen las ecuaciones del anterior teorema con $\log k=c$ determinado para que verifique la igualdad anterior, es decir que $G(i/j)$ representan los dos tests de mala clasificación.

Antes de señalar que si los dos tests $G(i/j)$ son iguales, entonces $c=0$ y la probabilidad de mala clasificación es

$$\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}y^2} dy.$$

En caso de tests desiguales, c deberá ser determinado mediante un método aproximado a partir de las tablas de la normal.

Observemos finalmente que los discriminantes que aparecen en la expresión que determina R_1 y R_2 involucran al vector

$$\delta = \Sigma^{-1} (c \mu^{(1)} - \mu^{(2)}).$$

Es interesante señalar que $x'\delta$ es la función lineal que maximiza

$$\frac{[E_1(x'd) - E_2(x'd)]^2}{\text{var}(x'd)}$$

para todas las posibles elecciones de d . El numerador de este cociente es

$$(1) [\mu^{(1)'}d - \mu^{(2)'}d]^2 = d'[(\mu^{(1)} - \mu^{(2)})(\mu^{(1)} - \mu^{(2)})']d$$

mientras el denominador

$$(2) d' E[(x - E(x))(x - E(x))'] d = d' \Sigma d.$$

Podemos maximizar (1) manteniendo (2) constante. Si λ es un multiplicador de Lagrange, estamos buscando el máximo de

$$d'[(\mu^{(1)} - \mu^{(2)})(\mu^{(1)} - \mu^{(2)})']d - \lambda(d' \Sigma d - 1).$$

Derivando respecto a los componentes de d e igualando a cero obtenemos

$$2[(\mu^{(1)} - \mu^{(2)})(\mu^{(1)} - \mu^{(2)})']d = 2\lambda \Sigma d.$$

Puesto que $(\mu^{(1)} - \mu^{(2)})'d$ es un escalar, k , escribimos la anterior igualdad como

$$\mu^{(1)} - \mu^{(2)} = \frac{\lambda}{k} \Sigma d$$

y la solución es proporcional a δ .

Señalemos finalmente que si tenemos una muestra de tamaño N de cada una de las poblaciones π_1 y π_2 , podemos utilizar la media muestral y clasificarla como procedente de $N[\mu^{(1)}, \frac{1}{N} \Sigma]$ o $N[\mu^{(2)}, \frac{1}{N} \Sigma]$.

5. CLASIFICACION EN UNA DE DOS POBLACIONES NORMALES MULTIVARIANTES CUANDO LOS PARAMETROS SON ESTIMADOS

El criterio de clasificación.- Hasta ahora hemos supuesto que las poblaciones eran conocidas exactamente. En muchas de las aplicaciones de esta técnica las poblaciones no son conocidas, pero deben ser inferidas a partir de las muestras, una de cada población. No ocuparemos a continuación del caso en que tenemos una muestra de cada una de las poblaciones, ambas normales multivariantes, y deseamos utilizar esta información para clasificar otras observaciones como provenientes de una de las dos poblaciones.

Supongamos que tenemos una muestra $x_1^{(1)} \dots x_{N_1}^{(1)}$ de $N(\mu^{(1)}, \Sigma)$ y una muestra $x_1^{(2)} \dots x_{N_2}^{(2)}$ de $N(\mu^{(2)}, \Sigma)$. La mejor estimación de $\mu^{(1)}$ es $\bar{x}^{(1)} = \sum_{i=1}^{N_1} x_{\alpha}^{(1)} / N_1$ y de $\mu^{(2)}$ es $\bar{x}^{(2)} = \sum_{i=1}^{N_2} x_{\alpha}^{(2)} / N_2$ y la Σ es definida mediante

$$(N_1 + N_2 - 2) S = \sum_{i=1}^{N_1} (x_{\alpha}^{(1)} - \bar{x}^{(1)}) (x_{\alpha}^{(1)} - \bar{x}^{(1)})' + \sum_{i=1}^{N_2} (x_{\alpha}^{(2)} - \bar{x}^{(2)}) (x_{\alpha}^{(2)} - \bar{x}^{(2)})'$$

Sustituimos estos valores por los correspondientes parámetros en la expresión obtenida anteriormente para obtener

$$x' S^{-1} (\bar{x}^{(1)} - \bar{x}^{(2)}) - \frac{1}{2} (\bar{x}^{(1)} + \bar{x}^{(2)})' S^{-1} (\bar{x}^{(1)} - \bar{x}^{(2)}).$$

El primer término de esta suma es la función discriminante basada en dos muestras [originada por Fisher (1936)]. Es una función lineal que tiene la mayor "varianza entre muestras" relativa a la "varianza dentro de las muestras". La idea es utilizar la anterior expresión como criterio de clasificación de la misma forma a como fue utilizado en el caso de conocer los parámetros de las poblaciones.

Cuando las poblaciones eran conocidas, reargumentaba que el criterio de clasificación es el mejor en el sentido de que minimizaba la pérdida esperada en el caso de probabilidades a priori conocidas y generaliza la clase de procedimientos admisibles cuando las probabilidades a priori no son conocidas. Ahora no podemos justificarlos en el mismo sentido. Sin embargo, para intuitivamente mostrar que el criterio proporciona buenos resultados. Otro criterio será indicado en un capítulo posterior.

Supongamos que tenemos una muestra $x_1 \dots x_N$ de una u otra población, π_1 o π_2 y deseamos llevar a cabo una clasificación como un todo y no elemento a elemento. Definimos entonces S como

$$(N_1 + N_2 + N - 3) S = \sum_{i=1}^{N_1} (x_{\alpha}^{(1)} - \bar{x}^{(1)}) (x_{\alpha}^{(1)} - \bar{x}^{(1)})' + \sum_{i=1}^{N_2} (x_{\alpha}^{(2)} - \bar{x}^{(2)}) (x_{\alpha}^{(2)} - \bar{x}^{(2)})' + \sum_{i=1}^N (x_{\alpha} - \bar{x}) (x_{\alpha} - \bar{x})'$$

donde

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_{\alpha}.$$

Entonces el criterio es

$$[\bar{x} - \frac{1}{2} (\bar{x}^{(1)} + \bar{x}^{(2)})]' S^{-1} (\bar{x}^{(1)} - \bar{x}^{(2)}).$$

Puede demostrarse que a mayor valor de N , menor son las probabilidades de mala clasificación.

La distribución del criterio.- Sea

$$V = \bar{x}' S^{-1} (\bar{x}^{(1)} - \bar{x}^{(2)}) - \frac{1}{2} (\bar{x}^{(1)} + \bar{x}^{(2)})' S^{-1} (\bar{x}^{(1)} - \bar{x}^{(2)}) = [\bar{x} - \frac{1}{2} (\bar{x}^{(1)} + \bar{x}^{(2)})]' S^{-1} (\bar{x}^{(1)} - \bar{x}^{(2)})$$

con \bar{x} , $\bar{x}^{(1)}$, $\bar{x}^{(2)}$ y S aleatorios.

La distribución de V es extremadamente complicada. Depende de los tamaños de las muestras y de los datos de los $x_{\alpha}^{(1)}$ y $x_{\alpha}^{(2)}$. Sean ahora

$$Z = \bar{x} - \frac{1}{2} (\bar{x}^{(1)} + \bar{x}^{(2)})$$

$$Y = \bar{x}^{(1)} - \bar{x}^{(2)},$$

entonces

$$V = Z' S^{-1} Y.$$

La esperanza de Y es $\mu^{(1)} - \mu^{(2)}$ y la matriz de varianzas es $[1/N_1 + 1/N_2] \Sigma$. Z se distribuye normalmente con media

$$E_1(z) = \frac{1}{2} (\mu^{(1)} - \mu^{(2)})$$

si Σ procede de Π_1 , y

$$E_2(z) = \frac{1}{2} (\mu^{(2)} - \mu^{(1)})$$

si Σ procede de Π_2 . En ambos casos la matriz de varianzas es $[1 + 1/4N_1 + 1/4N_2] \Sigma$. La matriz de varianzas entre E_1 y E_2

$$= \left(\frac{1}{2N_1} + \frac{1}{2N_2} \right) \Sigma.$$

Si $N_1 = N_2$ esta matriz de varianzas es nula. En este caso cabe fácilmente que la distribución de V para Σ procedente de Π_1 es la misma que la de $-V$ para Σ procedente de Π_2 . Así, si $V \geq 0$ la región de clasificación en Π_2 , la probabilidad de clasificar erróneamente a Σ cuando proviene de Π_1 es igual a la probabilidad de hacerlo, igualmente, mal cuando proviene de Π_2 .

La distribución de V ha sido estudiada por Anderson, Stitgears, y Wald.

La distribución asintótica del criterio. - En el caso de grandes muestras para ambos poblaciones podemos aplicar la teoría de las distribuciones límite. Puesto que $\bar{X}^{(i)}$ es la media de la muestra en N_i variables independientes (diversas muestras) todas ellas con distribución $N(\mu^{(i)}, \Sigma)$, tenemos que

$$\lim_{N_1 \rightarrow \infty} \bar{X}^{(1)} = \mu^{(1)} \quad (p).$$

analogamente

$$\lim_{N_2 \rightarrow \infty} \bar{X}^{(2)} = \mu^{(2)} \quad (p)$$

y

$$\lim S = \Sigma \quad (p)$$

cundo N_1, N_2 o ambos tienden a ∞ . y de esta última igualdad se obtiene

$$\lim S^{-1} = \Sigma^{-1} \quad (p)$$

límite en probabilidad

puesto que los ~~productos~~ de sumas, diferencias, productos y cocientes de variables aleatorias con los sumas, productos y cocientes de los límites en probabilidad siempre que el límite de cada denominador sea diferente de cero. Entonces, aplicando todo esto

$$\lim_{N_1, N_2 \rightarrow \infty} S^{-1} (\bar{X}^{(1)} - \bar{X}^{(2)}) = \Sigma^{-1} (\mu^{(1)} - \mu^{(2)}) \quad (p)$$

$$\lim_{N_1, N_2 \rightarrow \infty} (\bar{X}^{(1)} + \bar{X}^{(2)})' S^{-1} (\bar{X}^{(1)} - \bar{X}^{(2)}) = (\mu^{(1)} + \mu^{(2)})' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)}) \quad (p).$$

Se sigue que la distribución asintótica de V es la la distribución de U estudiada en el párrafo anterior. Para muestras suficientemente grandes de Π_1 y Π_2 podemos utilizar el criterio como si conociéramos exactamente las poblaciones cometiendo sólo un pequeño error.

TEOREMA. - Sea V definido como antes. La distribución asintótica de V cuando $N_1 \rightarrow \infty$, $N_2 \rightarrow \infty$ es $N(\frac{1}{2}\alpha, \alpha)$ si Σ procede de la población Π_1 y es $N(-\frac{1}{2}\alpha, \alpha)$ si Σ procede de la población Π_2 . Donde $\alpha = (\mu^{(1)} - \mu^{(2)})' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)})$.

Otra derivación del criterio. - Una obtención convenientemente memorizable del criterio es el uso de la regresión de una variable testigo dado por Fisher en su trabajo de 1936. Sea

$$y_{\alpha}^{(1)} = \frac{N_2}{N_1 + N_2}, \quad \alpha = 1, \dots, N_1, \quad y_{\alpha}^{(2)} = \frac{-N_1}{N_1 + N_2}, \quad \alpha = 1, \dots, N_2.$$

Se trata de encontrar la regresión sobre las variables $x_{\alpha}^{(i)}$ eligiendo b de manera que se minimice

$$\sum_{i=1}^2 \sum_{\alpha=1}^{N_i} [y_{\alpha}^{(i)} - b'(x_{\alpha}^{(i)} - \bar{x})]^2$$

donde

$$\bar{x} = [N_1 \bar{x}^{(1)} + N_2 \bar{x}^{(2)}] / (N_1 + N_2).$$

Las "ecuaciones normales" son

$$\sum_{i=1}^2 \sum_{\alpha=1}^{N_i} (x_{\alpha}^{(i)} - \bar{x})(x_{\alpha}^{(i)} - \bar{x})' b = \sum_{i=1}^2 \sum_{\alpha=1}^{N_i} y_{\alpha}^{(i)} (x_{\alpha}^{(i)} - \bar{x}) = \frac{N_1 N_2}{N_1 + N_2} [(\bar{x}^{(1)} - \bar{x}) - (\bar{x}^{(2)} - \bar{x})] = \frac{N_1 N_2}{N_1 + N_2} (\bar{x}^{(1)} - \bar{x}^{(2)}).$$

La matriz que multiplica a b puede escribirse como

$$\begin{aligned} \sum_{i=1}^2 \sum_{\alpha=1}^{N_i} (x_{\alpha}^{(i)} - \bar{x})(x_{\alpha}^{(i)} - \bar{x})' &= \sum_{i=1}^2 \sum_{\alpha=1}^{N_i} (x_{\alpha}^{(i)} - \bar{x}^{(i)})(x_{\alpha}^{(i)} - \bar{x}^{(i)})' + N_1 (\bar{x}^{(1)} - \bar{x})(\bar{x}^{(1)} - \bar{x})' + N_2 (\bar{x}^{(2)} - \bar{x})(\bar{x}^{(2)} - \bar{x})' = \\ &= \sum_{i=1}^2 \sum_{\alpha=1}^{N_i} (x_{\alpha}^{(i)} - \bar{x}^{(i)})(x_{\alpha}^{(i)} - \bar{x}^{(i)})' + \frac{N_1 N_2}{N_1 + N_2} (\bar{x}^{(1)} - \bar{x}^{(2)})(\bar{x}^{(1)} - \bar{x}^{(2)})'. \end{aligned}$$

En definitiva las anteriores ecuaciones normales pueden escribirse como

$$G b = (\bar{x}^{(1)} - \bar{x}^{(2)}) \left[\frac{N_1 N_2}{N_1 + N_2} - \frac{N_1 N_2}{N_1 + N_2} (\bar{x}^{(1)} - \bar{x}^{(2)})' b \right]$$

donde

$$G = \sum_{i=1}^2 \sum_{\alpha=1}^{N_i} (x_{\alpha}^{(i)} - \bar{x}^{(i)})(x_{\alpha}^{(i)} - \bar{x}^{(i)})'.$$

Puesto que $(\bar{x}^{(1)} - \bar{x}^{(2)})' b$ es un escalar, vemos que la solución b de la anterior expresión es proporcional a

$$S^{-1}(\bar{x}^{(1)} - \bar{x}^{(2)}).$$

El criterio de la razón de verosimilitud. - Otro criterio que puede ser utilizado en la clasificación es el de la razón de verosimilitud.

Consideremos que contrastamos la hipótesis nula de que $x_1, x_1^{(1)}, \dots, x_{N_1}^{(1)}$ son extraídas de $N(\mu^{(1)}, \Sigma)$ y $x_1^{(2)}, \dots, x_{N_2}^{(2)}$ lo son de $N(\mu^{(2)}, \Sigma)$, frente a la alternativa de que $x_1^{(1)}, \dots, x_{N_1}^{(1)}$ pertenecen a $N(\mu^{(1)}, \Sigma)$, mientras que $x_1, x_1^{(1)}, \dots, x_{N_2}^{(2)}$ lo hacen a $N(\mu^{(1)}, \Sigma)$, con $\mu^{(1)}, \mu^{(2)}$ y Σ no especificados. Bajo la primera hipótesis los estadísticos máximos verosímiles de $\mu^{(1)}, \mu^{(2)}$ y Σ son

$$\hat{\mu}_1^{(1)} = (N_1 \bar{x}^{(1)} + x) / (N_1 + 1)$$

$$\hat{\mu}_1^{(2)} = \bar{x}^{(2)}$$

$$\hat{\Sigma}_1 = \frac{1}{N_1 + N_2 + 1} \left[\sum_{\alpha=1}^{N_1} (x_{\alpha}^{(1)} - \hat{\mu}_1^{(1)})(x_{\alpha}^{(1)} - \hat{\mu}_1^{(1)})' + (x - \hat{\mu}_1^{(1)})(x - \hat{\mu}_1^{(1)})' + \sum_{\alpha=1}^{N_2} (x_{\alpha}^{(2)} - \hat{\mu}_1^{(1)})(x_{\alpha}^{(2)} - \hat{\mu}_1^{(1)})' \right]$$

puesto que

$$\begin{aligned} \sum_{\alpha=1}^{N_1} (x_{\alpha}^{(1)} - \hat{\mu}_1^{(1)})(x_{\alpha}^{(1)} - \hat{\mu}_1^{(1)})' + (x - \hat{\mu}_1^{(1)})(x - \hat{\mu}_1^{(1)})' &= \sum_{\alpha=1}^{N_1} (x_{\alpha}^{(1)} - \bar{x}^{(1)})(x_{\alpha}^{(1)} - \bar{x}^{(1)})' + N_1 (\bar{x}^{(1)} - \hat{\mu}_1^{(1)})(\bar{x}^{(1)} - \hat{\mu}_1^{(1)})' + (x - \hat{\mu}_1^{(1)})(x - \hat{\mu}_1^{(1)})' = \\ &= \sum_{\alpha=1}^{N_1} (x_{\alpha}^{(1)} - \bar{x}^{(1)})(x_{\alpha}^{(1)} - \bar{x}^{(1)})' + \frac{N_1}{N_1 + 1} (x - \bar{x}^{(1)})(x - \bar{x}^{(1)})' \end{aligned}$$

podemos escribir $\hat{\Sigma}_1$ de la forma

$$\hat{\Sigma}_1 = \frac{1}{N_1 + N_2 + 1} \left[G + \frac{N_1}{N_1 + 1} (x - \bar{x}^{(1)})(x - \bar{x}^{(1)})' \right]$$

donde G es la del apartado anterior. Bajo los supuestos de la hipótesis alternativa encontramos (por consideraciones de simetría)

que los estimadores máximo verosímiles de los parámetros son

5 (6)

$$\hat{\mu}_2^{(1)} = \bar{x}^{(1)}$$

$$\hat{\mu}_2^{(2)} = (N_2 \bar{x}^{(2)} + x) / (N_2 + 1)$$

$$\hat{\Sigma}_2 = \frac{1}{N_1 + N_2 + 1} \left[G + \frac{N_2}{N_2 + 1} (x - \bar{x}^{(2)})(x - \bar{x}^{(2)})' \right]$$

El criterio de la unión de verosimilitud es, portanto, la potencia $(N_1 + N_2 + 1)/2$ -sima de

$$\frac{|\hat{\Sigma}_2|}{|\hat{\Sigma}_1|} = \frac{\left| G + \frac{N_2}{N_2 + 1} (x - \bar{x}^{(2)})(x - \bar{x}^{(2)})' \right|}{\left| G + \frac{N_1}{N_1 + 1} (x - \bar{x}^{(1)})(x - \bar{x}^{(1)})' \right|}$$

que puede también ser escrito de la forma

$$\frac{1 + \frac{N_2}{N_2 + 1} (x - \bar{x}^{(2)})' G^{-1} (x - \bar{x}^{(2)})}{1 + \frac{N_1}{N_1 + 1} (x - \bar{x}^{(1)})' G^{-1} (x - \bar{x}^{(1)})}$$

la región de clasificación en Π_1 consistirá en aquellos puntos para los que el numerador es mayor que una cantidad dada.

6. CLASIFICACION EN UNA DE VARIAS POBLACIONES

Consideremos ahora el problema de clasificar una observación en una de varias poblaciones. Extenderemos las consideraciones hechas en los apartados anteriores a los casos de más de dos poblaciones. Sean Π_1, \dots, Π_m , m poblaciones con funciones de densidad $p_1(x), \dots, p_m(x)$ respectivamente. Dividamos el espacio de observación en m regiones exhaustivas y mutuamente excluyentes, R_1, \dots, R_m . Si una observación cae en R_i diremos que procede de Π_i . Sean $G(j/i)$ los costes de mala clasificación de una observación que procede de la población Π_i como procedente de Π_j . La probabilidad de esta clasificación errónea es

$$P(j/i, R) = \int_{R_j} p_i(x) dx.$$

Supongamos que tenemos probabilidades a priori de las poblaciones, q_1, \dots, q_m . Entonces la pérdida esperada es

$$\sum_{i=1}^m q_i \left\{ \sum_{\substack{j=1 \\ i \neq j}}^m G(j/i) P(j/i, R) \right\}.$$

Desearíamos elegir R_1, \dots, R_m para hacer mínima esta expresión.

Puesto que tenemos probabilidades a priori para las poblaciones, podemos definir la probabilidad condicional de que una observación proceda de la población Π_i dados los valores de los componentes del vector x . Esta probabilidad viene dada por

$$\frac{q_i \cdot p_i(x)}{\sum_{k=1}^m q_k \cdot p_k(x)}.$$

Si clasificamos la observación como perteneciente a Π_j , la pérdida esperada es

$$\sum_{\substack{i=1 \\ i \neq j}}^m \frac{q_i \cdot p_i(x)}{\sum_{k=1}^m q_k \cdot p_k(x)} \cdot G(j/i)$$

minimizaremos la pérdida esperada en estos indicadores, si elegimos j de manera que minimicela anterior expresión; es decir, consideramos

$$\sum_{\substack{i=1 \\ i \neq j}}^m q_i \cdot p_i(x) G(j/i)$$

para todos los j y seleccionamos aquel j que proporciona el mínimo. (En el caso de igualdad la elección es irrelevante). Este procedimiento asigna el punto x a una de las R_j . Siguiendo este procedimiento para cada x , definimos m regiones R_1, \dots, R_m . El procedimiento de clasificación entón, es clasificar una observación como procedente de π_j si cae en R_j .

TEOREMA - Si q_i es la probabilidad a priori de extraer una observación de una población π_i con densidad $p_i(x)$, $i=1, \dots, m$, y si el coste de mal clasificar una observación de π_i como procedente de π_j , es $G(j/i)$, entón las regiones de clasificación R_1, \dots, R_m , que minimizan el coste esperado, se definen asignando x a R_k si

$$\sum_{\substack{i=1 \\ i \neq k}}^m q_i \cdot p_i(x) \cdot G(k/i) < \sum_{\substack{i=1 \\ i \neq j}}^m q_i \cdot p_i(x) \cdot G(j/i), \quad j=1, \dots, m, j \neq k.$$

Si esta desigualdad se verifica para todos los valores, excepto para R_i de ellos en los que se convierte en una igualdad, el punto en cuestión puede ser asignado a cualquiera de las R_i+1 poblaciones correspondiente.

Si la probabilidad de que se verifique la igualdad, en lugar de la desigualdad stricta, es cero para cada k, j bajo π_i (para cada i), entón el procedimiento es único casi en todas partes.

Verifiquemos ahora este resultado. Sea

$$h_j(x) = \sum_{\substack{i=1 \\ i \neq j}}^m q_i \cdot p_i(x) \cdot G(j/i).$$

La pérdida esperada de un procedimiento R es

$$\sum_{j=1}^m \int_{R_j} h_j(x) dx = \int h(x) dx$$

donde $h(x) = h_j(x)$ para x en R_j . Para el procedimiento Bayes descrito en el teorema $h_i(x)$ es $h_i^*(x) = \min_i h_i(x)$. Así la diferencia entre la pérdida esperada para cualquier procedimiento R y la de R^* es

$$\int [h(x) - \min_i h_i(x)] dx = \sum_j \int_{R_j} [h_j(x) - \min_i h_i(x)] dx \geq 0,$$

verificándose la igualdad sólo si $h_j(x) = \min_i h_i(x)$ para x en R_j , excepto en hechos de conjunto de probabilidad cero.

Veamos como aplicar este método en aquella situación para la que $G(j/i) = 1$, $R_i, R_j, i \neq j$. Entón en R_k

$$\sum_{\substack{i=1 \\ i \neq k}}^m q_i \cdot p_i(x) < \sum_{\substack{i=1 \\ i \neq j}}^m q_i \cdot p_i(x) \quad j \neq k.$$

Restando $\sum_{\substack{i=1 \\ i \neq j, k}}^m q_i \cdot p_i(x)$ de ambos miembros de la desigualdad, tendremos

$$q_j \cdot p_j(x) < q_k \cdot p_k(x) \quad j \neq k.$$

En este caso el punto x está en R_k si k es el índice para el que $q_i \cdot p_i(x)$ es un máximo, es decir π_k es la población más probable.

Supongamos ahora que no tenemos probabilidades a priori, entón no podemos definir una pérdida esperada normal para un procedimiento de clasificación. Sin embargo, podemos definir una pérdida esperada sobre la condición de que las observaciones provienen de una población dada. La pérdida esperada si las observaciones provienen de π_i es

$$\sum_{\substack{j=1 \\ j \neq i}}^m G(j/i) \cdot P(j/i, R) = r(i, R)$$

Un procedimiento R es al menos tan bueno como R^* si $r(i, R) \leq r(i, R^*)$, $i=1, \dots, m$; R es mejor si al menos

una desigualdad estricta. R es admisible si no hay ningún procedimiento R^* que sea mejor. Una clase de procedimientos es completa si para cualquier procedimiento R exterior a la clase existe un procedimiento R^* en la clase que es mejor. 5 (7)

Veamos ahora que un procedimiento Bayes es admisible. Sea R un procedimiento Bayes y R^* otro procedimiento cualquiera. Puesto que R es Bayes

$$\sum_{i=1}^m q_i \cdot r(i, R) \leq \sum_{i=1}^m q_i \cdot r(i, R^*)$$

Supongamos que $r(i, R^*) \leq r(i, R)$ $i=1, \dots, m$ y $q_i > 0$. Entonces

$$q_1 [r(1, R) - r(1, R^*)] \leq \sum_{i=2}^m q_i [r(i, R^*) - r(i, R)] \leq 0$$

y $r(1, R^*) \geq r(1, R)$. Análogamente, si $q_i > 0$, y $r(i, R^*) \leq r(i, R)$ para $j \neq i$, entonces $r(i, R) \leq r(i, R^*)$. En definitiva, R^* no puede ser mejor que R y R es admisible.

TEOREMA.- Si $q_i > 0$, $i=1, \dots, m$, entonces R es un procedimiento admisible.

Respondemos ahora que $G(i/j) = 1$, $i \neq j$, y $P\{P_i(x) = 0/\pi_j\} = 0$. Esta última condición implica que las $p_i(x)$ son positivas sobre el mismo conjunto (exceptuando dicho conjunto de probabilidad nula). Supongamos ahora que $q_i = 0$ para $i=1, \dots, t$, $q_i > 0$ para $i=t+1, \dots, m$. Entonces para la reducción Bayes, R_i ($i=1, \dots, t$) es vacía (excepto para conjuntos de probabilidad nula) como se comprueba fácilmente de la definición dada para R_i bajo esta situación particular, anteriormente (lo significa $p_i(x) = 0, x \in R_i$).

Se sigue entonces que $r(i, R) = \sum_{j \neq i} P(j/i, R) = 1 - P(i/i, R) = 1$ para $i=1, \dots, t$. Entonces R_{t+1}, \dots, R_m es una reducción Bayes para el problema que involucra a $p_{t+1}(x), \dots, p_m(x)$ y q_{t+1}, \dots, q_m . Del teorema que acabamos de demostrar se sigue que no hay ningún procedimiento R^* para el que $P(i/i, R^*) = 0$, $i=1, \dots, t$, que sea mejor que el procedimiento Bayes. Consideremos ahora un procedimiento R^* tal que R^* incluya un conjunto de probabilidad distinta de cero tal que $P(1/1, R^*) > 0$. Para que R^* fuera mejor que R

$$P(i/i, R) = \int_{R_i} p_i(x) dx \leq \int_{R^*} p_i(x) dx = P(i/i, R^*) \quad i=1, \dots, m$$

En tal caso, un procedimiento R^{**} donde R_i^{**} es vacío para $i=1, \dots, t$, $R_i^{**} = R_i^*$, $i=t+1, \dots, m-1$ y $R_m^{**} = R_m^* \cup R_1^* \cup \dots \cup R_t^*$ daría riesgos tales como

$$\begin{aligned} P(i/i, R^{**}) &= 0 & i=1, \dots, t \\ P(i/i, R^{**}) &= P(i/i, R^*) \geq P(i/i, R) & i=t+1, \dots, m-1 \\ P(m/m, R^{**}) &> P(m/m, R^*) \geq P(m/m, R), \end{aligned}$$

entonces $(R_{t+1}^{**}, \dots, R_m^{**})$ sería mejor que (R_{t+1}, \dots, R_m) para el problema de decisión sobre $m-t$ problemas, lo que contradice la discusión precedente.

TEOREMA.- Si $G(i/j) = 1$, $i \neq j$ y $P\{P_i(x) = 0/\pi_j\} = 0$, entonces todo procedimiento Bayes es admisible.

Problemas ahora que los procedimientos admisibles son procedimientos Bayes. Prestaremos mucha atención al caso $m=3$. Supondremos que

$$P\left\{\frac{P_i(x)}{P_j(x)} = k/\pi_k\right\} = 0 \quad i \neq j, \quad 0 \leq k < \infty$$

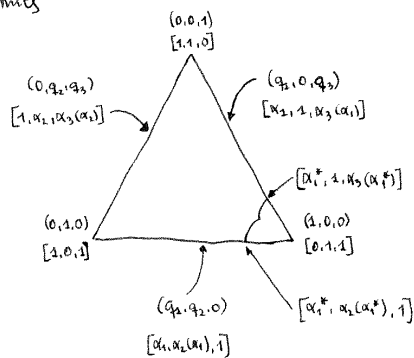
Esto implica que las ~~distribuciones~~ de probabilidad de $P_i(x)/P_j(x)$, para cualquier π_k , son continuas y que la familia de dens de los mismos, también continua.

Sea $\alpha_i(R) = 1 - P(i/i, R)$ la probabilidad de efectuar una decisión errónea cuando se utiliza el procedimiento R y se muestra en π_i . Cuando R es un procedimiento Bayes, $\alpha_i(R)$ es una función de q_1, q_2, q_3 , decir $\alpha(q_1, q_2, q_3)$. Es una función continua de q_1, q_2, q_3 ; por ejemplo

$$\alpha(q_1, q_2, q_3) = 1 - P\left\{\frac{P_2(x)}{P_1(x)} \leq \frac{q_1}{q_2}, \frac{P_3(x)}{P_1(x)} \leq \frac{q_1}{q_3} \mid \pi_1\right\},$$

viendo entonces las distribuciones de $\frac{P_2(x)}{P_1(x)}$ y $\frac{P_3(x)}{P_1(x)}$. Es conveniente pensar en (q_1, q_2, q_3) como las coordenadas baricéntricas de

en punto. los límites



del espacio de tripletas, y los valores de las funciones sobre dichos límites se indican en la figura.

Sea ahora R^* un procedimiento admisible, y sea $\alpha_i(R^*) = \alpha_i^*$. Demostremos que R^* es un procedimiento Bayes. Consideremos la totalidad de los procedimientos Bayes para los que (q_1, q_2, q_3) dan lugar a $\alpha_i(q_1, q_2, q_3) = \alpha_i^*$. Cuando $q_3 = 0$, tenemos en efecto un problema de dos decisiones y entonces $\alpha_2 = \alpha_2(\alpha_1^*)$ que es el menor de los α_2 dado $\alpha_1 = \alpha_1^*$ (obtenido de los resultados para el caso de dos poblaciones); así $\alpha_2(\alpha_1^*) \leq \alpha_2^*$; y $\alpha_3 = 1$. Análogamente, cuando $q_2 = 0$, $\alpha_3 = \alpha_3(\alpha_1^*) \leq \alpha_3^*$ y $\alpha_1 = 1$. El conjunto de puntos q_1, q_2, q_3 para los que $\alpha_i(q_1, q_2, q_3) = \alpha_i^*$ es una curva continua⁽¹⁾ desde el punto con $[\alpha_1^*, \alpha_2(\alpha_1^*), 1]$ hasta el punto con $[\alpha_1^*, 1, \alpha_3(\alpha_1^*)]$. Puesto que α_2 varía continuamente desde $\alpha_2(\alpha_1^*)$ hasta 1, existirá un punto donde $\alpha_2 = \alpha_2^*$. Existe un procedimiento Bayes \bar{R} tal que $\alpha_i(\bar{R}) = \alpha_i^*$, y $\alpha_2(\bar{R}) = \alpha_2^*$. Puesto que \bar{R} es admisible, por el teorema anterior, $\alpha_3(\bar{R}) \leq \alpha_3^*$. Pero puesto que R^* es admisible $\alpha_3(R^*) = \alpha_3^*$. Por la unicidad de las soluciones Bayes tenemos que $R^* = \bar{R}$.

TEOREMA. - Si $P \left\{ \frac{P_i(x)}{P_j(x)} = k / \pi_k \right\} = 0$ $\forall j, 0 \leq k < \infty$, entonces cualquier procedimiento admisible es un procedimiento Bayes.

La demostración del anterior teorema demuestra que la clase de los procedimientos Bayes es completa. Para cualquier procedimiento dado R^* , existe un procedimiento Bayes \bar{R} que es al menos tan bueno (lo que supone la competitividad secuencial). Pero si \bar{R} y R^* son igualmente buenos entonces son el mismo (excepto para conjuntos de probabilidad cero).

TEOREMA. - Si $P \left\{ \frac{P_i(x)}{P_j(x)} = k / \pi_k \right\} = 0$, la clase de los procedimientos Bayes es minimal completa.

Consideremos también ahora la solución minimax. Existe una solución Bayes para la cual $\alpha_1 = \alpha_2 = \alpha_3$. Puesto que el procedimiento es admisible, no existe otro que tenga una menor probabilidad máxima de error (es decir, cada uno de los riesgos sea menor). Esto proporciona el procedimiento minimax.

Se pueden consultar los libros antes citados para completar esta teoría. Añadiré tan solo que von Neuman (1945) encontró la solución del problema minimax por un camino diferente.

7. CLASIFICACION EN UNA DE VARIAS POBLACIONES NORMALES MULTIVARIANTES

Aplicaremos la teoría del anterior párrafo al caso en que cada una de las poblaciones tiene una distribución normal. Supondremos que las medias son diferentes y las matrices de varianzas coincidentes. Sea $N(\mu^{(i)}, \Sigma)$ la distribución de Π_i . En un principio supondremos conocidos los parámetros. Para más generalidad y propósitos a priori consideraremos formar m funciones como las del párrafo anterior y definir las regiones R_j como aquellas constituidas por puntos x tales que la j -ésima función es mínima.

Lo que sigue aprenderemos que los tests de mala clasificación son iguales. Entonces utilizaremos las funciones

$$U_{jk}(x) = \log \frac{P_j(x)}{P_k(x)} = [x - \frac{1}{2}(\mu^{(j)} + \mu^{(k)})]' \Sigma^{-1} (\mu^{(j)} - \mu^{(k)})$$

Si las probabilidades a priori son conocidas, las regiones R_j vienen definidas por aquellos x que satisfacen

$$R_j: U_{jk}(x) > \log \frac{\pi_k}{\pi_j}, \quad k=1, \dots, m; \quad k \neq j.$$

(1) A lo largo de cada eje $q_3 = (1-q_1)(1-q_2)$ y $q_2 = k(1-q_1)$, $0 < k \leq 1$, α_1 decrece continuamente y monótonamente de 1 a 0. Sea $q_1 = q_1(k)$ el valor de q_1 tal que $\alpha_1 = \alpha_1^*$; entonces $q_1(k)$ es una función continua de k [por la continuidad de $\alpha_1(q_1, q_2, q_3)$ y la monotonía de α_1 como función de q_1 dado k].

TEOREMA - Si g_i es la probabilidad a priori de obtener una observación de $\pi_i \sim N(\mu^{(i)}, \Sigma)$ ($i=1, \dots, m$) y si los costos de J (8) malclasificación son iguales, entonces las regiones de clasificación R_1, \dots, R_m , que minimizan el costo esperado vienen definidas por la relación anterior, siendo $U_{jk}(x)$ definida como antes.

Hay que señalar que cada $U_{jk}(x)$ es la función de clasificación relacionada con las poblaciones j -ésima y k -ésima, y $U_{jk}(x) = -U_{kj}(x)$. Puesto que se trata de funciones lineales, la región R_i está acotada por hiperplanos. Si las medias definen un hiperplano $(m-1)$ -dimensional (por ejemplo, si los vectores media $\mu^{(i)}$ son linealmente independientes y $p \geq m-1$), entonces R_i está acotada por $m-1$ hiperplanos.

En el caso en que no se conozcan probabilidades a priori, la región R_j viene definida mediante las desigualdades

$$U_{jk}(x) \geq c_j - c_k, \quad k=1, \dots, m; \quad k \neq j$$

las constantes c_k pueden tomarse no negativas. Este conjunto de regiones forman la clase de procedimientos admisibles. Para el procedimiento π minimax, estas constantes se eligen de manera que $P(\pi/i, R)$ sean iguales.

Veamos ahora como obtener las probabilidades de una clasificación correcta. Si x es una observación aleatoria, consideremos las variables aleatorias

$$U_{ji} = [x - \frac{1}{2}(\mu^{(i)} + \mu^{(j)})]' \Sigma^{-1} (\mu^{(i)} - \mu^{(j)}).$$

Con $U_{ji} = -U_{ij}$. Utilizamos $m(m-1)/2$ funciones de clasificación si las medias definen un hiperplano $(m-1)$ -dimensional. Si x procede de π_j , entonces U_{ji} se distribuye como $N(\frac{1}{2}\alpha_{jii}, \alpha_{jii})$ donde

$$\alpha_{jii} = (\mu^{(i)} - \mu^{(i)})' \Sigma^{-1} (\mu^{(i)} - \mu^{(i)}).$$

la variancia entre U_{ji} y U_{jk} es

$$\alpha_{jki} = (\mu^{(i)} - \mu^{(k)})' \Sigma^{-1} (\mu^{(i)} - \mu^{(k)}).$$

Para determinar las constantes c_j consideramos las integrales

$$P(j/j, R) = \int_{c_j - c_m}^{\infty} \dots \int_{c_j - c_1}^{\infty} f_j du_{j1} \dots du_{j,j-1} du_{j,j+1} \dots du_{jm}.$$

donde f_j es la densidad de U_{ji} ($i=1, 2, \dots, m$), ($i \neq j$).

TEOREMA - Si $\pi_i \sim N(\mu^{(i)}, \Sigma)$ y los costos de mala clasificación son iguales, las regiones de clasificación, R_1, \dots, R_m , que minimizan la máxima pérdida esperada condicional, están definidas mediante $U_{jk}(x) \geq c_j - c_k$, $k=1, \dots, m$; $k \neq j$, donde $U_{jk}(x)$ está definido como antes. Las constantes c_j se determinan de manera que las integrales que dan los valores de $P(j/j, R)$ sean iguales.

Como un ejemplo consideremos el caso $m=3$. No hay pérdida de generalidad si tomamos $p=2$, para densidades en mayor p podemos proyectar sobre el plano bidimensional determinado por las medias de las tres poblaciones o no son colineales (o sea, podemos transformar el vector x en u_1, u_2 y $p-2$ otras coordenadas, donde estas últimas $p-2$ componentes redistribuyen independientemente de u_1, u_2 con medias cero). Las regiones R_j están determinadas mediante tres vectores, tal como se muestra en la

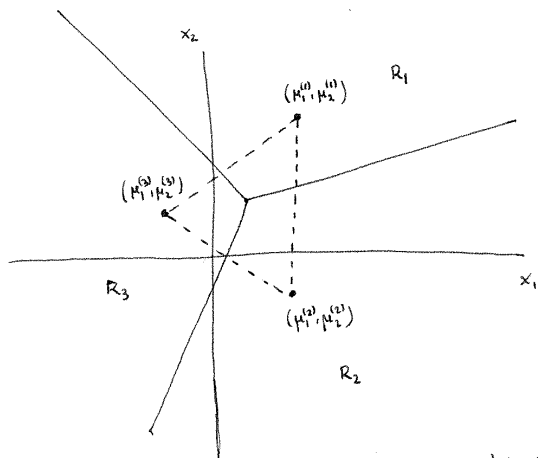


figura. Si el procedimiento es minimax, no podemos desplazar la recta entre R_1 y R_2 más cerca de $(\mu^{(1)}, \mu^{(2)})$, la recta entre R_2 y R_3 más cerca de $(\mu^{(2)}, \mu^{(3)})$ y la recta entre R_3 y R_1 más cerca de $(\mu^{(3)}, \mu^{(1)})$ y conservar todavía la igualdad $P(1/1, R) = P(2/2, R) = P(3/3, R)$, sin dejar un triángulo que no esté incluido en región alguna. Por tanto, como las regiones deben cubrir totalmente el espacio, las líneas deben encontrarse en un punto, y la igualdad de probabilidades determina $c_i = c_j$ universalmente.

Para realizar esto en un caso específico en el que tenemos valores numéricos para los vectores $\mu^{(1)}, \mu^{(2)}, \mu^{(3)}$ y la matriz Σ , consideraremos las tres ($\leq p+1$) distribuciones conjuntas, para los U_{ij} ($j \neq i$). Probaremos universalmente con los valores $c_i = 0$, y utilizando las tablas de Pearson de la muestra bivariante calcularemos $P(i/i, R)$, mediante

un método iterativo obtendremos los c_i que más se aproximan a la condición.

La teoría precedente parte del supuesto de conocer los parámetros de las poblaciones. Si no fueran conocidos y si tenemos disponible una muestra de cada población, substituiríamos en $U_{ij}(x)$ cada parámetro por su correspondiente estimación. Sean $x_1^{(1)}, \dots, x_{n_1}^{(1)}$

las observaciones procedentes de $N(\mu^{(i)}, \Sigma)$, $i=1, \dots, m$. Estimamos $\mu^{(i)}$ mediante

$$\bar{x}^{(i)} = \frac{1}{N_i} \sum_{\alpha=1}^{N_i} x_{\alpha}^{(i)}$$

y Σ mediante S definida por

$$\left(\sum_{i=1}^m N_i - m \right) S = \sum_{i=1}^m \sum_{\alpha=1}^{N_i} (x_{\alpha}^{(i)} - \bar{x}^{(i)}) (x_{\alpha}^{(i)} - \bar{x}^{(i)})'$$

Entonces, el análogo de $U_{ij}(x)$ es

$$U_{ij}(x) = [x - \frac{1}{2}(\bar{x}^{(i)} + \bar{x}^{(j)})]' S^{-1} (\bar{x}^{(i)} - \bar{x}^{(j)}).$$

Si las variables son aleatorias, las distribuciones son diferentes de aquellas de los U_{ij} . No obstante, cuando $N_i \rightarrow \infty$, la distribución conjunta se aproxima a la de U_{ij} . Por tanto, para muestras empíricamente grandes podemos utilizar la teoría anterior.

8. UN EJEMPLO DE CLASIFICACIÓN EN UNA DE VARIAS POBLACIONES NORMALES MULTIVARIANTES

Rao (1948) considera tres poblaciones consistentes en la casta de los Brahmin (π_1), la casta de los Artesanos (π_2) y la casta de los Korwa (π_3) de la India. Las medidas para cada uno de los individuos de una casta son, la estatura (x_1), altura del hombro (x_2), anchura nasal (x_3) y altura nasal (x_4). Las medidas de estas variables en las tres poblaciones aparecen en la tabla.

	Brahmin π_1	Artesanos π_2	Korwa π_3
Estatura (x_1)	164.51	160.53	158.17
Altura del hombro (x_2)	86.43	81.42	81.16
anchura nasal (x_3)	28.49	23.84	21.44
altura nasal (x_4)	51.24	48.62	46.72

La matriz de covarianzas para todas las poblaciones es

1.0000	0.5849	0.1774	0.1974
	1.0000	0.2094	0.2170
		1.0000	0.2910
			1.0000

Las desviaciones típicas son $\sigma_1 = 5.74$, $\sigma_2 = 3.20$, $\sigma_3 = 1.75$, $\sigma_4 = 3.50$. Suponemos que cada población es normal. Nuestro problema es dividir el espacio de las cuatro variables x_1, x_2, x_3, x_4 en tres regiones de clasificación. Supondremos que los costes de malclasificación son iguales. Encontraremos a) un conjunto de regiones bajo el supuesto de que extraer una muestra representativa de cada población es igualmente probable ($q_1 = q_2 = q_3 = 1/3$), y b) un conjunto de regiones tal que la mayor probabilidad de mala clasificación es minimizada (solución \minimax).

Calcularemos en primer lugar los coeficientes de $\Sigma^{-1}(\mu^{(1)} - \mu^{(2)})$ y $\Sigma^{-1}(\mu^{(1)} - \mu^{(3)})$. Luego $\Sigma^{-1}(\mu^{(2)} - \mu^{(3)}) = \Sigma^{-1}(\mu^{(1)} - \mu^{(3)}) - \Sigma^{-1}(\mu^{(1)} - \mu^{(2)})$. Luego calculamos $\frac{1}{2}(\mu^{(2)} + \mu^{(3)})' \Sigma^{-1}(\mu^{(1)} - \mu^{(2)})$. Los resultados obtenidos son las siguientes funciones discriminantes

$$u_{12}(x) = -0.0708x_1 + 0.4990x_2 + 0.3373x_3 + 0.0887x_4 - 43.13$$

$$u_{13}(x) = 0.0003x_1 + 0.3550x_2 + 1.1063x_3 + 0.1375x_4 - 62.49$$

$$u_{23}(x) = 0.0711x_1 - 0.1440x_2 + 0.7690x_3 + 0.0488x_4 - 19.36$$

Las otras tres funciones son $u_{21}(x) = -u_{12}(x)$, $u_{31}(x) = -u_{13}(x)$ y $u_{32}(x) = -u_{23}(x)$. Si existen probabilidades a priori y son iguales el mejor conjunto de ~~regiones de clasificación~~ regiones de clasificación es

$$R_1: u_{12}(x) \geq 0, u_{13}(x) \geq 0$$

$$R_2: u_{21}(x) \geq 0, u_{23}(x) \geq 0$$

$$R_3: u_{31}(x) \geq 0, u_{32}(x) \geq 0$$

Por ejemplo se obtienen un Brahmin individual x , con $u_{12}(x) \geq 0$ y $u_{13}(x) \geq 0$, lo clasificaremos como Brahmin.

Para encontrar las probabilidades de malclasificación cuando un individuo es extraído de la población π_i , necesitamos las

Población de X	u	medias	Desviaciones típicas	Conclusion
π_1	u_{12}	1.491	1.727	0.8658
	u_{13}	3.487	2.641	
π_2	u_{21}	1.491	1.727	-0.3894
	u_{23}	1.031	1.436	
π_3	u_{31}	3.487	2.641	0.3983
	u_{32}	1.031	1.436	

Las probabilidades de mala clasificación se obtienen mediante el uso de tablas para la normal bivariante. Estas probabilidades son 0.21 para π_1 , 0.42 para π_2 y 0.25 para π_3 . Por ejemplo, si las medidas son hechas sobre un Brahmin, la probabilidad de que sea clasificado como Antezano o Korwa es 0.21.

La solución minimax se obtiene encontrando las constantes c_1, c_2 y c_3 de manera que las probabilidades de una mala clasificación sean iguales. Las regiones de clasificación son:

$$R'_1: u_{12}(x) \geq 0.54, \quad u_{13}(x) \geq 0.29;$$

$$R'_2: u_{21}(x) \geq -0.54, \quad u_{23}(x) \geq -0.25;$$

$$R'_3: u_{31}(x) \geq -0.29, \quad u_{32}(x) \geq 0.25.$$

La probabilidad de mala clasificación común es 0.30. Así, la máxima probabilidad de error ha sido reducida de 0.42 a 0.30.

EJEMPLOS

- ① En un ejemplo del capítulo anterior hallabamos de un grupo de 49 ancianos que provenían de dos poblaciones, la π_1 , factor senil ausente, y la π_2 , factor senil presente. Recordemos que los datos obtenidos eran

	Población	
	Factor senil ausente	Factor senil presente
$\bar{x}_1^{(1)}$	12.87	8.75
$\bar{x}_2^{(1)}$	9.57	5.33
$\bar{x}_3^{(1)}$	11.49	8.50
$\bar{x}_4^{(1)}$	7.97	4.75
$N_1 = 37$		$N_2 = 12$

$$S = \begin{bmatrix} 11.2553 & 9.4042 & 7.1489 & 3.3830 \\ & 13.3830 & 7.3830 & 2.5532 \\ & & 11.5744 & 2.6170 \\ & & & 5.8085 \end{bmatrix}$$

hemos a calcular la función discriminante lineal para los dos grupos poblacionales. El vector diferencia de medias es

$$(\bar{x}_1 - \bar{x}_2)' = [3.82, 4.24, 2.99, 3.22]$$

La función discriminante lineal viene dada por

$$y = 0.030x_1 + 0.204x_2 + 0.010x_3 + 0.443x_4.$$

Como que los tests segundo y cuarto, similitud y figuras completas, respectivamente, dominan la función, mientras que los tests de información y abstracción tiene una influencia casi despreciable a la función y. El investigador podría centrarse en aquellos dos como indicadores de la calidad del factor senil.

El término independiente puede obtenerse, bien a partir de la expresión

$$\frac{1}{2} (\bar{x}_1 + \bar{x}_2)' S^{-1} (\bar{x}_1 - \bar{x}_2) = 4.76$$

o bien como el punto medio de los valores medio de la función discriminante en un grupo y otro. Así, $\bar{y}_1 = 5.97, \bar{y}_2 = 3.54$

$$\frac{\bar{y}_1 + \bar{y}_2}{2} = 4.76.$$

Si suponemos que nada sabemos acerca de las probabilidades a priori q_1 y q_2 de las poblaciones π_1, π_2 , respectivamente, ni además trabajamos en la hipótesis de var. iguales, $C(1/2) = C(2/1)$, la solución minimax conduce a un punto de discriminación $c=0$, por tanto:

$$R_1: y - \frac{1}{2}(\bar{x}_1 + \bar{x}_2)'S^{-1}(\bar{x}_1 - \bar{x}_2) \geq 0$$

$$R_2: \text{---} < 0$$

o bien

$$R_1: y(x) \geq 4.76$$

$$R_2: y(x) < 4.76$$

Aplicando esta regla a los individuos de las dos muestras utilizadas, obtendríamos la siguiente clasificación de los mismos.

		Diagnóstico Psiquiátrico		
		F.S.A.	F.S.P.	Total
Clasificación mediante la función discriminante	F.S.A.	29	4	33
	F.S.P.	8	8	16
Total		37	12	49

La tabla da una medida de la capacidad de la función discriminante lineal para reproducir el diagnóstico. Obsérvese que la estimación de las probabilidades teóricas de mala clasificación mediante las correspondientes proporciones muestrales no son muy adecuadas por cuanto en un caso, $P(2/1)$ ni concuerdan, pero no así en el otro $P(1/2)$. En efecto los valores de α , necesario para conocer la distribución de $U = [x - \frac{1}{2}(\mu^{(1)} + \mu^{(2)})]'S^{-1}[\mu^{(1)} - \mu^{(2)}]$, puede estimarse mediante

$$\hat{\alpha} = (\bar{x}^{(1)} - \bar{x}^{(2)})'S^{-1}(\bar{x}^{(1)} - \bar{x}^{(2)}) = 2.42$$

U tiene una distribución $N(\frac{1}{2}\alpha, \alpha)$ si $x \in \pi_1$ y $N(-\frac{1}{2}\alpha, \alpha)$ si $x \in \pi_2$. Utilizando la estimación encontrada para α , tendríamos

$$U \approx N(\frac{1}{2}\hat{\alpha}, \hat{\alpha}) \text{ si } x \in \pi_1$$

$$U \approx N(-\frac{1}{2}\hat{\alpha}, \hat{\alpha}) \text{ si } x \in \pi_2$$

entonces

$$P(2/1) = P(U \leq 0) = P\left(\frac{U - \frac{1}{2}\hat{\alpha}}{\sqrt{\hat{\alpha}}} \leq \frac{-1.21}{\sqrt{2.42}}\right) = P(Z \leq -0.78) = \Phi(-0.78) \approx 0.2177$$

analogamente

$$P(1/2) \approx 0.2177$$

Las estimaciones, mediante la falta de estas cantidades vienen

$$\hat{P}(2/1) = \frac{8}{37} \approx 0.2162$$

pero

$$\hat{P}(1/2) = \frac{4}{12} \approx 0.3333 \text{ que impone una estimación por exceso.}$$

De el ejemplo que sigue damos una alternativa a este método, para tratar de mejorar las probabilidades de mala clasificación tan dadas.

- ② La información acerca de las frecuencias relativas con que las observaciones de las dos poblaciones son recuentadas pueden ser de utilidad en clasificación de estas observaciones mediante funciones discriminantes. Vamos a reclassificar los individuos del ejemplo anterior haciendo uso de una regla Bayes. Necesitamos para ello probabilidades a priori para cada población, por lo que llegaremos a ello una estimación de las mismas mediante las correspondientes frecuencias relativas, es decir: $\hat{q}_1 = \frac{37}{49}$, $\hat{q}_2 = \frac{12}{49}$. En estos $G(1/2)$ y $G(2/1)$ se vuelve a clasificar los individuos result.

Osi las cosas, la regla de clasificación viene dada por:

$$R_1: y(x) - \frac{1}{2}(\bar{x}_1 + \bar{x}_2)' S^{-1}(\bar{x}_1 - \bar{x}_2) \geq \log \frac{q_2}{q_1}$$

$$R_2: y(x) - \frac{1}{2}(\bar{x}_1 + \bar{x}_2)' S^{-1}(\bar{x}_1 - \bar{x}_2) < \log \frac{q_2}{q_1}$$

$$\log \frac{q_2}{q_1} = \log \frac{12}{37} = -1.126$$

o sea

$$R_1: y(x) \geq 3.634$$

$$R_2: y(x) < 3.634$$

la aplicación de esta regla induce al siguiente cuadro de clasificación:

		Diagnostico Psiquiatrico		
		F.S.A	F.S.P	
clasificación mediante la función discrimi- nante	F.S.A	37	5	42
	F.S.P	0	7	7
		37	12	49

La introducción de las probabilidades a priori ha disminuido considerablemente, menos de la mitad, el número de individuos mal clasificados que se obtenían mediante la regla utilizada del ejemplo anterior.

Capítulo 6

La distribución de la matriz de covarianzas muestral y de la varianza muestral generalizada

VARIANZA MUESTRAL GENERALIZADA1. INTRODUCCIÓN

La matriz de covarianzas muestral es $S = [1/(N-1)] \sum (x_i - \bar{x})(x_i - \bar{x})'$ y es un estimador de la matriz de covarianzas poblacional Σ . Hemos estudiado anteriormente la distribución de probabilidad de $A = (N-1)S$ para el caso 2×2 . Lo que haremos ahora es generalizar este resultado para una matriz A de cualquier orden. Cuando $\Sigma = I$, esta distribución es en un sentido una generalización de la distribución χ^2 . La distribución de A (o de S), a menudo llamada distribución de Wishart, es fundamental en estadística multivariante. Nos ocuparemos de ella y estudiaremos algunas de sus propiedades.

La varianza muestral generalizada se define como $|S|$ y es una especie de medida de la dispersión de la muestra. Su distribución también es considerada en su momento.

2. LA DISTRIBUCIÓN DE WISHART

En este párrafo estudiaremos la distribución de $A = \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})'$, donde las x_i son independientes, cada una con distribución $N(\mu, \Sigma)$. Como vimos en un capítulo anterior, A se distribuye como $A = \sum_{i=1}^n z_i z_i'$, donde $n = N-1$ y las z_i son independientes, cada una con distribución $N(0, \Sigma)$. Demostraremos que la distribución (o mejor) la densidad de A , para A definida positiva es

$$\frac{|A|^{\frac{1}{2}(n-p-1)} \cdot \exp \left\{ -\frac{1}{2} \text{tr} \Sigma^{-1} A \right\}}{2^{\frac{1}{2}np} \prod_{i=1}^p \Gamma(p-i+1/2) \cdot |\Sigma|^{\frac{1}{2}n}} \quad (1)$$

Obtenemos en primer lugar (1) para $\Sigma = I$. Utilizaremos aquí repetidamente el siguiente caso especial de un teorema enunciado cuando estudiáramos la distribución del coeficiente de correlación parcial. Si los x_i son independientes y u_i tiene la distribución $N(0, \Sigma)$, entonces $\sum_{i=1}^n u_i u_i' - \sum_{i=1}^n u_i u_i' w_i' (\sum_{i=1}^n w_i w_i')^{-1} \sum_{i=1}^n u_i u_i' w_i$ se distribuye como $\sum_{i=1}^{n-q} v_i v_i'$, donde q es el número de componentes de w_i y las v_i son independientes, cada una con una distribución $N(0, \Sigma)$, e independientemente de $\sum_{i=1}^n u_i u_i' w_i' (\sum_{i=1}^n w_i w_i')^{-1} \sum_{i=1}^n u_i u_i' w_i$. En particular si $q=1$, entonces $\sum_{i=1}^n u_i u_i' - \sum_{i=1}^n u_i u_i' w_i' (\sum_{i=1}^n w_i w_i')^{-1} \sum_{i=1}^n u_i u_i' w_i$ tiene una distribución χ^2 con $n-q$ grados de libertad. cuya densidad es

$$\frac{1}{2^{\frac{1}{2}(n-q)} \Gamma(\frac{1}{2}(n-q))} e^{-\frac{1}{2}t}$$

Tenemos posteriormente que $\sum_{i=1}^n u_i u_i' w_i' (\sum_{i=1}^n w_i w_i')^{-1} \sum_{i=1}^n u_i u_i' w_i$ es nula; si $\Sigma = 0$, entonces $E(\sum_{i=1}^n u_i u_i' w_i' (\sum_{i=1}^n w_i w_i')^{-1} \sum_{i=1}^n u_i u_i' w_i) = 0$ y la matriz de covarianzas es

$$E \left(\sum_{i=1}^n u_i u_i' w_i' (\sum_{i=1}^n w_i w_i')^{-1} \sum_{i=1}^n u_i u_i' w_i \right) = \sum_{\alpha \beta} w_{\alpha} w_{\beta}' \delta_{\alpha \beta} = \sum_{\alpha} w_{\alpha} w_{\alpha}'$$

Sea $a_{ij} = (a_{i+1,1}, a_{i+1,2}, \dots, a_{i+1,p})$ y $A_{ij} = (a_{j1}, a_{j2}, \dots, a_{jp})$.

Entonces

$$A_{ii} = \begin{pmatrix} a_{ii} & a_{i,i+1} \\ a_{i,i+1}' & A_{i+1,i+1} \end{pmatrix}$$

Sea $a_{ii, i+1, \dots, p} = a_{ii} - a_{ii} A_{i+1, i+1}^{-1} a_{ii}'$. El conjunto (z_{i1}, \dots, z_{in}) se distribuye independientemente de (z_{j1}, \dots, z_{jn}) , $j \neq i$ (a causa de $\Sigma = I$), y portanto, condicionados a z_{j1}, \dots, z_{jn} ($j \neq i, \alpha = 1, \dots, n$). Los elementos de (z_{i1}, \dots, z_{in}) se distribuyen independientemente, cada uno $N(0, 1)$, que es de la forma $N(0, \Sigma)$ con $\Sigma = 0$ y $\phi = 1$. Sea $z_{i1}^{(i+1)} = (z_{i1}, z_{i2}, \dots, z_{ip})$. Entonces $a_{ii} = \sum_{i=1}^n z_{i1} z_{i1}^{(i+1)'} + \sum_{i=1}^n z_{i1} z_{i1}^{(i+1)'} w_i' (\sum_{i=1}^n w_i w_i')^{-1} \sum_{i=1}^n z_{i1} z_{i1}^{(i+1)'} w_i$. Aplicando el caso especial del teorema citado encontramos que condicionado a $z_{j1}^{(i+1)} = z_{j1}^{(i+1)}$, $a_{ii, i+1, \dots, p}$ tiene una distribución χ^2 con $n-(p-1)$ grados de libertad e independiente de a_{ii} que se distribuye condicionadamente como $N(0, A_{i+1, i+1})$.

Se observará que la distribución condicional depende de $z_{j1}^{(i+1)}$ sólo a través de $A_{i+1, i+1}$; y definir, la densidad de $a_{ii, i+1, \dots, p}$,

$a_{ii}, a_{ii, i+1, \dots, p}, a_{ii, i+1, \dots, p}, a_{ii, i+1, \dots, p}, a_{ii, i+1, \dots, p}, a_{ii, i+1, \dots, p}$

$$\begin{aligned} f_1(a_{ii, i+1, \dots, p}, a_{ii, i+1, \dots, p}) \dots f_{p-1}(a_{ii, i+1, \dots, p}, a_{ii, i+1, \dots, p}) \cdot f(a_{ii, i+1, \dots, p}) \cdot f(a_{ii, i+1, \dots, p}) \\ = \frac{a_{ii, i+1, \dots, p}^{\frac{1}{2}(n-p)} \cdot e^{-\frac{1}{2}a_{ii, i+1, \dots, p}}}{2^{\frac{1}{2}n} \Gamma(\frac{1}{2}n)} \cdot \prod_{i=1}^{p-1} \left\{ \frac{a_{ii, i+1, \dots, p}^{\frac{1}{2}(n-p-i)} \cdot e^{-\frac{1}{2}a_{ii, i+1, \dots, p}}}{2^{\frac{1}{2}(n-p-i)} \Gamma(\frac{1}{2}(n-p-i))} \cdot \frac{e^{-\frac{1}{2}a_{ii, i+1, \dots, p}}}{(2\pi)^{\frac{1}{2}(p-1)} |A_{i+1, i+1}|^{\frac{1}{2}}} \right\} \end{aligned}$$

Encontramos la densidad de $a_{ij}, a_{i1}, \dots, a_{ip}$ substituyendo en la expresión anterior $a_{ii, i+1, \dots, p} = a_{ii} - a_{(i)} A_{i+1, i+1}^{-1} a_{(i)}$, y multiplicando por el Jacobiano, que es 1, y multiplicando posteriormente la expresión. El exponente de e en la anterior igualdad es

$$-\frac{1}{2} \left[a_{pp} + \sum_{i=1}^{p-1} a_{ii, i+1, \dots, p} + \sum_{i=1}^{p-1} a_{(i)} A_{i+1, i+1}^{-1} a_{(i)} \right] = -\frac{1}{2} \left[a_{pp} + \sum_{i=1}^{p-1} (a_{ii} - a_{(i)} A_{i+1, i+1}^{-1} a_{(i)}) + \sum_{i=1}^{p-1} a_{(i)} A_{i+1, i+1}^{-1} a_{(i)} \right] =$$

$$= -\frac{1}{2} \sum_{i=1}^p a_{ii} = -\frac{1}{2} \text{tr} A.$$

Utilizando una propiedad de las matrices, que asegura que

$$a_{ii, i+1, \dots, p} = \frac{\begin{vmatrix} a_{ii} & a_{(i)} \\ a_{(i)} & A_{i+1, i+1} \end{vmatrix}}{\begin{vmatrix} A_{i+1, i+1} \end{vmatrix}} = \frac{|A_{ii}|}{|A_{i+1, i+1}|}$$

encontramos que ($a_{pp} = A_{pp}$)

$$a_{pp} \prod_{i=1}^{p-1} a_{ii, i+1, \dots, p} = a_{pp} \prod_{i=1}^{p-1} \frac{|A_{ii}|}{|A_{i+1, i+1}|} = |A_{11}| = |A|.$$

Entonces

$$a_{pp} \prod_{i=1}^{p-1} \frac{a_{ii, i+1, \dots, p}}{|A_{i+1, i+1}|^{\frac{1}{2}}} = |A|^{\frac{1}{2}} \cdot a_{pp}^{\frac{1}{2}} \cdot \prod_{i=1}^{p-1} \frac{|A_{ii}|^{\frac{1}{2}(i-1)}}{|A_{i+1, i+1}|^{\frac{1}{2}i}} = |A|^{\frac{1}{2}(n-p-1)}.$$

La potencia de π en el denominador es $\frac{1}{2} [(p-1) + (p-2) + \dots + 1] = \frac{1}{2} [p(p-1)/2]$. Puesto que $\pi \left(\left(\frac{1}{2} \pi \right)^{\prod_{i=1}^{p-1} \left(\frac{1}{2} [n-(p-i)] \right)} \right) =$
 $= \prod_{i=1}^p \pi \left[\frac{1}{2} (n-i+1) \right]$, encontramos que la densidad de $a_{11}, a_{22}, \dots, a_{pp}$ es

$$\frac{|A|^{\frac{1}{2}(n-p-1)} e^{-\frac{1}{2} \text{tr} A}}{2^{\frac{1}{2}np} \pi^{p(p-1)/4} \prod_{i=1}^p \pi \left[\frac{1}{2} (n+1-i) \right]}$$

que es la expresión de (1) cuando $\Sigma = I$.

Obtenemos la distribución para Σ arbitrario. Sea $A = \Sigma_{\alpha=1}^n Z_{\alpha} Z_{\alpha}'$, donde las Z_{α} son independientes, cada una con distribución $N(0, \Sigma)$, y sea $A^* = \Sigma_{\alpha=1}^n Z_{\alpha}^* Z_{\alpha}^{*'}'$, donde las Z_{α}^* son independientes, cada una con distribución $N(0, I)$. Entonces la densidad de A^* es la obtenida. Sea C una matriz triangular arbitraria ($c_{ij} = 0, i > j$) tal que $C[C'] = I$ (siempre existe C por una propiedad de las matrices definidas positivas). La distribución de $CAC' = \Sigma_{\alpha=1}^n (C Z_{\alpha}) (C Z_{\alpha})'$ es la misma que la de A^* puesto que la distribución de $C Z_{\alpha}$ es la de Z_{α}^* . Por tanto obtenemos la distribución de A substituyendo $A^* = CAC'$ en la densidad obtenida y multiplicando por el Jacobiano.

LEMA. - Sea A^* una matriz simétrica que transformamos en otra matriz simétrica A mediante $A^* = CAC'$, donde C es una matriz triangular no singular. El jacobiano de la transformación es $\text{mod } |C|^{p+1}$ (mod. indica módulo).

Demostración. - la transformación es

$$a_{ij}^* = \sum_{k \in l} c_{ik} a_{kk} c_{jl}.$$

las derivadas parciales son

$$\frac{\partial a_{ij}^*}{\partial a_{kk}} = c_{ik} c_{jk}$$

$$\frac{\partial a_{ij}^*}{\partial a_{kk}} = c_{ik} c_{jl} + c_{il} c_{jk} \quad l \neq k$$

Obtendremos la matriz de las derivadas parciales en el orden $a_{11}, a_{22}, \dots, a_{pp}, a_{12}, \dots, a_{pp}$; la posición de la fila corresponde a a_{ij}^* y la posición de la columna a a_{ij} . la matriz es

$$\begin{bmatrix} c_{11}^2 & 2c_{11}c_{12} & \dots & 2c_{11}c_{1p} & c_{12}^2 & \dots & c_{1p}^2 \\ 0 & c_{11}c_{22} & \dots & c_{11}c_{2p} & c_{12}c_{22} & \dots & c_{1p}c_{2p} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & c_{11}c_{pp} & c_{1p}c_{2p} & \dots & c_{1p}c_{pp} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & c_{22}^2 & \dots & c_{2p}^2 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & 0 & \dots & c_{pp}^2 \end{bmatrix}$$

que es una matriz triangular. El determinante es el producto de los elementos de la diagonal, que vale $\prod_{j=1}^p c_{jj}^{p+1} = |C|^{p+1}$. lo que demuestra el lema.

Si en la densidad de A^* reemplazamos A por AC' , de $CC' = I$ tenemos $1 = |CC'| = |C| \cdot |C'| = |C| \cdot |C|^{-1} = |C| \cdot |C|^{-1}$; así $|C'| = 1/|C|$ y $\text{mod } |C| = 1/\sqrt{|C|}$. Además $\Sigma = C'(C')^{-1} = (C'C)^{-1}$. Así $\text{tr } AC' = \text{tr } AC'C = \text{tr } A\Sigma^{-1}$, $|AC'| = |C| \cdot |A \cdot C'| = |A| \cdot |CC'| = |A| \cdot |C|^{-1}$. Estas substituciones y el Jacobiano ante obtenido dan por resultado la expresión (1) en que hemos iniciado el párrafo.

TEOREMA. - Supongamos los vectores p -dimensionales Z_1, \dots, Z_n ($n \geq p$), que son independientes, cada uno distribuido $N(0, \Sigma)$. Entonces la densidad de $A = \sum_{\alpha=1}^n Z_\alpha Z_\alpha'$ es

$$\frac{|A|^{-\frac{1}{2}(n-p-1)} e^{-\frac{1}{2} \text{tr } A \Sigma^{-1}}}{2^{\frac{1}{2}np} \pi^{p(p-1)/4} |\Sigma|^{\frac{1}{2}n} \prod_{i=1}^p \Gamma[\frac{1}{2}(n+1-i)]}$$

para A definida positiva, y 0 en el caso.

COROLARIO. - Supongamos los vectores p -dimensionales Z_1, \dots, Z_N ($N \geq p$), independientes, cada uno distribuido $N(\mu, \Sigma)$. Entonces la densidad de $A = \sum_{\alpha=1}^N (Z_\alpha - \bar{Z})(Z_\alpha - \bar{Z})'$ es la anterior con $n = N-1$.

Es costumbre denotar la densidad obtenida mediante $W(A/\Sigma, n)$ y la correspondiente distribución como $W(A/\Sigma, n)$ o bien $W(\Sigma, n)$. La primera obtención de esta distribución (Wishart, 1928) fue mediante un argumento geométrico, que está muy relacionado con la demostración que hemos presentado aquí. Sea $V_i' = (Z_{i1}, \dots, Z_{ip})$ un vector en el espacio R^p . Los elementos diagonales de A son los cuadrados de las longitudes de estos vectores, $a_{ii} = V_i' V_i$, y los elementos de fuera de la diagonal están relacionados con los ángulos entre los vectores puesto que $r_{ij} = a_{ij}/\sqrt{a_{ii}a_{jj}}$ es el coseno del ángulo entre V_i y V_j . La matriz A describe las longitudes y relaciones de estos vectores.

El elemento de probabilidad de $\sqrt{a_{11}}V_1, V_1'V_2, \dots, V_1'V_p$ dados V_2, \dots, V_p es aproximadamente la probabilidad de que V_1 caiga en la región para la que $\sqrt{a_{11}} < V_1'V_p < d\sqrt{a_{11}} + d\sqrt{a_{11}}$. El primer par de desigualdades define la región entre dos hiperplanos. En esta región la densidad $(2\pi)^{-\frac{1}{2}n} \exp(-\frac{1}{2} V_1'V_1)$ es aproximadamente constante. La intersección de las regiones es un casquete esférico en $n-(p-1)$ dimensiones con una sección transversal de $p-1$ dimensiones. El volumen de esta región es aproximadamente $d\sqrt{a_{11}} da_{11} da_{11} \dots da_{1p} / |A_{11,11}|^{1/2}$. El radio del casquete esférico es $a_{11} \cos \alpha_1, \dots, \alpha_p = a_{11} - a_{11} A_{11,11}^{-1} a_{11}$. El área (o volumen) de esta superficie es la potencia $[n-(p-1)-1]$ -ésima del radio por el área de la superficie de la esfera (en $n-(p-1)$ dimensiones) de radio unidad. El área de la superficie de la esfera unidad es $G(n-(p-1)) = 2\pi^{\frac{1}{2}[n-(p-1)]} / \Gamma[\frac{1}{2}[n-(p-1)]]$. Así la probabilidad de $\sqrt{a_{11}}V_1, V_1'V_2, \dots, V_1'V_p$ es

$$\frac{e^{-\frac{1}{2} V_1'V_1}}{(2\pi)^{\frac{1}{2}n}} \cdot \frac{2\pi^{\frac{1}{2}[n-(p-1)]}}{\Gamma[\frac{1}{2}[n-(p-1)]]} \cdot \frac{d\sqrt{a_{11}} da_{11} da_{11} \dots da_{1p}}{|A_{11,11}|^{1/2}}$$

El elemento de probabilidad de $V_1'V_2, V_1'V_3, \dots, V_1'V_p$ implica la substitución de $d\sqrt{a_{11}} = da_{11}/(2\sqrt{a_{11}})$. Esto conduce al mismo término del producto ante obtenido y que representaba la densidad de $a_{11}, a_{12}, \dots, a_{1p}, a_{22}, \dots, a_{2p}, a_{33}, \dots, a_{3p}, \dots, a_{p-1,p-1}, a_{p-1,p}, \dots, a_{pp}$.

Existen otros muchos métodos para llegar a este resultado (ver Anderson pag. 153).

Para terminar con este párrafo daremos la distribución conjunta de los variaciones y covarianzas muestrales. Hemos demostrado que $N\Sigma^{-1}$, para muestras de tamaño N que provienen de una $N(\mu, \Sigma)$ distribuye como $\sum_{\alpha=1}^N Z_\alpha Z_\alpha'$ donde las Z_α son independientes, cada una con

distribución $N(0, \Sigma)$ y $n = N - 1$, $E(N\hat{\Sigma}) = n\Sigma$, y $E(S) = \Sigma$, donde $S = (N/n)\hat{\Sigma}$.

TEOREMA. Supongamos que $Z_1, \dots, Z_N (N \geq p+1)$ se distribuyen independientemente, cada una como $N(\mu, \Sigma)$. Entonces la distribución de $S = (1/n) \sum_{\alpha=1}^n (Z_\alpha - \bar{Z})(Z_\alpha - \bar{Z})'$ es $W[(1/n)\Sigma, n]$, donde $n = N - 1$ y $W[(1/n)\Sigma, n]$ es la distribución de Wishart con matriz de varianzas $(1/n)\Sigma$ y n grados de libertad.

Demonstración. Claramente, $S = (1/n)A = \sum_{\alpha=1}^n [(1/n)Z_\alpha][(1/n)Z_\alpha]'$, donde los $(1/n)Z_\alpha$ son independientes, con distribución $N(0, (1/n)\Sigma)$. Basta aplicar el teorema anterior a estos mismos $Z_\alpha^* = (1/n)Z_\alpha$.

Si $n < p$, la matriz $A = \sum_{\alpha=1}^n Z_\alpha Z_\alpha'$ no tiene una densidad de probabilidad. No obstante, nos referiremos a la distribución correspondiente, como a la distribución de Wishart.

3. ALGUNAS PROPIEDADES DE LA DISTRIBUCIÓN DE WISHART

La Función Característica

La función característica de la distribución de Wishart puede obtenerse fácilmente a partir de la distribución de las observaciones. Supongamos que Z_1, \dots, Z_n se distribuyen independientemente, cada una en densidad

$$\frac{1}{(2\pi)^{\frac{1}{2}p}} \exp\left(-\frac{1}{2}z' \Sigma^{-1} z\right).$$

Sea

$$A = \sum_{\alpha=1}^n Z_\alpha Z_\alpha'.$$

Introducimos la matriz $\Theta = (\theta_{ij})$ con $\theta_{ij} = \theta_{ji}$. La función característica de $A_{11}, A_{22}, \dots, A_{pp}, 2A_{12}, 2A_{13}, \dots, 2A_{p-1,p}$ es

$$E \left[\exp(i \operatorname{tr}(A\Theta)) \right] = E \left[\exp(i \operatorname{tr} \sum_{\alpha=1}^n Z_\alpha Z_\alpha' \Theta) \right] = E \left[\exp(i \operatorname{tr} \sum_{\alpha=1}^n Z_\alpha' \Theta Z_\alpha) \right] = E \left[\exp(i \sum_{\alpha=1}^n Z_\alpha' \Theta Z_\alpha) \right]$$

en virtud del hecho de que $\operatorname{tr} EFG = \sum_{i,j} f_{ji} g_{ji} = \operatorname{tr} FGE$. Como los Z_α son independientes, tendremos

$$E \left[\exp(i \sum_{\alpha=1}^n Z_\alpha' \Theta Z_\alpha) \right] = \prod_{\alpha=1}^n E \left[\exp(i Z_\alpha' \Theta Z_\alpha) \right] = \left\{ E \left[\exp(i Z' \Theta Z) \right] \right\}^n.$$

Para Θ real existe una matriz B , real no singular tal que

$$B' \Sigma^{-1} B = I$$

$$B' \Theta B = D$$

donde D es una matriz diagonal. Si hacemos

$$Z = BY$$

entonces

$$E \left[\exp(i Z' \Theta Z) \right] = E \left[\exp(i Y' D Y) \right] = E \left[\prod_{j=1}^p \exp(i d_{jj} Y_j^2) \right] = \prod_{j=1}^p E \left[\exp(i d_{jj} Y_j^2) \right].$$

El j -ésimo término en el producto es la $E[\exp(i d_{jj} Y_j^2)]$ donde Y_j tiene una distribución $N(0, 1)$; es pues la función característica de la distribución χ^2 con un grado de libertad, es decir $(1 - 2i d_{jj})^{-1/2}$. Así

$$E \left[\exp(i Z' \Theta Z) \right] = \prod_{j=1}^p (1 - 2i d_{jj})^{-1/2} = |I - 2iD|^{-1/2}$$

puesto que $I - 2iD$ es una matriz diagonal. Remitiendo en cuenta lo que vale D y lo que vale I , tendremos

$$|I - 2iD| = |B' \Sigma^{-1} B - 2iB' \Theta B| = |B' (\Sigma^{-1} - 2i\Theta) B| = |B'| \cdot |\Sigma^{-1} - 2i\Theta| \cdot |B| = |B|^2 \cdot |\Sigma^{-1} - 2i\Theta|$$

$|B'| \cdot |\Sigma^{-1}| \cdot |B| = |I| = 1$, y $|B|^2 = 1/|\Sigma^{-1}|$. Combinando estos resultados, obtenemos

$$E \left[\exp(i \operatorname{tr}(A\Theta)) \right] = \frac{|\Sigma^{-1}|^{\frac{1}{2}n}}{|\Sigma^{-1} - 2i\Theta|^{\frac{1}{2}n}}.$$

Puede demostrarse que ste sustituido el símbolo con el que re $(\sigma_{jk}^2 - 2i\theta_{jk})$ sea definida positiva. En particular es cierto para toda θ real.

TEOREMA.- Si Z_1, \dots, Z_n son independientes, cada una con distribución $N(0, \Sigma)$, la función característica $A_{11}, \dots, A_{pp}, 2A_{12}, \dots, 2A_{p-1,p}$,

$$\text{donde } (A_{ij}) = A = \sum_{\alpha=1}^n Z_{\alpha} Z_{\alpha}^T, \text{ es}$$

$$f(\theta) = \frac{|\Sigma^{-1}|^{\frac{1}{2}n}}{|\Sigma^{-1} - 2i\theta|^{\frac{1}{2}n}}$$

Podemos obtener los momentos de A , bien a partir de la función característica, bien a partir de la distribución original. El valor esperado de A_{ij} es

$$E(A_{ij}) = E\left[\sum_{\alpha=1}^n Z_{i\alpha} Z_{j\alpha}\right] = \sum_{\alpha=1}^n \sigma_{ij} = n\sigma_{ij}.$$

Para obtener las variancias necesitamos

$$\begin{aligned} E(A_{ij} A_{kl}) &= E\left[\sum_{\alpha, \beta=1}^n Z_{i\alpha} Z_{j\alpha} Z_{k\beta} Z_{l\beta}\right] = E\left[\sum_{\alpha=1}^n Z_{i\alpha} Z_{j\alpha} Z_{k\alpha} Z_{l\alpha}\right] + E\left[\sum_{\substack{\alpha, \beta=1 \\ \alpha \neq \beta}}^n Z_{i\alpha} Z_{j\alpha} Z_{k\beta} Z_{l\beta}\right] = \\ &= n(\sigma_{ij}\sigma_{kl} + \sigma_{ik}\sigma_{jl} + \sigma_{il}\sigma_{jk}) + n(n-1)\sigma_{ij}\sigma_{kl} = n^2\sigma_{ij}\sigma_{kl} + n\sigma_{ik}\sigma_{jl} + n\sigma_{il}\sigma_{jk}, \end{aligned}$$

obteniendo los momentos de cuarto orden los definidos en un momento cuando estudiamos la normal multivariante. Así, la variancia entre A_{ij} y A_{kl} es

$$E[(A_{ij} - n\sigma_{ij})(A_{kl} - n\sigma_{kl})] = n\sigma_{ik}\sigma_{jl} + n\sigma_{il}\sigma_{jk}.$$

Para $i=k$ y $j=l$ obtenemos la variancia de A_{ij}

$$E[(A_{ij} - n\sigma_{ij})^2] = n(\sigma_{ij}^2 + \sigma_{ii}\sigma_{jj}).$$

La suma de matrices de Wishart

Supongamos que A_i ($i=1, 2$) se distribuyen independientemente como $W(\Sigma, n_i)$ respectivamente. Entonces A_1 se distribuye como $\sum_{\alpha=1}^{n_1} Z_{\alpha} Z_{\alpha}^T$ y A_2 como $\sum_{\alpha=n_1+1}^{n_1+n_2} Z_{\alpha} Z_{\alpha}^T$, donde las Z_{α} son independientes, cada una $N(0, \Sigma)$. Entonces

$$A = A_1 + A_2$$

se distribuye como $\sum_{\alpha=1}^n Z_{\alpha} Z_{\alpha}^T$ con $n = n_1 + n_2$. Es decir, A se distribuye $W(\Sigma, n)$. Obviamente la suma de q matrices, distribuidas independientemente, cada una con distribución Wishart con covarianza Σ , tiene una distribución Wishart con matriz Σ y número de grados de libertad igual a la suma de los grados de libertad de las matrices componentes.

Una cierta transformación lineal

Haremos frecuentemente la transformación

$$A = CBC^T$$

donde C es una matriz no singular $p \times p$. Si A se distribuye como $W(\Sigma, n)$, entonces B se distribuye como $W(\Phi, n)$ donde

$$\Phi = C^{-1}\Sigma C^{-T}.$$

Esto se prueba mediante el siguiente argumento:

$$A = \sum_{\alpha=1}^n Z_{\alpha} Z_{\alpha}^T,$$

donde Z_{α} se distribuyen independientemente, cada una como $N(0, \Sigma)$. Entonces

$$Y_{\alpha} = C^{-1}Z_{\alpha}$$

se distribuye como $N(0, \Phi)$. Por tanto,

$$B = \sum_{\alpha=1}^n y_{\alpha} y_{\alpha}' = C^{-1} \sum_{\alpha=1}^n z_{\alpha} z_{\alpha}' C^{-1} = C^{-1} A C^{-1}$$

re distribuye como $W(\Phi, n)$. El Jacobiano de la transformación, $\left| \frac{\partial(A)}{\partial(B)} \right|$, viene dado por

$$\left| \frac{\partial(A)}{\partial(B)} \right| = \frac{w(B, \Phi, n)}{w(A, \Sigma, n)} = \frac{|B|^{\frac{1}{2}(n-p-1)} |\Sigma|^{\frac{1}{2}n}}{|A|^{\frac{1}{2}(n-p-1)} |\Phi|^{\frac{1}{2}n}} = \text{mod } |C|^{p+1}.$$

Distribuciones marginales

Si A re distribuye $W(\Sigma, n)$, la distribución marginal de cualquier subconjunto de elementos de A puede ser difícil de obtener. No obstante, la distribución marginal de algunos conjuntos de elementos puede obtenerse con facilidad. Vamos a verlo con los siguientes teoremas

TEOREMA.- Sean A y Σ fraccionados en q y $p-q$ filas y columnas

$$A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

Si A re distribuye $W(\Sigma, n)$, entonces A_{11} lo hace $W(\Sigma_{11}, n)$.

Demostración.- A re distribuye como $\sum_{\alpha=1}^n z_{\alpha} z_{\alpha}'$, donde los z_{α} son independientes, cada uno $N(0, \Sigma)$. Fraccionamos z_{α} en dos subvectores de q y $p-q$ componentes respectivamente,

$$z_{\alpha} = \begin{pmatrix} z_{\alpha}^{(1)} \\ z_{\alpha}^{(2)} \end{pmatrix}.$$

Entonces $z_{\alpha}^{(1)}$ son independientes, cada uno en distribución $N(0, \Sigma_{11})$ y A_{11} re distribuye como $\sum_{\alpha=1}^n z_{\alpha}^{(1)} z_{\alpha}^{(1)'} =$ que tiene una distribución $W(\Sigma_{11}, n)$.

TEOREMA.- Sean A y Σ fraccionados en p_1, \dots, p_q filas y columnas ($p_1 + p_2 + \dots + p_q = p$)

$$A = \begin{pmatrix} A_{11} & \dots & A_{1q} \\ \vdots & & \vdots \\ A_{q1} & \dots & A_{qq} \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \Sigma_{11} & \dots & \Sigma_{1q} \\ \vdots & & \vdots \\ \Sigma_{q1} & \dots & \Sigma_{qq} \end{pmatrix}$$

Si $\Sigma_{ij} = 0$ para $i \neq j$ y si A re distribuye como $W(\Sigma, n)$, entonces $A_{11}, A_{22}, \dots, A_{qq}$ son independientes y re distribuyen $W(\Sigma_{jj}, n)$.

Demostración.- Fraccionamos z_{α} de la forma

$$z_{\alpha} = \begin{pmatrix} z_{\alpha}^{(1)} \\ \vdots \\ z_{\alpha}^{(q)} \end{pmatrix}$$

Puesto que $\Sigma_{ij} = 0$, $z_{\alpha}^{(i)}$ y $z_{\alpha}^{(j)}$ son independientes, $\forall \alpha, \beta$. Entonces $A_{ii} = \sum_{\alpha=1}^n z_{\alpha}^{(i)} z_{\alpha}^{(i)'}$ es independiente de $A_{jj} = \sum_{\alpha=1}^n z_{\alpha}^{(j)} z_{\alpha}^{(j)'}$. El uso del teorema re deriva del anterior.

4. EL TEOREMA DE COCHRAN

El teorema de Cochran es útil para probar que ciertas "formas cuadráticas vectoriales" re distribuyen como sumas de "cuadrados vectoriales". Es una afirmación estadística de un teorema algebraico. Demos en primer lugar la proposición algebraica concerniente a variables escalares.

LEMA.- Si

$$q_i = \sum_{\alpha, \beta=1}^N a_{\alpha\beta}^i y_{\alpha} y_{\beta}, \quad i=1, \dots, m$$

es de rango r_i y

$$\sum_{i=1}^m q_i = \sum_{\alpha=1}^N y_{\alpha}^2$$

Entonces una condición necesaria y suficiente para que exista una transformación ortogonal de $\{Y_\alpha\}$ a $\{Z_\alpha\}$ tal que

6(4)

$$q_i = \sum_{\alpha=r_{i-1}+1}^{r_i} z_\alpha^2$$

es que

$$r_1 + \dots + r_m = N$$

Demstración: - La necesidad de la condición es obvia porque la suma de los rangos no puede ser menor que N si se verifica que

$$\sum_{i=1}^m q_i = \sum_{\alpha=1}^N Y_\alpha^2$$

y no puede ser mayor que N en la transformación que pasa a $z_1, \dots, z_{r_1+r_2+\dots+r_m}$ ha de ser no singular.

Probermos pues la suficiencia. De una propiedad matricial sabemos que existe una matriz no singular D tal que

$$D' \begin{pmatrix} I & 0 & 0 \\ 0 & -I & 0 \\ 0 & 0 & 0 \end{pmatrix} D = A_2$$

donde $A_1 = (a_{\alpha\beta}^{(1)})$ y la suma de los ordenes de I y $-I$ es r_1 . Sea $a_{\alpha\beta} = b_{\alpha\beta}^{(1)}$, $\alpha = 1, \dots, r_1$, $\beta = 1, \dots, N$. Entonces

$Z_\alpha = \sum_{\beta=1}^N b_{\alpha\beta}^{(1)} Y_\beta$ ($\alpha = 1, \dots, r_1$) forman un conjunto de r_1 funciones lineales de $\{Y_\beta\}$ tal que

$$q_1 = \sum_{\alpha=1}^{r_1} c_\alpha z_\alpha^2$$

donde $c_\alpha = 1$ o -1 . En general existe un conjunto de r_i funciones lineales

$$Z_\alpha = \sum_{\beta=1}^N b_{\alpha\beta}^{(i)} Y_\beta, \quad \alpha = r_1 + \dots + r_{i-1} + 1, \dots, r_1 + \dots + r_i.$$

tal que

$$q_i = \sum_{\alpha=r_{i-1}+1}^{r_i+r_{i-1}} c_\alpha z_\alpha^2.$$

Así

$$\sum_{i=1}^m q_i = \sum_{\alpha=1}^N c_\alpha z_\alpha^2,$$

pero como la $\sum_{i=1}^m q_i$ es definida positiva los c_α han de ser todos positivos, $c_\alpha = 1$, $\forall \alpha$. Por tanto se verifica la igualdad

$$q_i = \sum_{\alpha=r_{i-1}+1}^{r_i+r_{i-1}} z_\alpha^2.$$

Además, substituyendo los Z_α por todos los i , tenemos

$$Z_\alpha = \sum_{\beta=1}^N b_{\alpha\beta} Y_\beta, \quad \alpha = 1, \dots, N$$

pero

$$\sum_{i=1}^m q_i = \sum_{\alpha=1}^N z_\alpha^2 = \sum_{\alpha=1}^N Y_\alpha^2,$$

lo que asegura que la transformación que pasa de los $\{Y_\alpha\}$ a los $\{Z_\alpha\}$ es ortogonal (a saber, $I = B'IB = B'B$).

Una expresión alternativa del lema es la siguiente: Sea r_i el rango de la matriz de orden N , A_i , matriz que es además, simétrica ($i = 1, \dots, m$), y supongamos $\sum_{i=1}^m A_i = I$. Una condición necesaria y suficiente para que exista una matriz ortogonal

$$B = \begin{pmatrix} B_1 \\ \vdots \\ B_m \end{pmatrix}$$

tal que $A_i = B_i' B_i$ y que $\sum_{i=1}^m r_i = N$.

Estableceremos a continuación el teorema de Cochran.

TEOREMA DE COCHRAN - Supongamos que Y_α se distribuye $N(0, \Sigma)$ independientemente de Y_β ($\alpha \neq \beta$). Supongamos que la matriz $A_i = (a_{\alpha\beta}^i)$ utilizada para formar

$$Q_i = \sum_{\alpha, \beta=1}^N a_{\alpha\beta}^i Y_\alpha Y_\beta' \quad i=1, \dots, m$$

es de rango r_i y supongamos

$$\sum_{i=1}^m Q_i = \sum_{\alpha=1}^N Y_\alpha Y_\alpha'$$

Entonces una condición necesaria y suficiente para que Q_i ($i=1, \dots, m$) se distribuya como

$$\sum_{\alpha=r_{i-1}+1}^{r_i+r_{i-1}+1} Z_\alpha Z_\alpha'$$

donde las Z_α son $N(0, \Sigma)$ e independientes y Q_i se distribuya independientemente de Q_j ($i \neq j$) y que

$$r_1 + \dots + r_m = N.$$

COROLARIO - Si $r_i \geq p$, entonces Q_i se distribuye como $W(\Sigma, r_i)$

Demostración - Si $r_1 + r_2 + \dots + r_m = N$ se verifica, existe una matriz ortogonal B tal que $A_i = B_i' B_i$. Puesto que B es ortogonal

$$Z_\alpha = \sum_{\beta=1}^N b_{\alpha\beta} Y_\beta$$

son independientes, cada uno con distribución $N(0, \Sigma)$, por una propiedad ya estudiada de este tipo de transformaciones, vemos que

$$Q_i = \sum_{\alpha, \beta=1}^N a_{\alpha\beta}^i Y_\alpha Y_\beta' = \sum_{\alpha, \beta=1}^N b_{\alpha\beta}^i b_{\alpha\beta}^i Y_\alpha Y_\beta' = \sum_{\gamma=r_{i-1}+1}^{r_i+r_{i-1}+1} Z_\gamma Z_\gamma'$$

La necesidad se deriva mediante una argumentación análoga a la del lema anterior.

Este teorema es útil para generalizar resultados de análisis de la varianza unidimensional, de ello nos ocuparemos más tarde. Como un ejemplo de aplicación del teorema, probaremos que ~~la suma de los cuadrados~~ ~~la~~ ~~media de los~~ ~~los~~ ~~matrices~~ ~~resultantes~~ al multiplicar una observación por su transpuesta, que es exactamente N veces el vector-media muestral por su transpuesta y un múltiplo de la matriz de covarianzas muestral se distribuyen independientemente con distribución Wishart singular y no singular respectivamente. Sea Y_1, \dots, Y_N vectores aleatorios independientes $N(0, \Sigma)$. Utilizaremos las matrices $(a_{\alpha\beta}^{(1)}) = (1/N)$ y $(a_{\alpha\beta}^{(2)}) = [\delta_{\alpha\beta} - (1/N)]$

Entonces

$$Q_1 = \sum_{\alpha, \beta=1}^N \frac{1}{N} Y_\alpha Y_\beta' = N \bar{Y} \bar{Y}'$$

$$Q_2 = \sum_{\alpha, \beta=1}^N (\delta_{\alpha\beta} - 1/N) Y_\alpha Y_\beta' = \sum_{\alpha=1}^N Y_\alpha Y_\alpha' - N(\bar{Y} \bar{Y}') = \sum_{\alpha=1}^N (Y_\alpha - \bar{Y})(Y_\alpha - \bar{Y})'$$

Obviamente

$$Q_1 + Q_2 = \sum_{\alpha=1}^N Y_\alpha Y_\alpha'$$

la primera matriz es de rango 1; la segunda matriz es de rango $(N-1)$ (puesto que el rango de la suma de dos matrices es menor o igual que la suma de los rangos y el rango de la segunda matriz es menor que N). Las condiciones del teorema se satisfacen; por tanto Q_1

se distribuye como χ^2 , donde $\chi^2 \sim N(0, \Sigma)$ y Q_2 se distribuye independientemente como $W(\Sigma, N-2)$.

(5)

5. LA VARIANZA GENERALIZADA

El concepto multivariante equivalente a la varianza, σ^2 , de una distribución univariante es la matriz de varianzas Σ . Otro concepto multivariante análogo es el scalar $|\Sigma|$, que se le llama varianza generalizada de la distribución multivariante. Análogamente, la varianza generalizada de la muestra de vectores x_1, \dots, x_N es

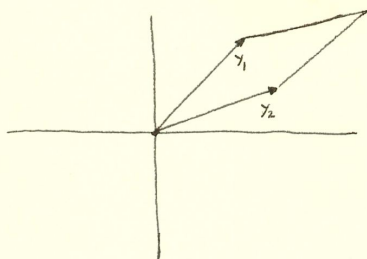
$$|S| = \left| \frac{1}{N-1} \sum_{\alpha=1}^N (x_\alpha - \bar{x})(x_\alpha - \bar{x})' \right|.$$

En algún sentido es una medida de extensión. Consideramos aquí este concepto porque recurrimos a él en muchos contextos de hipótesis basados en la razón de similitud.

Una interpretación geométrica de la varianza muestral generalizada puede darse en términos de p puntos en el espacio N -dimensional. Sea $Z_\alpha = x_\alpha - \bar{x}$, y hagamos

$$\begin{pmatrix} y_1' \\ \vdots \\ y_p' \end{pmatrix} = (Z_1 \dots Z_N).$$

Las N componentes del vector y_i son las i -ésimas componentes de Z_1, \dots, Z_N . Entonces $a_{ij} = y_i' y_j$ es el cuadrado de la longitud del i -ésimo vector y $a_{ij} = y_i' y_j$ es el producto del seno del ángulo entre y_i e y_j por las longitudes de y_i e y_j . Consideremos ahora



la figura basada en dos de estos vectores. Si $p=2$, tenemos un paralelogramo con y_1 e y_2 como ejes principales; si $p=3$ tenemos un paralelepípedo con y_1, y_2 e y_3 como ejes principales. Para cualquiera, la figura es un paralelepípedo en el hiperplano p -dimensional definido mediante y_1, \dots, y_p . Está delimitado por pares de hiperplanos $(p-1)$ -dimensionales, un hiperplano del par siendo definido mediante $p-1$ vectores y el otro pasando por el extremo del vector que falta. El determinante $|A| = |(N-1)S| = (N-1)^p |S|$ es el cuadrado del volumen de este paralelepípedo.

TEOREMA. Sea $Y = (y_1, \dots, y_p)$, donde los y_i son vectores N -dimensionales. Entonces el cuadrado del volumen p -dimensional del paralelepípedo con y_1, \dots, y_p como ejes principales es $|Y'Y| = |A| = (N-1)^p |S|$.

Demostración. El teorema es cierto si $p=1$, puesto que $|Y'Y| = y_1' y_1$, cuadrado de la longitud de y_1 . Supongámonos cierto para $p=k-1$ y probémoslo para $p=k$. Señalamos primero que si dos paralelepípedos k -dimensionales tienen bases consistentes en paralelepípedos $(k-1)$ -dimensionales de igual volumen e igual altura, sus volúmenes (k -dimensionales) serán iguales (por cuanto dicho volumen es la integral del volumen $(k-1)$ -dimensional). Puesto que es cierto para un paralelepípedo rectangular, el volumen k -dimensional es el producto de la altura por el volumen $(k-1)$ -dimensional de la base. Una combinación lineal de y_1, \dots, y_{k-1} , por ejemplo, $c_1 y_1 + \dots + c_{k-1} y_{k-1}$ es un vector en la base del paralelepípedo y la mínima longitud de $v = y_k - (c_1 y_1 + \dots + c_{k-1} y_{k-1})$ es la altura del paralelepípedo con y_1, \dots, y_k como ejes principales. La longitud de v se minimiza eligiendo c_1, \dots, c_{k-1} de manera que

$$0 = y_j' v = y_j' y_k - (c_1 y_j' y_1 + \dots + c_{k-1} y_j' y_{k-1}), \quad j=1, \dots, k-1. \quad \text{Sea } Y_k = (y_1, \dots, y_k) \text{ e } Y_{k-1} = (y_1, \dots, y_{k-1}),$$

$$G = \begin{pmatrix} 1 & 0 & \dots & 0 & -c_1 \\ 0 & 1 & \dots & 0 & -c_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 1 & -c_{k-1} \\ 0 & 0 & \dots & 0 & 1 \end{pmatrix}$$

Entonces $|G|=1$ y $(Y_{k-1} v) = Y_k G$. Tenemos

$$|Y_k' Y_k| = |G'| \cdot |Y_k' Y_k| \cdot |G| = |G' Y_k' Y_k G| = \left| \begin{pmatrix} y_{k-1}' \\ v' \end{pmatrix} (Y_{k-1} v) \right| = \begin{vmatrix} y_{k-1}' Y_{k-1} & 0 \\ 0 & v' v \end{vmatrix} = |Y_{k-1}' Y_{k-1}| v' v.$$

Puesto que hemos supuesto que $|Y_{k-1}' Y_{k-1}|$ es el cuadrado del volumen $(k-1)$ -dimensional de la base del paralelepípedo, $v' v$ es el cuadrado de la altura, el producto es el cuadrado del volumen del paralelepípedo k -dimensional. Hemos probado el teorema por inducción.

Vemos más tarde que muchos conceptos de estadística multivariante tienen una interpretación en términos de estos volúmenes. Estos volúmenes son análogos a las distancias que surgen en casos especiales cuando $p=1$.

Consideremos ahora una interpretación geométrica de $|A|$ en términos de N puntos en un espacio p -dimensional. Sean z_1, \dots, z_N definidos como antes, representen N puntos en el espacio p -dimensional. Cuando $p=1$, $|A| = \sum_{\alpha} z_{1\alpha}^2$, que es la suma de los cuadrados de las distancias de cada punto al origen. En general $|A|$ es la suma de los volúmenes de todos los paralelepípedos formados tomando como ejes principales puntos del conjunto z_1, \dots, z_N .

Veamos que

$$|A| = \begin{vmatrix} \sum_{\alpha} z_{1\alpha}^2 & \dots & \sum_{\alpha} z_{1\alpha} z_{p-1,\alpha} & \sum_{\alpha} z_{1\alpha} z_{pp} \\ \vdots & & \vdots & \vdots \\ \sum_{\alpha} z_{p-1,\alpha} z_{1\alpha} & \dots & \sum_{\alpha} z_{p-1,\alpha}^2 & \sum_{\alpha} z_{p-1,\alpha} z_{pp} \\ \sum_{\alpha} z_{p\alpha} z_{1\alpha} & \dots & \sum_{\alpha} z_{p\alpha} z_{p-1,\alpha} & \sum_{\alpha} z_{p\alpha}^2 \end{vmatrix} = \sum_{\beta} \begin{vmatrix} \sum_{\alpha} z_{1\alpha}^2 & \dots & \sum_{\alpha} z_{1\alpha} z_{p-1,\alpha} & z_{1\beta} z_{pp} \\ \vdots & & \vdots & \vdots \\ \sum_{\alpha} z_{p-1,\alpha} z_{1\alpha} & \dots & \sum_{\alpha} z_{p-1,\alpha}^2 & z_{p-1,\beta} z_{pp} \\ \sum_{\alpha} z_{p\alpha} z_{1\alpha} & \dots & \sum_{\alpha} z_{p\alpha} z_{p-1,\alpha} & z_{p\beta}^2 \end{vmatrix}$$

por la regla del desarrollo de un determinante. Aplicando este procedimiento repetidas veces a las columnas, tendremos

$$|A| = \sum_{\alpha_1, \dots, \alpha_p=1}^N |z_{\alpha_1}, \dots, z_{\alpha_p}|.$$

Por el teorema anterior el cuadrado del volumen del paralelepípedo con ejes principales $z_{\alpha_1}, \dots, z_{\alpha_p}$ es

$$V_{\alpha_1, \dots, \alpha_p}^2 = \left| \sum_{\beta} z_{\beta} z_{\beta} \right|$$

donde la suma sobre β se extiende a todos los $(\alpha_1, \dots, \alpha_p)$. Si desarrollamos este determinante de la misma forma que lo hicimos con $|A|$ obtenemos

$$V_{\alpha_1, \dots, \alpha_p}^2 = \left| z_{\beta_1}, \dots, z_{\beta_p} \right|$$

donde la suma es para cada β_j sobre el rango $(\alpha_1, \dots, \alpha_p)$. Sumando ahora sobre todos los posibles conjuntos $(\alpha_1, \dots, \alpha_p)$ obtenemos precisamente $|A|$. Así pues $|A|$ es la suma de los cuadrados de los volúmenes de todos los paralelepípedos formados por conjuntos de p vectores de los z_{α} como ejes principales. Si reemplazamos z_{α} mediante $x_{\alpha} - \bar{x}$, podemos enunciar el siguiente teorema:

TEOREMA. Sean x_1, \dots, x_N una muestra de tamaño N y sea $|S|$ el determinante de la matriz de varianzas muestrales. Entonces $|S|$ es proporcional a la suma de los cuadrados de los volúmenes de todos los paralelepípedos formados utilizando como ejes principales puntos de la forma $(x_i - \bar{x})$, elegidos los x_i entre los x_1, \dots, x_N , siendo el factor de proporcionalidad $1/(N-1)^p$.

El equivalente poblacional de $|S|$ es $|S|$, del que también puede darse una interpretación. Sabemos que

$$P \{ x' S^{-1} x \leq \chi_p^2(\alpha) \} = 1 - \alpha$$

si x se distribuye $N(0, S)$, es decir, la probabilidad de que x caiga en el interior del elipsoide $x' S^{-1} x = \chi_p^2(\alpha)$, es $1 - \alpha$. El volumen de este elipsoide viene dado por $C(p) |S|^{1/2} [\chi_p^2(\alpha)]^{1/2} / p$, con $C(p) = \frac{2\pi^{p/2}}{\Gamma(p/2)}$.

Distribución de la varianza muestral generalizada

La distribución de $|S|$ es la misma que la de $|A|/(N-1)^p$, donde

$$A = \sum_{\alpha=1}^n z_{\alpha} z_{\alpha}'$$

y z_{α} se distribuye, independientemente de z_{β} ($\alpha \neq \beta$), como $N(0, I)$, y $n = N-1$. Hagamos

$$W_{\alpha} = G z_{\alpha}$$

donde G ha sido elegida de manera que $G S G' = I$. Entonces W_{α} se distribuye, independientemente de W_{β} ($\alpha \neq \beta$), como $N(0, I)$ y

$$|B| = |A|/|S|$$

donde

$$B = \sum_{\alpha=1}^n W_{\alpha} W_{\alpha}' = \sum_{\alpha=1}^n G z_{\alpha} z_{\alpha}' G' = G A G'.$$

$$\text{y } |C| \cdot |S| \cdot |G'| = 1.$$

Como hicimos en un apartado anterior, sea $B_{ii} = (b_{jk})$, $j, k = i, \dots, p$, $b_{(i)} = (b_{i11}, b_{i12}, \dots, b_{ip})$.
 $b_{(i), i+1, \dots, p} = b_{ii} - b_{(i)} B_{i+1, i+1}^{-1} b'_{(i)}$. Entonces

$$|B| = b_{11, 2, \dots, p} b_{22, 3, \dots, p} \dots b_{pp}.$$

Como vimos en el apartado 2 de este capítulo, los $b_{11, 2, \dots, p}$, $b_{22, 3, \dots, p}$, \dots , b_{pp} son independientes, y $b_{ii, i+1, \dots, p}$ tiene una distribución χ^2 con $n - (p - i)$ grados de libertad. Así $|B|$ se distribuye como $\chi^2_n \cdot \chi^2_{n-1} \dots \chi^2_{n-p+1}$.

TEOREMA. - La distribución de la varianza generalizada $|S|$ de una muestra X_1, \dots, X_N de una $N(\mu, \Sigma)$ es la misma que la distribución de $|S|/(N-1)^p$ por el producto de p factores independientes, siendo la distribución del factor i -ésimo una χ^2 con $N-i$ grados de libertad.

Podemos dar nuevamente, una interpretación geométrica de este teorema. Sea $X_i = (W_{i1}, \dots, W_{in})$ un vector en un espacio n -dimensional. Los puntos X_1, \dots, X_p se distribuyen independientemente; cada componente de X_i se distribuye como $N(0, 1)$. $|B|$ es el cuadrado del volumen del paralelepípedo definido mediante X_1, \dots, X_p . Sea U_i el vector de X_i ortogonal a X_{i+1}, \dots, X_p . Entonces el cuadrado del volumen del paralelepípedo definido mediante X_{i+1}, \dots, X_p, X_i es el volumen al cuadrado del paralelepípedo definido mediante X_{i+1}, \dots, X_p por U_i , $U_i = a_{ii, i+1, \dots, p}$ (longitud al cuadrado de U_i). Se sigue que $|B| = U_1' U_1 \cdot U_2' U_2 \dots U_p' U_p$. U_i es ortogonal a U_{i+1}, \dots, U_p , por tanto la distribución condicional de U_i es normal en $(n - p + i)$ dimensiones y las $n - (p - i)$ coordenadas se distribuyen $N(0, 1)$. Así $U_i' U_i$ tiene una χ^2 con $n - (p - i)$ grados de libertad (independiente de las otras U_i).

Para $p = 1$ o 2 , la distribución exacta de $|S|$ puede obtenerse con relativa facilidad, pero para valores de p más elevados los cálculos involucrados en el proceso no son sencillos. No obstante, los momentos de $|S|$ pueden obtenerse fácilmente a partir de la igualdad

$$|S| = |A| / (N-1)^p$$

7 $|A|$ puede escribirse

$$|A| = |S| \cdot |B| = |S| \cdot \chi^2_{N-1} \cdot \chi^2_{N-2} \dots \chi^2_{N-p}.$$

Puesto que el k -ésimo momento de una χ^2 con m grados de libertad es $2^k \Gamma(\frac{1}{2}m + k) / \Gamma(\frac{1}{2}m)$ y el momento de un producto de variables aleatorias independientes es el producto de los momentos, el k -ésimo momento de $|A|$ es

$$|Z|^k \prod_{i=1}^p \left\{ 2^k \frac{\Gamma(\frac{1}{2}(N-i) + k)}{\Gamma(\frac{1}{2}(N-i))} \right\} = 2^{kp} |Z|^k \frac{\prod_{i=1}^p \Gamma(\frac{1}{2}(N-i) + k)}{\prod_{i=1}^p \Gamma(\frac{1}{2}(N-i))}.$$

Así

$$E(|A|) = |Z| \prod_{i=1}^p (N-i),$$

$$\text{var}(|A|) = |Z|^2 \prod_{i=1}^p (N-i) \left[\prod_{j=1}^p \Gamma(N-j+2) - \prod_{j=1}^p \Gamma(N-j) \right].$$

En el caso $p = 1$ y $p = 2$, podemos dar la distribución de $V = |A|/|S|$. Para $p = 1$, V se distribuye como una χ^2 con $N-1$ grados de libertad. Para $p = 2$, encontramos que el momento $k/2$ de V es

$$E[V^{k/2}] = 2^k \frac{\Gamma(\frac{1}{2}(N-1) + k/2) \Gamma(\frac{1}{2}(N-2) + k/2)}{\Gamma(\frac{1}{2}(N-1)) \cdot \Gamma(\frac{1}{2}(N-2))}$$

utilizando la fórmula de la duplicación para las funciones gamma,

$$2^{2\alpha-1} \Gamma(\alpha) \cdot \Gamma(\alpha + \frac{1}{2}) = \sqrt{\pi} \Gamma(2\alpha)$$

obtenemos como expresión para el momento $k/2$ de V

$$E[V^{k/2}] = \frac{\Gamma(N-2+k)}{\Gamma(N-2)}$$

que es el momento k -ésimo de $1/2$ de una variable con la distribución χ^2 con $2N-4$ grados de libertad. Revertimos el siguiente teorema:

TEOREMA DE CARLEMAN. - Si $\{\mu_i\}$ ($i=1,2,\dots$) es una sucesión de momentos tal que

$$\sum_{k \geq 1} \left(\frac{1}{\mu_{2k}} \right)^{1/2k}$$

es divergente, entonces al menos una distribución tiene la sucesión de momentos $\{\mu_i\}$.

La demostración puede encontrarse en Shohat y Tamarkin (1943) (The problem of the moments N.Y., A.M.S.). Utilizando el criterio de divergencia

$$\left(\frac{1}{\mu_{2k}} \right)^{1/2k} > \frac{1}{N-3+2k} > C \frac{1}{k}$$

para un $C > 0$ adecuado y para k suficientemente grande, podemos verificar que las condiciones del teorema anterior se satisfacen para los momentos de \sqrt{V} . Así $2\sqrt{V}$ se distribuye como una χ^2 con $2N-4$ grados de libertad.

Una aproximación a la densidad de $V^{1/2}$ viene dada por (Hoel, 1937)

$$\frac{e^{\frac{1}{2}p(N-p)} y^{\frac{1}{2}p(N-p)-1} e^{-cy}}{\Gamma\left(\frac{1}{2}p(N-p)\right)}$$

donde

$$c = \frac{p}{2} \left(1 - \frac{(p-1)(p-2)}{2N} \right)^{1/p}.$$

La distribución asintótica de la varianza generalizada

Sea $|B|/n^p = V_1(n), V_2(n), \dots, V_p(n)$, donde las V_i 's se distribuyen independientemente y $nV_i(n) = \chi_{n-p+i}^2$. Puesto que χ_{n-p+i}^2 se distribuye como $\sum_{\alpha=1}^{n-p+i} W_{\alpha}^2$, una vez las W_{α} independientes y $N(0,1)$, el teorema central del límite (aplicado a W_{α}^2) establece que

$$\frac{nV_i(n) - (n-p+i)}{\sqrt{2(n-p+i)}} = \sqrt{n} \frac{V_i(n) - 1 + \frac{p-i}{n}}{\sqrt{2} \sqrt{1 - \frac{p-i}{n}}}$$

se distribuye asintóticamente como $N(0,1)$. Entonces $\sqrt{n} [V_i(n) - 1]$ se distribuye asintóticamente como $N(0,2)$. Sea ahora

$$U(n) = \begin{pmatrix} V_1(n) \\ \vdots \\ V_p(n) \end{pmatrix}, \quad b = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}, \quad y \quad \sqrt{n} (U(n) - b)$$

que se distribuye asintóticamente como $N(0, T)$, siendo $T = 2I$. Definimos ahora $|B|/n^p = W = f(U_1, \dots, U_p) = U_1 U_2 \dots U_p$, siendo $\frac{\partial f}{\partial U_i} \Big|_{U=b} = 1$ y $\Phi_b' T \Phi_b = 2p$. Aplicando ahora un teorema límite ya utilizando cuando sea necesario la transformación de Fisher para coeficientes de correlación, obtenemos

$$\sqrt{n} \left(\frac{|B|}{n^p} - 1 \right)$$

se distribuye asintóticamente como $N(0, 2p)$.

TEOREMA. - Sea S una matriz de varianzas mutual de orden p con n grados de libertad. Entonces $\sqrt{n} (|S|/|S_1| - 1)$ es asintóticamente normal con media 0 y varianza $2p$.

POBLACIONAL DIAGONAL.

En su momento estudiaremos la distribución para un coeficiente de correlación simple cuando la correspondiente covarianza poblacional sea cero. Aquí vamos a encontrar la densidad para el conjunto r_{ij} $i \neq j$ $ij = 1 \dots p$ cuando $\rho_{ij} = 0$ $i \neq j$.

Comenzaremos estudiando la distribución de A cuando Σ es diagonal, es decir $W(\{\sigma_{ii}\delta_{ij}\}, n)$. La densidad de A es

$$\frac{|A_{ij}|^{\frac{1}{2}(n-p-1)} \exp\left(-\frac{1}{2} \sum_{i=1}^p \frac{a_{ii}}{\sigma_{ii}}\right)}{2^{\frac{1}{2}np} \pi^{\frac{1}{4}p(p-1)} \prod_{i=1}^p \sigma_{ii}^{\frac{1}{2}n} \prod_{i=1}^p \Gamma\left(\frac{1}{2}(n+1-i)\right)}$$

puesto que

$$|\Sigma| = \begin{vmatrix} \sigma_{11} & 0 & \dots & 0 \\ 0 & \sigma_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_{pp} \end{vmatrix} = \prod_{i=1}^p \sigma_{ii}$$

Hacemos la transformación

$$(1) \quad a_{ij} = \sqrt{a_{ii}} \sqrt{a_{jj}} r_{ij} \quad i \neq j$$

$$(2) \quad a_{ii} = a_{ii}$$

El Jacobiano de esta transformación viene dado por el producto del Jacobiano de (2) por el Jacobiano de (1) para a_{ii} fijos. El Jacobiano de (1) es el determinante de una matriz diagonal de orden $p(p-1)/2$ cuyos elementos en la diagonal principal son $\sqrt{a_{ii}}\sqrt{a_{jj}}$. Puesto que cada índice, por ejemplo k , aparece en el conjunto r_{ij} $i < j$, $p-1$ veces, el Jacobiano es

$$J = \prod_{i=1}^p a_{ii}^{\frac{1}{2}(p-1)}$$

Sustituyendo en la densidad de A obtendremos la densidad conjunta de $\{a_{ii}\}$ y $\{r_{ij}\}$

$$\frac{|\sqrt{a_{ii}} \sqrt{a_{jj}} r_{ij}|^{\frac{1}{2}(n-p-1)} \exp\left(-\frac{1}{2} \sum_{i=1}^p \frac{a_{ii}}{\sigma_{ii}}\right)}{2^{\frac{1}{2}np} \pi^{\frac{1}{4}p(p-1)} \prod_{i=1}^p \sigma_{ii}^{\frac{1}{2}n} \prod_{i=1}^p \Gamma\left(\frac{1}{2}(n+1-i)\right)} \prod_{i=1}^p a_{ii}^{\frac{1}{2}(p-1)} = \frac{|r_{ij}|^{\frac{1}{2}(n-p-1)}}{\prod_{i=1}^p \Gamma\left(\frac{1}{2}(n+1-i)\right) \prod_{i=1}^p \left\{ \frac{a_{ii}^{\frac{1}{2}n-1} \exp\left(-\frac{1}{2} \frac{a_{ii}}{\sigma_{ii}}\right)}{2^{\frac{1}{2}n} \sigma_{ii}^{\frac{1}{2}n}} \right\}}$$

puesto que

$$|\sqrt{a_{ii}} \sqrt{a_{jj}} r_{ij}| = \left(\prod_{i=1}^p a_{ii} \right) |r_{ij}|$$

donde $r_{ii} = 1$. En el último término del producto del miembro derecho de la igualdad obtenida después de la sustitución, hagamos

$a_{ii}/(\sigma_{ii}) = u_i$; entonces la integral de este término es

$$\int_0^\infty \frac{u_i^{\frac{1}{2}n-1} \exp\left(-\frac{1}{2} u_i\right)}{2^{\frac{1}{2}n} \sigma_{ii}^{\frac{1}{2}n}} da_{ii} = \int_0^\infty u_i^{\frac{1}{2}n-1} e^{-u_i} du_i = \Gamma\left(\frac{1}{2}n\right)$$

por la definición de la función Γ (obten por el hecho de que a_{ii}/σ_{ii} tiene una densidad χ^2 con n grados de libertad). Por tanto la densidad de r_{ij} es

$$\frac{[\Gamma(\frac{1}{2}n)]^p |r_{ij}|^{\frac{1}{2}(n-p-1)}}{\prod_{i=1}^p \Gamma\left(\frac{1}{2}(n+1-i)\right) \pi^{\frac{1}{4}p(p-1)}}$$

TEOREMA.- Si X_1, \dots, X_n son independientes, cada una con distribución $N(\mu_j, (\sigma_{ii}\delta_{ij}))$, entonces la densidad de los coeficientes de correlación muestrales viene dada por la fórmula anterior, donde $n = N-1$.