

UNIVERSIDAD DE SEVILLA
FACULTAD DE MATEMATICAS

UNIVERSIDAD DE SEVILLA
FACULTAD DE MATEMATICAS
SECRETARIA

3-6-80

ENTRADA N.º 213

Tesis
20

R. 8751

" ALGUNAS TECNICAS SOBRE
DETECCION DE OUTLIERS "

JOAQUIN J. M. MUÑOZ GARCIA

Visado en Sevilla

Mayo de 1980

Fdo. Rafael Infante Macias

Fdo. Antonio Pascual Acosta

Tesis que presenta Joaquin J. M.
Muñoz Garcia para optar al grado
de Doctor en Ciencias Matematicas

Fdo. Joaquin J. M. Muñoz Garcia



UNIVERSIDAD DE SEVILLA
FACULTAD DE MATEMATICAS

" ALGUNAS TECNICAS SOBRE
DETECCION DE OUTLIERS "

JOAQUIN J. M. MUÑOZ GARCIA

Memoria para optar al grado de
Doctor en Ciencias Matematicas
realizada bajo la direccion de
los.

Prof. Dr. D. RAFAEL INFANTE MACIAS y
prof. Dr. D. ANTONIO PASCUAL ACOSTA.

Sevilla Mayo 1.980



ALGUNAS TECNICAS SOBRE

DETECCION DE OUTLIERS

JOAQUIN J. M. MUÑOZ GARCIA



Quiero expresar mi mas sincera gratitud al Pr. Dr. D. Rafael Infante Macias, director de esta memoria; por haber despertado en mi el interés por las Matematicas Aplicadas y por su constante estimulo y ayuda en la preparacion de esta.

Asimismo quiero expresar mi profundo agradecimiento al Pr. Dr. D. Antonio Pascual Acosta, a quien debo gran parte de mi formación sobre Calculo de Probabilidades y Estadística Matemática y por contribuir en la realización y corrección de esta memoria.

Por ultimo quiero hacer llegar mi reconocimiento a todos aquellos que de un modo u otro, han contribuido a la realizacion de este trabajo , en especial a D.^o Luis Parras Guijosa.

Sevilla, Mayo de 1.980



A mi abuelo



CONTENIDO.

PROLOGO.

INTRODUCCION A LAS TECNICAS DE DETECCION DE OUTLIERS

1.1 Planteamiento general del problema	3
1.2 Comentarios históricos sobre técnicas de detección de outliers	7

DETECCION DE OUTLIERS MEDIANTE LA FUNCION INFLUENCIA

2.1 Introducción	21
2.2 Calculo de los parametros caracteristicos de la distribución perturbada	23
2.3 Funcion influencia para la distancia S^2	25
2.4 Función influencia para la media discriminante ...	28
2.5 Función influencia para los coeficientes de la función discriminante	37



2.6 Detección de outliers	40
---------------------------------	----

DETECCION DE OUTLIERS MEDIANTE EL COCIENTE DE VEROSIMILITUDES

3.1 Introducción	44
3.2 Contraste de cociente de verosimilitudes	45
3.3 Detección de outliers	53
3.4 Aplicación	56

DETECCION DE OUTLIERS BASADO EN LA DISTANCIA ENTRE MATRICES

SIMETRICAS Y DEFINIDAS POSITIVAS

4.1 Introducción	70
4.2 Distancia entre matrices de suma de cuadrados y suma de productos	70
4.3 Acotación	79
4.4 Detección de outliers	90
4.5 Aplicación	96

REFERENCIAS BIBLIOGRAFICAS	99
----------------------------------	----



PROLOGO

Uno de los problemas que preocupan hoy día fundamentalmente, a un estadístico cuando va a analizar un conjunto de datos, es la aparición de outliers o conjunto de observaciones que "parecen ser inconsistentes" con el resto del conjunto de datos. Lo que verdaderamente caracteriza a un outliers, es el "impacto" que produce en el estadístico cuando va a analizar la muestra.

Es evidente que la presencia de outliers en un conjunto de datos puede conducir a errores en nuestro intento de hacer inferencias acerca de la población de la que se han extraído, debido a que podrían falsear fuertemente las estimaciones o contrastes que hagamos sobre los parámetros poblacionales.

La presencia de outliers plantea por tanto un problema fundamental en el análisis de los datos, por lo que es natural buscar medios para



detectar e interpretar la presencia de dichas observaciones , rechazandolas en algunas ocasiones para así restablecer las propiedades de los datos , o por lo menos teniendo en cuenta su presencia en todos los análisis estadísticos.

El objetivo de esta Memoria se centra en la obtención de nuevos criterios para la detección de outliers en poblaciones normales multivariantes.

En el primer capítulo despues de plantear de forma general el problema que conlleva la presencia de outliers en un conjunto de datos, se pasa a describir el esquema que ha de seguir cualquier regla o técnica de detección de outliers. En el último apartado del capítulo incluimos un breve resumen con comentarios históricos sobre estas técnicas primero para el caso de poblaciones univariantes y en el segundo epigrafe para poblaciones multivariantes.

El capítulo segundo de la Memoria se dedica a la obtención de criterios para la detección de outliers basados en el concepto de función influencia. Estas reglas estan dadas para situaciones en las que se requiere aplicar una técnica de analisis discriminante, y han de aplicarse una vez obtenida la muestra, antes del calculo de la función discriminante donde ya no podrán ser incluidas las observaciones consideradas outliers.

Utilizando como técnica el test de cociente de verosimilitudes, el capítulo tercero se dedica al estudio de un criterio para detectar outliers en muestras procedentes de poblaciones normales multivariantes. Se incluye al final un programa de ordenador para detección de outliers mediante este metodo y se aplica en el último epigrafe a una situación práctica



En el último capítulo de la Memoria y utilizando como estadístico el máximo del cuadrado de la distancia entre las matrices de sumas de cuadrado y sumas de productos de observaciones muestrales, damos un procedimiento para detectar la presencia de outliers en muestras procedentes de poblaciones normales bivariantes. Después de dar un programa de ordenador para la resolución de esta técnica se aplica dicho procedimiento en una situación práctica.

CAPITULO PRIMERO

INTRODUCCION A LAS TECNICAS DE

DETECCION DE OUTLIERS

CONTENIDO.

1.1 - Planteamiento general del problema

1.2 - Comentarios históricos sobre técnicas
de detección de outliers.

- Outliers en poblaciones univariantes

- Outliers en poblaciones multivariantes

1.1. PLANTEAMIENTO GENERAL DEL PROBLEMA.

Según ANSGOMBE (1.960) existen tres fuentes de variabilidad sobre los elementos de una muestra o de un conjunto de datos extraídos de una cierta población y que se pueden clasificar de la siguiente forma:

Variabilidad intrínseca.- Es la variación inherente a la población. Tal variación surge de una forma natural y no puede ser reducida sin modificar la población objeto del estudio.

Variabilidad debida al error en las medidas.- Es el error que lleva consigo la falta de precisión de los instrumentos de medida utilizados para realizar el experimento. Dentro de esta fuente de variabilidad podemos incluir el error de transcripción de las observaciones.

Variabilidad debida al error de ejecución.- Es debida a una recolección imperfecta de los datos. Puede ocurrir que tomemos una muestra sesgada o que incluyamos en la misma observaciones no representativas de la población que estamos estudiando.

Basandose en estas fuentes de variabilidad Anscombe distingue dos tipos de observaciones

- Outliers: que son aquellas observaciones que presentan una gran variabilidad de tipo intrínseco.
- Falsas observaciones: que son debidas a un gran error de ejecución y de medida.

Sin embargo al igual que BARNETT (1.978) discrepamos de esta terminología debido al hecho que ante un conjunto de datos extraídos de una cierta población, generalmente no podemos determinar hasta donde llega una fuente de variabilidad y donde comienza la otra. Por ello seguimos en esta memoria el criterio de llamar outliers a aquella o aquellas observaciones que presentan una gran variabilidad de cualquiera de los tipos descritos por Anscombe, con respecto al conjunto principal de datos.

De entre las diversas definiciones del concepto "outliers" que aparecen a lo largo de las numerosas publicaciones estadísticas sobre el tema extraemos a continuación, las que consideramos mas interesantes.

Para GUMBEL (1.960) " los outliers son aquellos valores que parecen demasiado grandes o demasiado pequeñas comparadas con el resto de las observaciones".

En 1.961 FERGUSON da la siguiente idea de outliers " En una muestra extraída de una cierta población aparecen una o varias observaciones que sorprendentemente , se encuentran lejos del grupo principal de datos"

GRUBBS (1.969) afirma que un outliers " es una observación que se presenta fuertemente desviada de los otros miembros de la muestra".

Por ultimo Barnett (1.978) , considera un outliers " como una observación o conjunto de observaciones que parecen ser inconsistentes con el resto del conjunto de datos".

De estas definiciones se deduce que lo que caracteriza a una observación outliers es el "impacto" que produce en el estadístico cuando va a analizar los datos.

Es evidente que la presencia de outliers en un conjunto de datos, puede conducirnos a errores en nuestro intento de hacer inferencias acerca de la población de la que se han extraído, debido a que falsearan fuertemente las estimaciones o contrastes que hagamos sobre los parámetros poblacionales.

La presencia de outliers plantea por tanto un problema fundamental en el análisis de datos.

Podemos distinguir para la resolución de este problema dos grandes áreas: Reglas o técnicas de detección de outliers y estimación robusta. El objetivo de esta memoria se encuentra dentro de la primera area y por ello trataremos de una manera sucinta el problema de la estimación ro-



busta, en este primer capítulo.

Por técnica de detección de outliers entendemos un procedimiento que permite detectar de una manera objetiva la presencia de outliers en una muestra y determinar qué observaciones son outliers. En esta definición hemos de destacar la palabra "objetiva", debido a que existen autores, que no consideran la existencia de outliers sino simplemente de falsas observaciones en el sentido dado por Anscombe y que piensan que éstas deben ser detectadas de forma "subjetiva" por la experiencia del observador.

Para dar cualquier técnica de detección de outliers, en general se ha de seguir el siguiente esquema:

Se determina un estadístico que debe incluir las observación u observaciones que suponemos que son outliers, a continuación se halla la distribución que sigue dicho estadístico bajo la hipótesis de que todas las observaciones constituyen una muestra procedente de una misma población.

Para un nivel de significación α prefijado, que según GRUBBS (1.969) debe variar entre el 1% y el 5%, y que puede interpretarse como el riesgo de considerar una buena observación como outliers, se determina el valor crítico o percentil de la distribución teórica del estadístico correspondiente a esa probabilidad α . Esto nos permite determinar la región crítica de esta regla de detección.

Si para la muestra dada, el valor del estadístico, cae en la región crítica, diremos que existen observaciones outliers, para un nivel de significación α . Una vez detectada la presencia de outliers la técnica o regla utilizada nos permitirá en cada caso determinar qué observación u observaciones pueden ser consideradas outliers.

Una vez que se sabe que observaciones son outliers, el procedimiento a seguir podría ser uno de los siguientes.

Una primera solución sería desprestigiar dichas observaciones y tomar unas nuevas en su lugar, si esto es posible, y aplicar otra vez la técnica de detección de outliers a este nuevo conjunto de datos. No obstante aconsejamos no desprestigiar las observaciones consideradas outliers y reservarlas para un análisis posterior de las mismas ya que este nos puede conducir a determinar las fuentes de variabilidad que han dado lugar a la presencia de outliers, pudiéndolas corregir para un futuro experimento.

Otra solución y pensamos que esta se ha de realizar cuando detectamos un número de outliers próximo a $\left[\frac{n}{2} \right]$ (parte entera de $\frac{n}{2}$), es tomar una nueva muestra y realizar de nuevo el estudio. Si no es posible tomar esta nueva muestra, entonces hemos de cambiar el modelo de distribución supuesto al principio por otro modelo que puede ser incluso una mixtura de distribuciones.

Como indicamos anteriormente otra forma posible de abordar el problema es basándose en la robustez. Puede ocurrir que no sea posible realizar una técnica para detectar outliers o aunque existan outliers, que queramos retenerlos en la muestra porque pensemos que los mismos contienen una información que nos interesa o por cualquier otra circunstancia. Estos problemas pueden resolverse mediante la estimación robusta de los parámetros poblacionales sobre los que queremos realizar la inferencia.

1.2. COMENTARIOS HISTORICOS SOBRE TECNICAS DE DETECCION DE OUTLIERS.

Incluimos en este apartado un resumen histórico sobre las diferentes reglas o métodos para la detección de outliers. Conviene indicar que aunque los inicios de estas técnicas para poblaciones univariantes pueden remontarse hasta mediados del siglo XVIII no es hasta la década de los cincuenta, cuando aparecen las primeras referencias sobre detección de outliers para poblaciones multivariantes, debido a los problemas que surgen en el caso multivariante y a los que haremos referencia mas adelante.

1.2.1. Outliers en poblaciones univariantes.

Los inicios de estas técnicas se pueden encontrar en el siglo XVIII cuando Boscovich en 1755, intenta determinar la elipticidad de la Tierra promediando las medidas de exceso de los grados polares alrededor del ecuador, descartando los valores extremos. Daniel Bernouilli en su trabajo " *Dijudicatio maxime probabilis plurium observationum discrepantium atque verisimilima inductio inde formata* " de 1777, aprueba el rechazo de observaciones extremas. Y Legendre en su obra " *Nouvelles methodes pour la determination des orbites des cometes* " en 1.805 recomienda tambien el rechazo de observaciones.

Sin embargo la escuela germana de astrónomos en un trabajo publicado por Bessel hacia 1.838, afirma al efecto, que nunca una observación será rechazada por su gran residuo pues el conjunto de todas las observaciones contribuirán al resultado y afirman al mismo tiempo. " Creemos que únicamente mediante una observación estricta de esta regla, se elimina la arbitrariedad de los resultados".

Vemos como en estos comienzos el rechazo de observaciones se realiza únicamente a criterio del observador. Así Saunder en su trabajo de 1.90

llega a escribir que la experiencia práctica del individuo ha de contar sobremanera para el rechazo de observaciones.

Sin embargo el primer intento de dar un criterio de rechazo basado en algun tipo de razonamiento probabilístico es debido a Peirce (1.852) El argumento de Peirce es reproducido por Chauvenet (1.863) el cual da una regla similar basada en un argumento mas simple. Stone hacia 1.868 introduce un test de rechazo basado en un módulo de descuido o negligencia, lo cual suponía un error al tomar las observaciones el experimentador. Pero podemos decir que el trabajo más serio en estos inicios es el Glaisher (1.872) " On the rejection of discordant observations " en el que trabaja ya con observaciones procedentes de poblaciones normales univariantes dando un procedimiento de medias ponderadas. Un año mas tarde Stone hace una critica a este metodo dando un procedimiento de ponderación ., basado en la maximización de la función de verosimilitud y posteriormente este mismo criterio fue dado tambien por Edgeworth en 1.883.

En esta misma linea podemos citar el trabajo de Wri^gth de 1.884 " A treatise on the adjustment of observations by the method of least squares " en el que rechaza las observaciones que se desvian de la media en mas de tres veces la desviación típica; y el procedimiento de Goodwin (1.913) que rechaza una observacion en una muestra de tamaño n si la desviación de la media de las restantes $n - 1$ observaciones excede cuatro veces a la desviación promedio de las $n - 1$ observaciones

Hemos de señalar que todos estos procedimientos en cierto modo rudimentarios adolecen de un defecto general que cònsiste en no distinguir la varianza poblacional y la varianza muestral. Es IRWIN (1.925) quien da el primer procedimiento para detección de outliers en el que ya

distingue entre varianza muestral y varianza poblacional. Para el caso en que σ es conocido propone el test estadístico,

$$\frac{[X_{(n)} - X_{(n-1)}]}{\sigma} \quad \text{y} \quad \frac{[X_{(n-1)} - X_{(n-2)}]}{\sigma}$$

donde $X_{(i)}$ denota el estadístico ordenado de rango i . Conviene citar también a STUDENT (1.927) que da un criterio basado en el recorrido muestral y THOMPSON (1.935) que da un test exacto studentizado,

$$\frac{x - \bar{x}}{S}$$

A partir de este año los estadísticos prestan una mayor atención a estos problemas y las técnicas de detección de outliers unidimensionales se hacen más rigurosas.

GRUBBS (1.950) da un criterio basado en el cociente entre la suma de cuadrados de desviaciones para una muestra reducida y la suma de cuadrados de las desviaciones para la muestra completa, DIXON (1.950) da varias reglas basadas en la razón entre la observación "sospechosa" y su más próxima y el recorrido o rango muestral y DAVID, PEARSON y HARTLEY (1.954) dan un test para la detección de outliers en poblaciones normales univariantes basado en la razón entre el recorrido muestral y la desviación típica muestral.

KUDO (1.956), da unas ciertas reglas para la detección de outliers en poblaciones normales univariantes tanto para el caso de conocer la varianza poblacional como cuando no se conoce. Dichos métodos se basan en el ordenamiento de las observaciones. Kudo describe cinco criterios por los que una regla de rechazo puede considerarse óptima.

Como caso particular de este trabajo de Kudo se obtiene el estadístico de PEARSON y CHANDRASEKAR (1.936) que ha tenido una gran importancia

en el desarrollo de las reglas de rechazo para outliers, probando que este estadístico es óptimo según los criterios por él dados y demostrando que el de NAIR (1.948) no posee esta propiedad de optimalidad.

En el año 1.960, aparece un trabajo de ANSCOMBE que ha tenido una importancia fundamental en el desarrollo ulterior de esta teoría. Comienza dando los tipos de variabilidad que se pueden presentar en un conjunto de datos y a los que ya hemos hecho referencia en el apartado anterior, introduciendo a continuación dos conceptos que él denomina "premium" y "protection" con los que trata de dar un criterio para poder estudiar y comparar todas las reglas de detección de outliers. Define el "premium" como el porcentaje en que se incrementa la varianza del error de estimación debido al uso del criterio de detección cuando de hecho todas las observaciones proceden de una misma población normal y la "protection" como la reducción en la varianza o error cuadrático medio cuando aparecen observaciones outliers.

ANSCOMBE Y TUKEY (1.963) definen el residuo de una observación de la siguiente forma.

$$(\text{residuo}) = (\text{valor observado}) - (\text{valor ajustado})$$

En este trabajo, dan algunas técnicas de tipo gráfico y otras de tipo numérico para discutir los residuos que resultan de las medias por filas y columnas en un análisis de varianza de dos factores o más. Afirman que: "La razón más importante para el cálculo de los residuos es detectar outliers, observaciones que tienen un gran residuo en comparación con la mayor parte de las otras y esto debe ser tratado especialmente". En este mismo trabajo dan ciertos estadísticos para la detección de outliers.

Esta definición de residuo es generalizada por COX y SNELL (1.968), indicando los campos de aplicabilidad de esta definición. Estudian las



propiedades de los residuos, haciendo especial hincapié en los residuos de poblaciones Binomial y Poisson; es interesante la discusión que de dicho trabajo hacen destacados especialistas en el tema como , Anscombe F. J., Barnett V., Mallows C. L. , Pearce S. C., Harrison P. J., entre otros.

Basandose en los trabajos de Anscombe , TIAO y GUTTMAN (1.967) trata de dar una estimación para la media de una población normal unidimensional a partir de una muestra de tamaño n , con la posibilidad de que en dicha muestra existan uno o más outliers. Dan unas reglas de rechazo basadas en el analisis de los residuos, que suponen estan correlacionados. Tambien hacen un estudio del " premium " y " protection " de los metodos de detección de outliers que proponen.

En esta linea se encuentra el trabajo de ANDREWS (1.971), que basandose en la distribución conocida de los residuos en el modelo de regresión lineal, da tests de significación exactos, aplicando a continuación dichos resultados para contrastar la existencia de uno o mas outliers.

TIETJE Y MOORE (1.972), basandose en la idea de outliers dada por KENDALL Y BUCKLAND (1.957) definen el efecto de "enmascaramiento" , como la imposibilidad de una regla de rechazo de identificar una observación outliers en presencia de varios valores "sospechosos", ilustrando dicho efecto con un ejemplo sobre isotopos del uranio, aplicandole el estadístico de GRUBBS (1.950).

BROWN (1.975) da tests estadísticos para la detección de outliers basandose en los residuos, que resultan en el analisis de varianza de dos factores, pero a diferencia de otros tests basados en residuos, Brown no se basa en el valor del residuo, sino en el signo de los residuos a que da lugar cada una de las filas y columnas. Obteniendo

así un estadístico que se distribuye aproximadamente como una ley χ^2 .

En el año 1.975 ROSNER publica un artículo de bastante importancia en la detección de outliers en muestras univariantes, pues en él compara la potencia de varias reglas de rechazo.

Comienza exponiendo el problema de detección de outliers en muestras de poblaciones normales. Así afirma que.

" Si x_1, x_2, \dots, x_n es una muestra aleatoria de una muestra compuesta de dos subconjuntos $J_1 = \{x_{11}, x_{12}, \dots, x_{1m}\}$, $J_2 = \{x_{j1}, x_{j2}, \dots, x_{jk}\}$ donde $x_{iq} \in N(\mu, \sigma^2)$ $q = 1, 2, \dots, m$ y $x_{jl} \in N(\mu_{jl}, \sigma_{jl}^2)$ $l = 1, 2, \dots, k$ donde todos los parámetros son desconocidos. Llamaremos J_1 el subconjunto de las verdaderas observaciones y a J_2 subconjunto de los outliers ($0 \leq k < n$ y $m + k = n$). El problema de la detección de outliers consiste en identificar los subconjuntos J_1 y J_2 .

Generaliza cuatro tipos de estadísticos para la detección de un outlier, modificando después los estadísticos para que sirvan para detectar más de un outlier. Estos cuatro estadísticos son.

Desviación extrema studentizada (ESD).

$$ESD = \max_{i=1, \dots, n} \frac{|x_i - \bar{x}|}{S}$$

Recorrido studentizado (STR).

$$STR = \frac{(x_{(n)} - x_{(1)})}{S}$$

Kurtosis (KUR)

$$KUR = \frac{n \sum_{i=1}^n (x_i - \bar{x})^4}{\left[\sum_{i=1}^n (x_i - \bar{x})^2 \right]^2}$$

R - estadístico (RST)

$$RST = \max_{i: 1, \dots, n} \frac{|x_i - a|}{b}$$

donde a y b son la media y varianza muestral corregidas.

Y de estos estadísticos estudia su potencia mediante procedimientos de simulación, llegando a las conclusiones siguientes:

- 1) Entre los procedimientos para la detección de más de un outliers el procedimiento STR es claramente el inferior.
- 2) Los otros tres procedimientos se pueden considerar iguales, aunque el ESD es ligeramente mejor.
- 3) Para efectos computacionales da el siguiente orden de menos a más complejos, RST, ESD, STR, KUR.
- 4) El procedimiento preferido será el ESD, que es generalmente mejor y computacionalmente es muy razonable.

Hemos dado hasta aquí un breve esquema sobre reglas de detección de outliers para poblaciones normales univariantes. Una de las líneas de investigación en este campo del Análisis de Datos consiste en encontrar reglas de rechazo para muestras procedentes de poblaciones con diferentes leyes de probabilidad. Así el trabajo de SHAPIRO - WILK (1.972) para el caso de una muestra de tipo exponencial o el más reciente de COLLET (1.980) para conjunto de datos extraídos de una población con distribución de Von Mises.

1.2.2. Outliers en poblaciones multivariantes.

Según acabamos de ver la mayoría de los criterios de detección de outliers, para datos unidimensionales están basados en estadísticos ordenados o en los tamaños de los residuos.

En el caso de poblaciones multidimensionales surgen problemas no existentes en el caso anterior y que hacen aumentar la complejidad del estudio de reglas de rechazo para datos multidimensionales. Así por ejemplo; la imposibilidad de ordenar las observaciones, que una observación puede ser considerada outliers porque exista una gran variabilidad en una de sus componentes o porque existan pequeñas variabilidades en todas sus componentes etc. Conviene también señalar aunque parece obvio que una muestra p - variante no puede ser estudiada como un conjunto de p - muestras univariantes pues se perdería el sentido de dependencia o correlación entre las variables e incluso pueden aparecer outliers unidimensionales y sin embargo el vector del que es componente no lo sea.

Como GNANADESIKAN y KETTERING (1.972) señalan. " La complejidad del caso multivariante , lleva consigo que sea infructuoso buscar procedimientos de detección de outliers , que sean válidos para todas las situaciones. Mas razonable parece buscar procedimientos de detección de outliers contra situaciones de tipo específico, debiéndose construir un gran conjunto de técnicas con diferentes sensibilidades, de forma que un outlier para una cierta situación puede no serlo para otra."

Podemos decir que las primeras publicaciones sobre detección de outliers en este campo, son una generalización del trabajo realizado por FERGUSON en 1.961 sobre detección de outliers en poblaciones normales univariantes. Todos ellos tienen como base el exponente de

la distribución normal multivariante

$$Q(x, \mu, \Sigma) = Q(x) = (X - \mu)' \Sigma^{-1} (X - \mu)$$

y se considera este estadístico como una distancia dentro del espacio de las observaciones normales multivariantes. Dicha distancia está estudiada bajo los supuestos de que se conozcan los dos parámetros μ y Σ uno de ellos o ninguno de los dos.

Para todos estos casos se realiza un ordenamiento de las observaciones basado en la distancia antes descrita y se toma generalmente como estadístico para decidir si una observación es o no outliers.

$$T = \max_i Q(x_i)$$

Así si μ y Σ son conocidos, tendremos estadísticos ordenados de una distribución gamma, y en este caso se compara el estadístico T con un punto crítico dado por un nivel de significación fijado de antemano. Dicho punto crítico podemos obtenerlo de las tablas dadas por GUPTA (1.960), para la distribución de los estadísticos ordenados de una distribución gamma.

En el caso de que μ y Σ o uno de los dos no es conocido la distribución de T no está determinada de forma exacta, sin embargo SIOTANI (1.959) da una tabulación de una distribución aproximada a la de T .

Quizás el estudio más riguroso sobre detección de outliers, es el realizado por WILKS (1.963), quien da un estadístico para detectar k outliers $k = 1, 2, \dots$, basándose para ello en ideas geométricas. Así si se tiene una distribución normal p - dimensional, de la que extraemos una muestra de tamaño n ($n > p$), como estadístico para $k = 1$ considera lo que él llama razón de dispersión uni - outliers, y que

representa por

$$R_i = \frac{|S_{(n-1)}^{(i)}|}{|S_{(n)}|} \quad i = 1, 2, \dots, n$$

siendo la razón de dispersión k - outliers

$$R_i = \frac{|S_{(n-k)}^{(i)}|}{|S_{(n)}|} \quad i = 1, 2, \dots, \binom{n}{k}$$

donde $S_{(n)}$ es la matriz de sumas de cuadrados y sumas de productos de n observaciones muestrales y $S_{(n-1)}^{(i)}$ es la matriz de s. c. y s. p. de n-1 observaciones de la muestra de tamaño n donde se ha quitado la observación i - ésima.

El estadístico R_i se puede interpretar como la razón de los volúmenes entre paralelotojos, ya que WILKS (1.962), relaciona $|S_{(n)}|$ con el volumen de los paralelotojos formados con p - puntos de los n y el centro de gravedad de la muestra $\bar{X}' = (\bar{X}_1, \bar{X}_2, \dots, \bar{X}_p)$.

Con este estadístico Wilks considera que el mayor "candidato" a ser un outliers es aquel elemento muestral que alcanza el $\min_i R_i$ y será considerado outliers si es menor que un cierto punto crítico, determinado bajo un nivel de significación dado, mediante la distribución del $\min_i R_i$.

En 1.972 GNANADESIKAN Y KETTERING realizan un trabajo de recopilación, en el que uno de sus puntos trata sobre detección de outliers en poblaciones normales multivariantes. Estas técnicas están basadas en procedimientos gráficos entre los que podemos citar uno para detección de outliers en análisis discriminante y el basado en el análisis de las componentes principales.

Sobre este último método HAWKINS (1.974) hace algunas modificaciones

realizando aplicaciones sobre datos geológicos y dando una comparación entre diferentes estadísticos que surgen a lo largo de su trabajo.

Es obvio que una forma gráfica para detectar outliers en distribuciones bidimensionales o tridimensionales sería representar la nube de puntos de las observaciones y ver que valores de la variable se aleja de la misma, sin embargo para dimensiones mayores esto nos resulta imposible, ANDREWS (1.972) propone una técnica para representar observaciones p - dimensionales sin ninguna restricción para p . Afirma que un dato p - dimensional $X' = (X_1, X_2, \dots, X_p)$ define una función.

$$f_X(t) = \frac{X_1}{\sqrt{2}} + X_2 \sin t + X_3 \cos t + X_4 \sin 2t + X_5 \cos 2t + \dots$$

siendo dibujada esta función en el recorrido $-\pi < t < \pi$. Andrews aplica este método a un trabajo de antropología. Esta representación gráfica de puntos p - dimensionales es aprovechada por GNANADESIKAN (1.973) para dar otra técnica de tipo gráfico para la detección de outliers.

DEVLIN y otros (1.975) describen un procedimiento para detección de outliers, al utilizar el coeficiente de correlación, basándose en la función influencia dada por HAMPEL (1.974) aunque ellos utilizan la función influencia muestral dada por

$$I_i(y_i, \hat{\theta}) = (n-1) (\hat{\theta} - \hat{\theta}_i) \quad i = 1, 2, \dots, n.$$

Este método también está basado en representaciones gráficas de esta función influencia. En este mismo trabajo podemos ver una aplicación al ejemplo que toman de FISHER (1.936) sobre la longitud y anchura de los sepalos de cincuenta plantas de Iris Setosa.

Por ultimo para acabar esta breve reseña histórica citaremos las obras de GNANADESIKAN (1.977) en la que expone diversas técnicas para la detección de outliers multivariantes y la más reciente de BARNETT (1.978) dedicada toda ella al estudio de outliers en el Analisis Estadístico de Datos.



CAPITULO SEGUNDO

DETECCION DE OUTLIERS MEDIANTE

LA FUNCION INFLUENCIA

CONTENIDO.

- 2.1 - Introduccion.
- 2.2 - Calculo de los parametros caracteristicos de la distribución perturbada.
- 2.3 - Función influencia para la distancia S^2
 - Calculo de $S^2(F_1^*)$
 - Funcion influencia para S^2
 - Distribución de $CI_{S^2 F_1}$
- 2.4 - Funcion influencia para la media discriminante
 - Calculo de $b^{\mu_1}(F_1^*)$ y $b^{\mu_2}(F_1^*)$
 - Funcion influencia para b^{μ_1} y b^{μ_2}
 - Distribucion de $CI_{b^{\mu_1} F_1}$ y $CI_{b^{\mu_2} F_1}$
- 2.5 - Funcion influencia para los coeficientes de la funcion discriminante.
 - Calculo de $b^*b(F_1^*)$
 - Funcion influencia para b^*b
 - Distribución de $CI_{b^*b F_1}$
- 2.6 - Deteccion de outliers

2.1. INTRODUCCION.

El objetivo de este capítulo se centra en la obtención de criterios para la detección de outliers, basandonos en el concepto de función influencia introducida por HAMPEL (1.974). Estos criterios están dados, para situaciones en las que se requiera aplicar una técnica de análisis discriminante, y han de aplicarse una vez obtenida la muestra; antes del cálculo de la función discriminante, donde ya no serán incluidas las observaciones consideradas outliers.

A continuación introducimos algunos de los conceptos, que se van a utilizar a lo largo del capítulo.

ANDERSON y BAHADUR (1.962), dan un procedimiento para discriminar entre dos poblaciones normales multivariantes, $N_p(\mu_1, \Sigma_1)$ y $N_p(\mu_2, \Sigma_2)$ no singulares, con vectores medias y matrices de covarianza distintas. Formalmente el problema se resuelve de la siguiente manera:

Sea b un vector p - dimensional de componentes reales y c una constante real, una observación x p - dimensional es clasificada como perteneciente a la primera población si $b'x \leq c$ y de la segunda si $b'x > c$.

Las posibles probabilidades de errores en el procedimiento de clasificación son.

$$\epsilon_1 = P[X \in N_p(\mu_2, \Sigma_2) / X \in N_p(\mu_1, \Sigma_1)] = 1 - F_{N_1(0,1)}\left(\frac{c - b'\mu_1}{(b'\Sigma_1 b)^{1/2}}\right)$$

$$\epsilon_2 = P[X \in N_p(\mu_1, \Sigma_1) / X \in N_p(\mu_2, \Sigma_2)] = 1 - F_{N_1(0,1)}\left(\frac{b'\mu_2 - c}{(b'\Sigma_2 b)^{1/2}}\right)$$

Se trata de minimizar las probabilidades ϵ_1 y ϵ_2 o lo que es lo mismo maximizar los argumentos.

$$Z_1 = \frac{c - b'\mu_1}{(b'\Sigma_1 b)^{1/2}}$$

$$Z_2 = \frac{b'\mu_2 - c}{(b'\Sigma_2 b)^{1/2}}$$

Anderson y Bahadur demuestran que el procedimiento definido por.

$$b = (\lambda_1 \Sigma_1 + \lambda_2 \Sigma_2)^{-1} (\mu_2 - \mu_1)$$

con

$$c = b' \mu_2 - \lambda_2 b' \Sigma_2 b = b' \mu_1 + \lambda_1 b' \Sigma_1 b$$

para λ_1 y λ_2 arbitrarios tales que $\lambda_1 \Sigma_1 + \lambda_2 \Sigma_2$ sea definida positiva, es un procedimiento admisible.

El problema surge al determinar λ_1 y λ_2 de forma que se cumplan las condiciones anteriores. Para su resolución existen varios métodos. Nos basaremos en el procedimiento minimax que consiste en igualar Z_1 y Z_2 por lo que las probabilidades de clasificación correcta serían mayores que $1/2$ y se verifica $\lambda_1 = 1 - \lambda_2 = \lambda$.

$$0 = Z_1^2 - Z_2^2 = b' [\lambda^2 \Sigma_1 - (1-\lambda)^2 \Sigma_2] b$$

y esta ecuación tiene solución única ya que Z_1^2 es creciente con λ y Z_2^2 es decreciente con λ .

Por otro lado CHERNOFF (1.972, 1973) utiliza como medida para discriminar entre dos poblaciones $N_p(\mu_1, \Sigma_1)$, $N_p(\mu_2, \Sigma_2)$ no singulares, la distancia S^2 .

$$S^2 = \lambda(1-\lambda) \delta' (\lambda \Sigma_1 + (1-\lambda) \Sigma_2) \delta$$

donde t verifica las mismas condiciones que las del procedimiento minimax descrito anteriormente y por δ representamos el vector diferencia

$$\delta = \mu_2 - \mu_1$$

Otro concepto de gran importancia a lo largo de este capítulo es el de función influencia que podríamos definir en la forma siguiente.

Sea Ω un espacio métrico separable y completo, sea T una función vectorial definida en un subconjunto del espacio de todas las medidas de probabilidad sobre Ω y con valores en el espacio euclideo k -dimensional \mathbb{R}^k y sea F una función del dominio de T . Si notamos δ_w la distribución degenerada en un punto arbitrario $w \in \Omega$, la curva de influencia de T en F se define, mediante el límite.

$$CI_{TF}(w) = \lim_{\varepsilon \rightarrow 0} \left\{ \frac{T[(1-\varepsilon)F + \varepsilon\delta_w] - T(F)}{\varepsilon} \right\}$$

si este limite existe para cada $w \in \Omega$.

La funcion influencia trata pues de reflejar la posible "influencia" o "perturbacion" que experimentan los parametros poblacionales, bajo la hipotesis de que la variable proceda de una poblacion con funcion de distribucion $(1-\varepsilon)F + \varepsilon\delta_w$ en lugar de F . $CI_{TF}(\cdot)$ va a depender del punto en el que suponemos està definida la distribucion degenerada, por lo que dicha funcion puede considerarse como una transformacion de la variable aleatoria, y tendrà por tanto una distribucion de probabilidad.

El procedimiento a seguir consiste en determinar la distribucion de la funcion influencia bajo la hipotesis de que no existen outliers es decir que la variable no se distribuye segun la distribucion $(1+\varepsilon)F + \varepsilon\delta_w$. Para lo cual serà necesario que estudiemos, como son perturbados los parametros poblacionales al considerar la distribucion $(1-\varepsilon)F + \varepsilon\delta_w$, calculando a continuacion la funcion influencia CI_{TF} y por ultimo su distribucion.

2.2. CALCULO DE LOS PARAMETROS CARACTERISTICOS DE LA DISTRIBUCION PERTURBADA.

Sean dos poblaciones $N_p(\mu_1, \Sigma_1)$ y $N_p(\mu_2, \Sigma_2)$ no singulares.

Representemos por $F_1^*(w)$, la funcion de distribucion perturbada de la funcion de distribucion $F_1(w)$ de la poblacion $N_p(\mu_1, \Sigma_1)$.

Definida por

$$F_1^*(w) = (1-\varepsilon)F_1(w) + \varepsilon\delta_w$$

donde $\varepsilon \in \mathbb{R}$ con $0 \leq \varepsilon \leq 1$.

Pasemos ahora a calcular los parametros caracteristicos de la distribución F_1^* que van a ser necesarios para la determinación de la curva influencia.

Así la media de esta distribución será.

$$\mu_1(F_1^*) = \int w dF_1^*(w) = \mu_1 + \epsilon(w - \mu_1)$$

Por otro lado la matriz de covarianza de la población perturbada, que notaremos por $\Sigma_1(F_1^*)$, toma la forma

$$\Sigma_1(F_1^*) = \int (w - \mu_1)(w - \mu_1)' dF_1^*(w) = (1 - \epsilon)\Sigma_1 + \epsilon(w - \mu_1)(w - \mu_1)'$$

Si notamos

$$w - \mu_1 = X$$

tendremos

$$\mu_1(F_1^*) = \mu_1 + \epsilon X$$

$$\Sigma_1(F_1^*) = (1 - \epsilon)\Sigma_1 + \epsilon X X'$$

Veamos a continuación la forma que tendría la matriz $\Sigma = \lambda \Sigma_1 + (1 - \lambda)\Sigma_2$, bajo la distribución F_1^* , ya que dicha matriz va a ser utilizada en desarrollos posteriores. Así

$$\Sigma(F_1^*) = \lambda \Sigma_1(F_1^*) + (1 - \lambda)\Sigma_2 = \Sigma - \lambda \epsilon \Sigma_1 + \epsilon \lambda X X'$$

cuya inversa es:

$$\{\Sigma(F_1^*)\}^{-1} = \{\Sigma - \lambda \epsilon \Sigma_1 + \epsilon \lambda X X'\}^{-1} = \{\Sigma - \epsilon \lambda \Sigma_1\}^{-1} \cdot$$

$$\cdot \{I_p + \epsilon \lambda (\Sigma - \epsilon \lambda \Sigma_1)^{-1} X X'\}^{-1} = (I_p + \epsilon \lambda X X' (\Sigma - \epsilon \lambda \Sigma_1)^{-1})^{-1} (\Sigma - \epsilon \lambda \Sigma_1)^{-1}$$

y aplicando la igualdad matricial (BELLMAN 1.965)

$$(I + AB)^{-1} = I - A(I + AB)^{-1}B$$

se obtiene

$$\{\Sigma(F_1^*)\}^{-1} = \{I_p - \epsilon \lambda (\Sigma - \epsilon \lambda \Sigma_1)^{-1} X (I_p + \epsilon \lambda X X' (\Sigma - \epsilon \lambda \Sigma_1)^{-1} X X'\}^{-1} \cdot$$

$$\cdot (\Sigma - \epsilon \lambda \Sigma_1)^{-1} =$$

$$= \frac{(\Sigma - \epsilon \lambda \Sigma_1)^{-1} + \epsilon \lambda X' (\Sigma - \epsilon \lambda \Sigma_1)^{-1} X (\Sigma - \epsilon \lambda \Sigma_1)^{-1}}{1 + \epsilon \lambda X' (\Sigma - \epsilon \lambda \Sigma_1)^{-1} X} -$$

$$- \frac{\epsilon \lambda (\Sigma - \epsilon \lambda \Sigma_1)^{-1} X X' (\Sigma - \epsilon \lambda \Sigma_1)^{-1}}{1 + \epsilon \lambda X' (\Sigma - \epsilon \lambda \Sigma_1)^{-1} X} \quad (1)$$

Desarrollando (1) en serie de Mac - Laurin en función de ϵ hasta el término de primer orden ya que los restantes vendrían dados en función de ϵ^H ($H > 1$), y estos términos se harían cero al efectuar el cálculo de la función influencia, resulta

$$\left\{ \Sigma(F_1^*) \right\}^{-1} = (I_p + \epsilon \lambda \Sigma^{-1} \Sigma_1) \Sigma^{-1} - \epsilon \lambda \Sigma^{-1} X X' \Sigma^{-1}$$

habiéndose utilizado para ello

$$\frac{d}{d\epsilon} (\Sigma - \epsilon \lambda \Sigma_1)^{-1} = (\Sigma - \epsilon \lambda \Sigma_1)^{-1} \lambda \Sigma_1 (\Sigma - \epsilon \lambda \Sigma_1)^{-1}$$

Por último el vector δ se transforma a su vez en

$$\delta(F_1^*) = \mu_2 - \mu_1(F_1^*) = \delta - \epsilon X$$

2.3. FUNCION INFLUENCIA PARA LA DISTANCIA S^2 .

En este apartado se estudia en primer lugar como sería esta distancia S^2 , bajo la hipótesis de que la distribución de la población es de la forma $(1-\epsilon)F_1(w) + \epsilon \delta_w$, para ello bastará utilizar los parámetros perturbados, que hemos obtenido en el apartado 2.2. A continuación determinamos la función influencia para este parámetro, y por último calculamos la distribución de dicha función influencia.

2.3.1. Cálculo de $S^2(F_1^*)$.

Por definición

$$S^2 = \lambda(1-\lambda) \delta' \Sigma^{-1} \delta$$

y basándonos en los resultados del apartado anterior

$$\begin{aligned}
S^2(F_1^*) &= \lambda(1-\lambda)(\delta - \varepsilon X)' \left\{ (I_p + \varepsilon \lambda \Sigma^{-1} \Sigma_1) \Sigma^{-1} - \varepsilon \lambda \Sigma^{-1} X X' \Sigma^{-1} \right\} \\
(\delta - \varepsilon X) &= \lambda(1-\lambda) \left\{ \delta' (I_p + \varepsilon \lambda \Sigma^{-1} \Sigma_1) \Sigma^{-1} \delta - \varepsilon \lambda \delta' \Sigma^{-1} X X' \Sigma^{-1} \delta - \right. \\
&- \varepsilon X' (I_p + \varepsilon \lambda \Sigma^{-1} \Sigma_1) \Sigma^{-1} \delta + \varepsilon^2 \lambda X' \Sigma^{-1} X X' \Sigma^{-1} \delta - \\
&- \varepsilon \delta' (I_p + \varepsilon \lambda \Sigma^{-1} \Sigma_1) \Sigma^{-1} X + \varepsilon^2 \lambda \delta' \Sigma^{-1} X X' \Sigma^{-1} X + \\
&+ \varepsilon^2 X' (I_p + \varepsilon \lambda \Sigma^{-1} \Sigma_1) \Sigma^{-1} X - \varepsilon^3 \lambda X' \Sigma^{-1} X X' \Sigma^{-1} X \left. \right\}
\end{aligned}$$

Y debido a que en el calculo de la funcion influencia tomaremos limite en ε para cuando $\varepsilon \rightarrow 0$ podemos despreciar en $S^2(F_1^*)$, los terminos con potencias de ε mayores que la unidad. Por lo que esta distancia perturbada quedaria

$$\begin{aligned}
S^2(F_1^*) &= \lambda(1-\lambda) \left\{ \delta' \Sigma^{-1} \delta + \varepsilon \lambda (\delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta) - \right. \\
&- \left. \varepsilon \lambda (\delta' \Sigma^{-1} X)^2 - 2\varepsilon \delta' \Sigma^{-1} X \right\}
\end{aligned}$$

2.3.2. Funcion influencia para S^2 .

$$\begin{aligned}
CI_{S^2; F_1}(w) &= \lim_{\varepsilon \rightarrow 0} \left\{ \frac{S^2(F_1^*) - S^2(F_1)}{\varepsilon} \right\} = \lambda(1-\lambda) \cdot \\
&\cdot \left\{ \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta - \lambda (\delta' \Sigma^{-1} X)^2 - 2\delta' \Sigma^{-1} X \right\}
\end{aligned}$$

Y si realizamos el cambio de variable $Y = \delta' \Sigma^{-1} X$ nos resultaria

$$CI_{S^2; F_1}(w) = \lambda(1-\lambda) \left\{ \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta - \lambda Y^2 - 2Y \right\}$$



2.3.3. Distribucion de CI_{S^2, F_1}

Calculamos ahora la distribución de la variable $CI_{S^2, F_1}(w)$, bajo la hipótesis de que $W \in N_p(\mu_1, \Sigma_1)$. Esta distribución que será utilizada como veremos en el apartado 2.6 para la detección de outliers.

si $W \in N_p(\mu_1, \Sigma_1)$

entonces

$$X \in N_p(0, \Sigma_1) \text{ e } Y \in N_1(0, \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta)$$

Y si llamamos Z a la variable Y tipificada se tiene

$$Z = \frac{Y}{(\delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta)^{1/2}}$$

La función influencia en función de Z nos quedaría de la forma

$$CI_{S^2, F_1}(w) = \lambda(1-\lambda) \left\{ \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta - \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta Z^2 - 2(\delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta)^{1/2} Z \right\} = \lambda(1-\lambda) \left\{ \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta + \frac{1}{\lambda} - \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta \left(Z + \frac{1}{\lambda(\delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta)^{1/2}} \right)^2 \right\}$$

y llamando

$$CI'_{S^2, F_1}(w) = \left\{ Z + \frac{1}{\lambda(\delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta)^{1/2}} \right\}^2$$

y dado que $Z \in N_1(0, 1)$ tendremos que

$$CI'_{S^2, F_1}(w) \in \chi^2 \left(1, \frac{1}{\lambda^2 \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta} \right)$$

A continuación damos un lema con la intención de aproximar la ley chi - cuadrado no centrada mediante una distribución gamma.

2.3.3.1. Lema. Sea X una variable aleatoria distribuida según una ley chi - cuadrado no centrada con r grados de libertad y parámetro de descentralización η . La distribución de la variable cX donde c es una constante real cualquiera, se puede aproximar por una ley gamma

$$G \left(\frac{(r+\eta)^c}{2(r+\eta)} ; \frac{r+\eta}{2(r+\eta)c} \right)$$

La demostración es trivial basta tener en cuenta la aproximación de PATNAIK (1949).

Por tanto la función influencia de S^2 nos quedaría de la forma

$$CI_{S^2, F_1}(w) = \lambda(1-\lambda) \left\{ \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta + \frac{1}{\lambda} \right\} - CI'_{S^2, F_1}(\bar{w})$$

donde $CI'_{S^2, F_1}(w)$ se distribuirá en virtud del lema anterior como una gamma

$$CI'_{S^2, F_1}(w) \in \Gamma \left(\frac{(\lambda^2 \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta + 1)^2}{2 \lambda^2 \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta (\lambda^2 \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta + 2)}; \frac{\lambda^2 \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta + 1}{2 \lambda^2 (1-\lambda) \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta (\lambda^2 \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta + 2)} \right)$$

2.4. FUNCION INFLUENCIA PARA LA MEDIA DISCRIMINANTE.

Obtenemos en este apartado la función influencia y su distribución, para la media discriminante de la primera población $b' \mu_1$ y para la de la segunda $b' \mu_2$.

2.4.1 Calculo de $b' \mu_1(F_1^*)$ y $b' \mu_2(F_2^*)$.

Teniendo en cuenta los resultados del apartado 2.2

$$\begin{aligned} b' \mu_1(F_1^*) &= \delta'(F_1^*) \{ \Sigma(F_1^*) \}^{-1} \mu_1(F_1^*) = (\delta - \epsilon X)' \{ (I_p + \epsilon \lambda \Sigma^{-1} \Sigma_1) \Sigma^{-1} - \\ &- \epsilon \lambda \Sigma^{-1} X X' \Sigma^{-1} \} (\mu_1 + \epsilon X) = \delta' (I_p + \epsilon \lambda \Sigma^{-1} \Sigma_1) \Sigma^{-1} \mu_1 - \\ &- \epsilon \lambda \delta' \Sigma^{-1} X X' \Sigma^{-1} \mu_1 - \epsilon X' \{ I_p + \epsilon \lambda \Sigma^{-1} \Sigma_1 \} \Sigma^{-1} \mu_1 + \\ &+ \epsilon^2 \lambda X' \Sigma^{-1} X X' \Sigma^{-1} \mu_1 + \epsilon \delta' \{ I_p + \epsilon \lambda \Sigma^{-1} \Sigma_1 \} \Sigma^{-1} X - \\ &- \epsilon^2 \lambda \delta' \Sigma^{-1} X X' \Sigma^{-1} X - \epsilon^2 X' \{ I_p + \epsilon \lambda \Sigma^{-1} \Sigma_1 \} \Sigma^{-1} X + \\ &+ \epsilon^3 \lambda X' \Sigma^{-1} X X' \Sigma^{-1} X \end{aligned}$$

$$\begin{aligned}
 b_{\mu_2}(\beta_3^*) &= \delta'(\beta_3^*) \{ \Sigma(\beta_3^*) \}^{-1} \mu_2 = (\delta + \varepsilon x)' \{ (I_p + \varepsilon \lambda \Sigma^{-1} \Sigma_1) \Sigma^{-1} - \\
 &- \varepsilon \lambda x x' \Sigma^{-1} \} \mu_2 = \delta' (I_p + \varepsilon \lambda \Sigma^{-1} \Sigma_1) \Sigma^{-1} \mu_2 - \varepsilon \lambda \delta' \Sigma^{-1} x x' \Sigma^{-1} \mu_2 - \\
 &- \varepsilon x' \{ I_p + \varepsilon \lambda \Sigma^{-1} \Sigma_1 \} \Sigma^{-1} \mu_2 + \varepsilon^2 \lambda x' \Sigma^{-1} x x' \Sigma^{-1} \mu_2.
 \end{aligned}$$

y dada la definición de función influencia podemos despreciar, los terminos en ε de orden mayor que uno. Luego las expresiones anteriores nos quedarían.

$$\begin{aligned}
 b_{\mu_1}(\beta_3^*) &= \delta' \Sigma^{-1} \mu_1 + \varepsilon \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_1 - \varepsilon \lambda \delta' \Sigma^{-1} x x' \Sigma^{-1} \mu_1 - \\
 &- \varepsilon x' \Sigma^{-1} \mu_1 + \varepsilon \delta' \Sigma^{-1} x
 \end{aligned}$$

$$\begin{aligned}
 b_{\mu_2}(\beta_3^*) &= \delta' \Sigma^{-1} \mu_2 + \varepsilon \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_2 - \varepsilon \lambda \delta' \Sigma^{-1} x x' \Sigma^{-1} \mu_2 - \\
 &- \varepsilon x' \Sigma^{-1} \mu_2.
 \end{aligned}$$

2.4.2. Función influencia para b_{μ_1} y b_{μ_2} .

Aplicando la definición dada de función influencia resulta

$$CI_{b_{\mu_1}, \beta_3}(\omega) = \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_1 - \lambda \delta' \Sigma^{-1} x x' \Sigma^{-1} \mu_1 - x' \Sigma^{-1} \mu_1 + \delta' \Sigma^{-1} x$$

$$CI_{b_{\mu_2}, \beta_3}(\omega) = \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_2 - \lambda \delta' \Sigma^{-1} x x' \Sigma^{-1} \mu_2 - x' \Sigma^{-1} \mu_2.$$

y notando

$$Y = \delta' \Sigma^{-1} X \quad \text{y} \quad U_i = \mu_i' \Sigma^{-1} X \quad i = 1, 2.$$

se obtiene

$$y \quad CI_{b_{\mu_1}, \beta_3}(\omega) = \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_1 - \lambda Y U_1 - U_1 + Y$$

$$CI_{b_{\mu_2}, \beta_3}(\omega) = \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_2 - \lambda Y U_2 - U_2.$$

2.4.3. Distribución de $CI_{B/\mu_1, \sigma_1}$ y $CI_{B/\mu_2, \sigma_2}$.

En este apartado vamos a calcular la distribución de la función influencia de las medias discriminantes, distribuciones que serán necesarias para la detección de outliers.

Para ello necesitamos de los siguientes lemas:

2.4.3.1. Lema. Sea \mathbf{V} un vector aleatorio bidimensional distribuido según una $N_2(\mu, \Sigma)$ no singular. Entonces la función generatriz de momentos de la variable ξ obtenida al multiplicar las componentes del vector aleatorio \mathbf{V} viene dada por.

$$\psi_{\xi}(\theta) = \left(1 - \frac{2\theta \sigma_1 \sigma_2 (1+\rho)}{2}\right)^{-\frac{1}{2}} \exp \left\{ \frac{\frac{1}{2} \frac{(1-\rho)(\mu_2 \sigma_1 + \mu_1 \sigma_2)^2}{\sigma_1^2 \sigma_2^2 (1-\rho^2)} \theta \frac{\sigma_1 \sigma_2 (1+\rho)}{2}}{1 - 2\theta \frac{\sigma_1 \sigma_2 (1+\rho)}{2}} \right\} \times$$

$$\times \left(1 - 2(-\theta) \frac{\sigma_1 \sigma_2 (1-\rho)}{2}\right)^{-\frac{1}{2}} \exp \left\{ \frac{\frac{1}{2} \frac{(1+\rho)(\mu_2 \sigma_1 - \mu_1 \sigma_2)^2}{\sigma_1^2 \sigma_2^2 (1-\rho^2)} (-\theta) \frac{\sigma_1 \sigma_2 (1-\rho)}{2}}{1 - 2(-\theta) \frac{\sigma_1 \sigma_2 (1-\rho)}{2}} \right\}$$

donde ρ es el coeficiente de correlación entre las componentes del vector .

En efecto :

$$\mathbf{V} = \begin{pmatrix} X \\ Y \end{pmatrix} \in N_2 \left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} ; \begin{pmatrix} \sigma_1^2 & \rho \sigma_1 \sigma_2 \\ \rho \sigma_1 \sigma_2 & \sigma_2^2 \end{pmatrix} \right)$$

Calculemos la función generatriz de momentos de la variable $\xi = X \cdot Y$.

$$\psi_{\xi}(\theta) = E[e^{\theta \xi}] = E[e^{\theta X Y}] = E_Y [E_X [e^{(\theta Y) X} / Y = y]]$$

$$E_X [e^{(\theta Y) X} / Y = y] = \exp \left\{ \theta y \left(\mu_1 + \rho \frac{\sigma_1}{\sigma_2} (y - \mu_2) \right) + \frac{1}{2} \sigma_1^2 (1-\rho^2) \theta^2 y^2 \right\}$$

Luego

$$\Psi_{\frac{1}{2}}(\theta) = E_Y \left[\exp \left\{ \left(\mu_1 + \rho \frac{\sigma_1}{\sigma_2} (Y - \mu_2) \right) \theta Y + \frac{1}{2} \sigma_1^2 (1 - \rho^2) \theta^2 Y^2 \right\} \right]$$

ya que

$$Y \in N_1(\mu_2, \sigma_2)$$

se tiene

$$\Psi_{\frac{1}{2}}(\theta) = \frac{\exp \left\{ -\frac{1}{2} \frac{\mu_2^2}{\sigma_2^2} \right\} \cdot \exp \left\{ \frac{1}{2} \frac{\left(\theta \mu_1 + \frac{\mu_2}{\sigma_2^2} - \rho \frac{\sigma_1}{\sigma_2} \mu_2 \theta \right)^2}{\frac{1}{\sigma_2^2} - \sigma_1^2 (1 - \rho^2) \theta^2 - 2 \rho \frac{\sigma_1}{\sigma_2} \theta} \right\}}{\sigma_2 \sqrt{\frac{1}{\sigma_2^2} - 2 \rho \frac{\sigma_1}{\sigma_2} \theta - \sigma_1^2 (1 - \rho^2) \theta^2}}$$

Realizando las siguientes descomposiciones.

$$1 - 2 \rho \sigma_1 \sigma_2 \theta - \sigma_1^2 \sigma_2^2 (1 - \rho^2) \theta^2 = (1 - \sigma_1 \sigma_2 \theta (1 + \rho)) (1 + \sigma_1 \sigma_2 \theta (1 - \rho))$$

y

$$\frac{1}{2} \frac{\left(\theta \mu_1 + \frac{\mu_2}{\sigma_2^2} - \rho \frac{\sigma_1}{\sigma_2} \mu_2 \theta \right)^2}{\frac{1}{\sigma_2^2} - \sigma_1^2 (1 - \rho^2) \theta^2 - 2 \rho \frac{\sigma_1}{\sigma_2} \theta} = \frac{1}{2} \frac{\mu_2^2}{\sigma_2^2} + \frac{2 \rho \sigma_1 \sigma_2 \mu_1 \mu_2 - \mu_1^2 \sigma_2^2 - \mu_2^2 \sigma_1^2}{2 \sigma_1^2 \sigma_2^2 (1 - \rho^2)} +$$

$$+ \frac{(1 - \rho) (\mu_1 \sigma_2 + \mu_2 \sigma_1)^2}{4 \sigma_1^2 \sigma_2^2 (1 - \rho^2) (1 - \theta \sigma_1 \sigma_2 (1 + \rho))} + \frac{(1 + \rho) (\mu_2 \sigma_1 - \mu_1 \sigma_2)^2}{4 \sigma_1^2 \sigma_2^2 (1 - \rho^2) (1 + \theta \sigma_1 \sigma_2 (1 - \rho))}$$

La función generatriz de momentos que se obtiene es.

$$\Psi_{\frac{1}{2}}(\theta) = \left(1 - \frac{2 \theta \sigma_1 \sigma_2 (1 + \rho)}{2} \right)^{-\frac{1}{2}} \exp \left\{ \frac{\frac{1}{2} \frac{(1 - \rho) (\mu_2 \sigma_1 + \mu_1 \sigma_2)^2}{\sigma_1^2 \sigma_2^2 (1 - \rho^2)} - \frac{\theta \sigma_1 \sigma_2 (1 + \rho)}{2}}{1 - 2 \theta \frac{\sigma_1 \sigma_2 (1 + \rho)}{2}} \right\}$$

$$\times \left(1 - 2(-\theta) \frac{\sigma_1 \sigma_2 (1 - \rho)}{2} \right)^{-\frac{1}{2}} \exp \left\{ \frac{\frac{1}{2} \frac{(1 + \rho) (\mu_2 \sigma_1 - \mu_1 \sigma_2)^2}{\sigma_1^2 \sigma_2^2 (1 - \rho^2)} - \frac{(-\theta) \sigma_1 \sigma_2 (1 - \rho)}{2}}{1 - 2(-\theta) \frac{\sigma_1 \sigma_2 (1 - \rho)}{2}} \right\}$$

2.4.3.2. Lema. Sea X una variable aleatoria distribuida según una chi cuadrado no centrada con un grado de libertad y parámetro de descentralización $\frac{(1 - \rho) (\mu_2 \sigma_1 + \mu_1 \sigma_2)^2}{2 \sigma_1^2 \sigma_2^2 (1 - \rho^2)}$ y sea Y una variable aleatoria independiente de la anterior y distribuida según una chi cuadrado no centrada con un grado de libertad y parámetro de descentralización

$$\frac{(1 + \rho) (\mu_2 \sigma_1 - \mu_1 \sigma_2)^2}{2 \sigma_1^2 \sigma_2^2 (1 - \rho^2)}$$

. Entonces la variable.

$$H = \frac{\sigma_1 \sigma_2 (1+\rho)}{2} X - \frac{\sigma_1 \sigma_2 (1-\rho)}{2} Y$$

tiene como funcion generatriz de momentos.

$$\psi_H(\theta) = \left(1 - \frac{2\theta \sigma_1 \sigma_2 (1+\rho)}{2}\right)^{-1/2} \exp\left\{\frac{\frac{1}{2} \frac{(1-\rho)(\mu_2 \sigma_1 + \mu_1 \sigma_2)^2}{\sigma_1^2 \sigma_2^2 (1-\rho^2)} \theta \cdot \frac{\sigma_1 \sigma_2 (1+\rho)}{2}}{1 - \frac{2\theta \sigma_1 \sigma_2 (1+\rho)}{2}}\right\} \\ \times \left(1 - \frac{2(-\theta) \sigma_1 \sigma_2 (1-\rho)}{2}\right) \exp\left\{\frac{\frac{1}{2} \frac{(1+\rho)(\mu_2 \sigma_1 - \mu_1 \sigma_2)^2}{\sigma_1^2 \sigma_2^2 (1-\rho^2)} (-\theta) \frac{\sigma_1 \sigma_2 (1-\rho)}{2}}{1 - \frac{2(-\theta) \sigma_1 \sigma_2 (1-\rho)}{2}}\right\}$$

La demostracion es trivial, basta utilizar la funcion generatriz de momentos de una distribucion chi - cuadrado no centrada (PATEL 1976)

2.4.3.3. Lema. Sea U una variable aleatoria distribuida segun una $\chi^2(1, \eta_1)$ y sea V una variable aleatoria distribuida segun una $\chi^2(1, \eta_2)$ independiente estocasticamente de U. Sean r, s dos constantes reales positivas arbitrarias. Con estas condiciones la funcion de densidad de la variable aleatoria $T = rU - sV$ viene dada por:

$$h(t) = \sum_{k=0}^{\infty} \sum_{j=0}^{\infty} \frac{e^{-\eta_1/2} \left(\frac{\eta_1}{2}\right)^k}{k!} \frac{e^{-\eta_2/2} \left(\frac{\eta_2}{2}\right)^j}{j!} \int_0^{\infty} \frac{g_{\chi^2_{2k+1}}\left(\frac{t+x}{r}\right)}{r} \cdot \frac{g_{\chi^2_{2j+1}}\left(\frac{x}{s}\right)}{s} dx$$

donde g_{2i+1} es la funcion de densidad de una variable chi cuadrado centrada con $2i+1$ grados de libertad.

Ademas la funcion de densidad de T puede aproximarse por.

$$h(t) = \begin{cases} \left(\frac{s}{r+s}\right)^{1/2} e^{-\eta_2/2 \left(1 - \frac{s}{r+s}\right)} \frac{1}{s} f_{\chi^2(1, \eta_2)}\left(\frac{-t}{s}\right) & t \leq 0 \\ \left(\frac{r}{r+s}\right)^{1/2} e^{-\eta_1/2 \left(1 - \frac{r}{r+s}\right)} \frac{1}{r} f_{\chi^2(1, \eta_1)}\left(\frac{t}{r}\right) & t \geq 0 \end{cases}$$

En efecto:

Para hallar la funcion de densidad de la variable aleatoria T, hemos de calcular la convolucion de las variables rU y sV.

Por ser

$$U \in \chi^2(1, n_1)$$

esta puede expresarse como mixturas de variables distribuidas segun chi cuadrados centradas (JOHNSON, 1970). Asi la funcion de densidad de U se puede expresar

$$f_U(\cdot) = \sum_{H=0}^{\infty} e^{-n_1/2} \frac{(n_1/2)^H}{H!} g_{\chi^2(2H+1)}(\cdot)$$

donde $g_{\chi^2(2H+1)}(\cdot)$ es la funcion de densidad de una variable aleatoria chi - cuadrado centrada con $2H+1$ grados de libertad.

Por tanto $Z = r \cdot U$ tendrá como función de densidad

$$f_{rU}(z) = \sum_{k=0}^{\infty} e^{-n/2} \frac{(n/2)^k}{k!} \frac{1}{r} g_{\chi^2(2k+1)}\left(\frac{z}{r}\right)$$

Analogamente se podrian expresar la funcion de densidad de V y sV.

Luego la funcion de densidad de T vendrá dada por

$$h(t) = \int_0^{\infty} f_{rU}(t+x) f_{sV}(x) dx = \sum_{k=0}^{\infty} \sum_{j=0}^{\infty} \frac{e^{-n/2} (n/2)^k e^{-1/2} (1/2)^j}{k! j!} \int_0^{\infty} \frac{g_{\chi^2(2k+1)}\left(\frac{t+x}{r}\right) g_{\chi^2(2j+1)}\left(\frac{x}{s}\right)}{r \cdot s} dx$$

y la integral que aparece en esta ultima expresion representa la función de densidad de la variable $T' = rU' + sV'$ donde $U' \in \chi^2_{(2k+1)}$ y $V' \in \chi^2_{(2j+1)}$ que viene expresada por

$$f_{2k+1, 2j+1}(t) = \begin{cases} \frac{c(2k+1, 2j+1)}{\Gamma\left(\frac{2k+1}{2}\right)} t^{k+j} e^{-\frac{t}{2r}} U\left(\frac{2j+1}{2}, k+j, \frac{r+s}{2rs} t\right) & t \geq 0 \\ \frac{c(2k+1, 2j+1)}{\Gamma\left(\frac{2j+1}{2}\right)} (-t)^{k+j} e^{\frac{t}{2s}} U\left(\frac{2k+1}{2}, k+j, \left[-\frac{r+s}{2rs} t\right]\right) & t \leq 0 \end{cases}$$

donde $c^{-1}(2k+1, 2j+1) = 2^{k+j+1} \frac{\Gamma(2k+1) \Gamma(2j+1)}{\Gamma\left(\frac{2k+1}{2}\right) \Gamma\left(\frac{2j+1}{2}\right)}$ y $U(a, b, z)$ es la conocida funcion de Kummer (ABRAMOWITZ, 1972).

Y utilizando ahora la aproximacion dada por PRESS (1966), resulta para la funcion $h(t)$ la expresion



$$h(t) = \begin{cases} \sum_{k=0}^{\infty} \sum_{j=0}^{\infty} \frac{e^{-\eta/2} (\eta/2)^k e^{-\eta/2} (\eta/2)^j}{k! j!} \frac{c(2k+1, 2j+1)}{\Gamma(\frac{2j+1}{2})} (-t)^{k+j} e^{\frac{t}{2s}} \left(-\frac{r+s}{2rs} t\right)^{-(k+1/2)} & t < 0 \\ \sum_{k=0}^{\infty} \sum_{j=0}^{\infty} \frac{e^{-\eta/2} (\eta/2)^k e^{-\eta/2} (\eta/2)^j}{k! j!} \frac{c(2k+1, 2j+1)}{\Gamma(\frac{2k+1}{2})} t^{k+j} e^{-\frac{t}{2r}} \left(\frac{r+s}{2rs} t\right)^{-(j+1/2)} & t \geq 0 \end{cases}$$

y mediante sucesivos calculos se obtiene finalmente:

$$h(t) = \begin{cases} \left(\frac{s}{r+s}\right)^{1/2} e^{-\eta/2(1-\frac{s}{r+s})} \frac{1}{s} f_{X^2(1, \eta_2)}\left(\frac{-t}{s}\right) & t < 0 \\ \left(\frac{r}{r+s}\right)^{1/2} e^{-\eta/2(1-\frac{r}{r+s})} \frac{1}{r} f_{X^2(1, \eta_1)}\left(\frac{t}{r}\right) & t \geq 0 \end{cases}$$

Veamos ahora la distribución de las funciones influencias para las medias discriminantes. En el epigrafe 2.4.2 se ha obtenido que

$$CI_{B\mu_1, F_1}(w) = \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_1 - (\lambda Y + 1)(U_1 - \frac{1}{\lambda}) - \frac{1}{\lambda}$$

$$CI_{B\mu_2, F_2}(w) = \lambda \delta' \Sigma^{-1} \Sigma_2 \Sigma^{-1} \mu_2 - (\lambda Y + 1) U_2$$

Ahora bien bajo la hipótesis inicial de que W no es outlier, las variables Y , U_1 y U_2 se distribuyen

$$Y \in N_1(0, \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta)$$

$$U_1 \in N_1(0, \mu_1' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_1)$$

$$U_2 \in N_1(0, \mu_2' \Sigma^{-1} \Sigma_2 \Sigma^{-1} \mu_2)$$

Consideremos las funciones

$$CI'_{B\mu_1, F_1}(w) = \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_1 + \frac{1}{\lambda} - CI_{B\mu_1, F_1}(w) = (\lambda Y + 1)(U_1 - \frac{1}{\lambda})$$

$$CI'_{b'_{\mu_2}, F_1}(w) = \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_2 - CI'_{b'_{\mu_2}, F_1}(w) =$$

$$= (\lambda \gamma + 1) U_2$$

Ya que si toda combinación lineal de las componentes de un vector aleatorio es una variable aleatoria normal, dicho vector también se distribuye normalmente (DUMAS DE RAULY (1.966)), por la definición de las variables

$$(\lambda \gamma + 1) \text{ y } (U_1 - \frac{1}{\lambda})$$

y en virtud de este teorema de caracterización, se tiene

$$\begin{pmatrix} \lambda \gamma + 1 \\ U_1 - \frac{1}{\lambda} \end{pmatrix} \in N_2 \left(\begin{pmatrix} 1 \\ -\frac{1}{\lambda} \end{pmatrix}; \begin{pmatrix} \lambda^2 \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta & \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_2 \\ \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_2 & \mu_2' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_2 \end{pmatrix} \right)$$

luego se verifican las condiciones de los lemas anteriormente demostrados y por tanto la función de densidad de $CI'_{b'_{\mu_2}, F_1}(w)$ será:

$$h(CI') = \begin{cases} \left(\frac{1-\rho}{2} \right)^{\frac{1}{2}} \exp \left\{ -\frac{1}{2} \left[\frac{(\sigma_2 - \sigma_1 \lambda)^2}{2 \sigma_1^2 \sigma_2^2 (1+\rho)} \right] \frac{1+\rho}{2} \right\} \times \frac{2}{\sigma_1 \sigma_2 (1-\rho)} \chi^2 \left(1, \frac{(\sigma_2 + \sigma_1 \lambda)^2}{2 \sigma_1^2 \sigma_2^2 (1+\rho)} \right) \left(\frac{-2 CI'_{b'_{\mu_2}, F_1}(w)}{\sigma_1 \sigma_2 (1-\rho)} \right) & CI'_{b'_{\mu_2}, F_1}(w) < 0 \\ \left(\frac{1+\rho}{2} \right)^{\frac{1}{2}} \exp \left\{ -\frac{1}{2} \left[\frac{(\sigma_1 \lambda + \sigma_2)^2}{2 \sigma_1^2 \sigma_2^2 (1-\rho)} \right] \frac{1-\rho}{2} \right\} \times \frac{2}{\sigma_1 \sigma_2 (1+\rho)} \chi^2 \left(1, \frac{(\sigma_2 - \sigma_1 \lambda)^2}{2 \sigma_1^2 \sigma_2^2 (1+\rho)} \right) \left(\frac{2 CI'_{b'_{\mu_2}, F_1}(w)}{\sigma_1 \sigma_2 (1+\rho)} \right) & CI'_{b'_{\mu_2}, F_1}(w) \geq 0 \end{cases}$$

donde

$$\sigma_1 = \sqrt{\lambda^2 \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta}$$

$$\sigma_2 = \sqrt{\mu_2' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_2}$$

$$\rho = \frac{\lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_2}{\sigma_1 \sigma_2}$$

Análogamente todas las combinaciones lineales de las variables

$$(t \gamma + 1) \text{ y } U_2$$

son normales, luego

$$\begin{pmatrix} Y \\ U_2 \end{pmatrix} \in N_2 \left(\begin{pmatrix} 1 \\ 0 \end{pmatrix}; \begin{pmatrix} \lambda^2 \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta & \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_2 \\ \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_2 & \mu_2' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_2 \end{pmatrix} \right)$$

y por los lemas anteriores, la función de densidad de $CI'_{b\mu_2, F_1}(w)$ es

$$h(CI') = \begin{cases} \left(\frac{1-\rho}{2}\right)^{1/2} \exp\left\{-\frac{1}{2} \left[\frac{1}{2\sigma_1^2(1+\rho)} \right] \frac{1+\rho}{2} \right\} \frac{2}{\sigma_1 \sigma_2 (1-\rho)} f_{\chi^2(1, \frac{1}{2\sigma_1^2(1-\rho)})} \left(\frac{-2 CI'_{b\mu_2, F_1}}{\sigma_1 \sigma_2 (1-\rho)} \right) \\ CI'_{b\mu_2, F_1}(w) < 0 \\ \left(\frac{1+\rho}{2}\right)^{1/2} \exp\left\{-\frac{1}{2} \left[\frac{1}{2\sigma_1^2(1-\rho)} \right] \frac{1-\rho}{2} \right\} \frac{2}{\sigma_1 \sigma_2 (1+\rho)} f_{\chi^2(1, \frac{1}{2\sigma_1^2(1-\rho)})} \left(\frac{2 CI'_{b\mu_2, F_1}}{\sigma_1 \sigma_2 (1+\rho)} \right) \\ CI'_{b\mu_2, F_1}(w) > 0 \end{cases}$$

donde

$$\sigma_1 = \sqrt{\lambda^2 \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta}$$

$$\sigma_2 = \sqrt{\mu_2' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_2}$$

$$\rho = \frac{\lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \mu_2}{\sigma_1 \sigma_2}$$

Y teniendo en cuenta el lema 2.3.3.1. la distribución de estas variables se pueden aproximar

$$CI'_{b\mu_2, F_1}(w) \propto G_1 \left(\frac{\left(1 + \frac{(\sigma_2/\lambda + \sigma_2)^2}{2\sigma_1^2\sigma_2^2(1-\rho)}\right)^2}{2\left(1 + 2\frac{(\sigma_2 + \sigma_1/\lambda)^2}{2\sigma_1^2\sigma_2^2(1+\rho)}\right)} ; \frac{1 + \frac{(\sigma_2/\lambda + \sigma_2)^2}{2\sigma_1^2\sigma_2^2(1-\rho)}}{2\left(1 + 2\frac{(\sigma_1/\lambda + \sigma_2)^2}{2\sigma_1^2\sigma_2^2(1-\rho)}\right) \frac{\sigma_1 \sigma_2 (1-\rho)}{2}} \right)$$

$$CI'_{b\mu_2, F_1}(w) < 0$$

$$CI'_{b\mu_2, F_1}(w) \propto G_1 \left(\frac{\left(1 + \frac{(\sigma_2 - \sigma_1/\lambda)^2}{2\sigma_1^2\sigma_2^2(1+\rho)}\right)^2}{2\left(1 + 2\frac{(\sigma_2 - \sigma_1/\lambda)^2}{2\sigma_1^2\sigma_2^2(1+\rho)}\right)} ; \frac{1 + \frac{(\sigma_2 - \sigma_1/\lambda)^2}{2\sigma_1^2\sigma_2^2(1+\rho)}}{2\left(1 + 2\frac{(\sigma_2 - \sigma_1/\lambda)^2}{2\sigma_1^2\sigma_2^2(1+\rho)}\right) \frac{\sigma_1 \sigma_2 (1+\rho)}{2}} \right)$$

$$CI'_{b\mu_2, F_1}(w) > 0$$

$$CI'_{b^2, F_1}(w) \propto G\left(\frac{\left(1 + \frac{1}{2\sigma_1^2(1-\rho)}\right)^2}{2\left(1 + 2\frac{1}{2\sigma_1^2(1-\rho)}\right)} ; \frac{1 + \frac{1}{2\sigma_1^2(1-\rho)}}{2\left(1 + 2\frac{1}{2\sigma_1^2(1-\rho)}\right)\frac{\sigma_1\sigma_2(1-\rho)}{2}}\right)$$

$$CI'_{b^2, F_1}(w) < 0$$

$$CI'_{b^2, F_1}(w) \propto G\left(\frac{\left(1 + \frac{1}{2\sigma_1^2(1+\rho)}\right)^2}{2\left(1 + 2\frac{1}{2\sigma_1^2(1+\rho)}\right)} ; \frac{1 + \frac{1}{2\sigma_1^2(1+\rho)}}{2\left(1 + 2\frac{1}{2\sigma_1^2(1+\rho)}\right)\frac{\sigma_1\sigma_2(1+\rho)}{2}}\right)$$

$$CI'_{b^2, F_1}(w) \geq 0$$

2.5. FUNCION INFLUENCIA PARA LOS COEFICIENTES DE LA FUNCION DISCRIMINANTE .

Al igual que en los apartados anteriores se estudia la funcion influencia para el producto $\hat{b}\hat{b}$, o modulo del vector \hat{b} y a continuacion determinamos su distribucion.

2.5.1. Calculo de $\hat{b}\hat{b}(F_1^*)$.

Se tiene que:

$$\begin{aligned} \hat{b}\hat{b}(F_1^*) &= \delta'(F_1^*) \{ \Sigma(F_1^*) \}^{-1} \{ \Sigma(F_1^*) \}^{-1} \delta(F_1^*) = \delta' \Sigma^{-1} \Sigma^{-1} \delta + \\ &+ \varepsilon \lambda \delta' \Sigma^{-1} \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta - \varepsilon t \delta' \Sigma^{-1} \Sigma^{-1} X X' \Sigma^{-1} \delta - \varepsilon \delta' \Sigma^{-1} \Sigma^{-1} X - \\ &- \varepsilon^2 \lambda \delta' \Sigma^{-1} \Sigma^{-1} \Sigma_1 \Sigma^{-1} X + \varepsilon^2 \lambda \delta' \Sigma^{-1} \Sigma^{-1} X X' \Sigma^{-1} X + \\ &+ \varepsilon \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \Sigma^{-1} \delta + \varepsilon^2 \lambda^2 \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta - \\ &- \varepsilon^2 \lambda^2 \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \Sigma^{-1} X X' \Sigma^{-1} \delta - \varepsilon^2 \lambda \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \Sigma^{-1} X - \\ &- \varepsilon^3 \lambda^2 \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \Sigma^{-1} \Sigma_1 \Sigma^{-1} X + \varepsilon^3 \lambda^2 \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \Sigma^{-1} X X' \Sigma^{-1} X - \\ &- \varepsilon \lambda \delta' \Sigma^{-1} X X' \Sigma^{-1} \Sigma^{-1} \delta - \varepsilon^2 \lambda^2 \delta' \Sigma^{-1} X X' \Sigma^{-1} \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta + \\ &+ \varepsilon^2 \lambda^2 \delta' \Sigma^{-1} X X' \Sigma^{-1} \Sigma^{-1} X X' \Sigma^{-1} \delta + \varepsilon^3 t \delta' \Sigma^{-1} X X' \Sigma^{-1} \Sigma^{-1} X + \end{aligned}$$

$$\begin{aligned}
& + \varepsilon^2 \lambda^2 \delta' \Sigma^{-1} X X' \Sigma^{-1} \Sigma^{-1} \Sigma^{-1} X - \varepsilon^3 \lambda^2 \delta' \Sigma^{-1} X X' \Sigma^{-1} \Sigma^{-1} X X' \Sigma^{-1} X - \\
& - \varepsilon X' \Sigma^{-1} \Sigma^{-1} \delta - \varepsilon^2 \lambda X' \Sigma^{-1} \Sigma^{-1} \Sigma^{-1} \delta + \varepsilon^2 \lambda X' \Sigma^{-1} \Sigma^{-1} X X' \Sigma^{-1} \delta \\
& + \varepsilon^2 X' \Sigma^{-1} \Sigma^{-1} X + \varepsilon^3 \lambda X' \Sigma^{-1} \Sigma^{-1} \Sigma^{-1} X - \varepsilon^3 \lambda X' \Sigma^{-1} \Sigma^{-1} X X' \Sigma^{-1} X - \\
& - \varepsilon^2 \lambda X' \Sigma^{-1} \Sigma^{-1} \Sigma^{-1} \delta - \varepsilon^3 \lambda^2 X' \Sigma^{-1} \Sigma^{-1} \Sigma^{-1} \Sigma^{-1} \delta + \\
& + \varepsilon^3 \lambda^2 X' \Sigma^{-1} \Sigma^{-1} \Sigma^{-1} X X' \Sigma^{-1} \delta + \varepsilon^3 \lambda X' \Sigma^{-1} \Sigma^{-1} \Sigma^{-1} X + \\
& + \varepsilon^4 \lambda^2 X' \Sigma^{-1} \Sigma^{-1} \Sigma^{-1} \Sigma^{-1} X - \varepsilon^4 \lambda^2 X' \Sigma^{-1} \Sigma^{-1} \Sigma^{-1} X X' \\
& \cdot \Sigma^{-1} X + \varepsilon^2 \lambda X' \Sigma^{-1} X X' \Sigma^{-1} \Sigma^{-1} \delta + \varepsilon^3 \lambda^2 X' \Sigma^{-1} X X' \Sigma^{-1} \Sigma^{-1} \\
& \Sigma^{-1} \delta - \varepsilon^3 \lambda^2 X' \Sigma^{-1} X X' \Sigma^{-1} \Sigma^{-1} X X' \Sigma^{-1} \delta - \varepsilon^3 \lambda X' \Sigma^{-1} X X' \\
& \Sigma^{-1} \Sigma^{-1} X - \varepsilon^4 \lambda^2 X' \Sigma^{-1} X X' \Sigma^{-1} \Sigma^{-1} \Sigma^{-1} X + \\
& + \varepsilon^4 \lambda^2 X' \Sigma^{-1} X X' \Sigma^{-1} \Sigma^{-1} X X' \Sigma^{-1} X
\end{aligned}$$

y despreciando las potencias en ε de orden mayor que uno, ya que estas no intervienen en la función influencia. Se obtiene

$$\begin{aligned}
b'b(\beta_1^*) &= \delta' \Sigma^{-1} \Sigma^{-1} \delta + \varepsilon \lambda \delta' \Sigma^{-1} \Sigma^{-1} \Sigma^{-1} \delta - \varepsilon \lambda \delta' \Sigma^{-1} \Sigma^{-1} X X' \Sigma^{-1} \delta - \\
&- \varepsilon \delta' \Sigma^{-1} \Sigma^{-1} X + \varepsilon \lambda \delta' \Sigma^{-1} \Sigma^{-1} \Sigma^{-1} \delta - \varepsilon \lambda \delta' \Sigma^{-1} X X' \Sigma^{-1} \Sigma^{-1} \delta - \\
&- \varepsilon X' \Sigma^{-1} \Sigma^{-1} \delta
\end{aligned}$$

2.5.2.- Función influencia para $b'b$.

$$\begin{aligned}
CI_{b'b, \beta_1}(\omega) &= \lim_{\varepsilon \rightarrow 0} \left\{ \frac{b'b(\beta_1^*) - b'b(\beta_1)}{\varepsilon} \right\} = 2 \left[\lambda \delta' \Sigma^{-1} \Sigma^{-1} \Sigma^{-1} \delta - \right. \\
&\left. - \lambda \delta' \Sigma^{-1} \Sigma^{-1} X X' \Sigma^{-1} \delta - \delta' \Sigma^{-1} \Sigma^{-1} X \right]
\end{aligned}$$

Y si notamos

$$Y = \delta' \Sigma^{-1} X = X' b \quad (1) \quad \text{y} \quad \eta = b' \Sigma^{-1} X \quad (2)$$

obtenemos que la función influencia se puede representar por

$$CI_{\delta b, F_1}(w) = 2 \left[\lambda b' \Sigma^{-1} \Sigma_1 b - (\lambda Y + 1) \xi \right]$$

2.5.3. Distribución de $CI_{\delta b, F_1}(w)$.

Bajo la hipótesis de que W no es outlier, las variables Y y ξ definidas en (1) y (2) se distribuyen

$$Y \in N_1(0, \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta)$$

$$N_1(0, b' \Sigma^{-1} \Sigma_1 \Sigma^{-1} b)$$

Consideremos la función

$$CI'_{\delta b, F_1}(w) = (\lambda Y + 1) \xi = \lambda b' \Sigma^{-1} \Sigma_1 b - \frac{CI_{\delta b, F_1}(w)}{2}$$

Ya que cualquier combinación lineal de las variables $(\lambda Y + 1)$ y ξ es normal univariante, en virtud del teorema de caracterización de la ley normal se verifica que

$$\begin{pmatrix} \lambda Y + 1 \\ \xi \end{pmatrix} \in N_2 \left(\begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} \lambda^2 \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta & \lambda b' \Sigma^{-1} \Sigma_1 b \\ \lambda b' \Sigma^{-1} \Sigma_1 b & b' \Sigma^{-1} \Sigma_1 \Sigma^{-1} b \end{pmatrix} \right)$$

Estamos así en las condiciones de los lemas 2.4.3.1, 2.4.3.2 y 2.4.3.3 y por el lema 2.3.3.1, se tendrá que la variable se distribuye

$$CI'_{\delta b, F_1}(w) \propto G \left(\frac{\left(1 + \frac{1}{2\sigma_1^2(1-\rho)}\right)^2}{2\left(1 + 2\frac{1}{2\sigma_1^2(1-\rho)}\right)} ; \frac{1 + \frac{1}{2\sigma_1^2(1-\rho)}}{2\left(1 + 2\frac{1}{2\sigma_1^2(1-\rho)}\right) \frac{\sigma_1 \sigma_2 (1-\rho)}{2}} \right)$$

$$CI'_{\delta b, F_1}(w) < 0$$

$$CI'_{\delta b, F_1}(w) \propto G \left(\frac{\left(1 + \frac{1}{2\sigma_1^2(1+\rho)}\right)^2}{2\left(1 + 2\frac{1}{2\sigma_1^2(1+\rho)}\right)} ; \frac{1 + \frac{1}{2\sigma_1^2(1+\rho)}}{2\left(1 + 2\frac{1}{2\sigma_1^2(1+\rho)}\right) \frac{\sigma_1 \sigma_2 (1+\rho)}{2}} \right)$$

donde

$$\sigma_1 = \sqrt{\lambda^2 \delta' \Sigma^{-1} \Sigma_1 \Sigma^{-1} \delta}$$

$$\sigma_2 = \sqrt{b' \Sigma^{-1} \Sigma_1 \Sigma^{-1} b}$$

$$\rho = \frac{\lambda b' \Sigma^{-1} \Sigma_1 b}{\sigma_1 \sigma_2}$$



2.6. DETECCION DE OUTLIERS.

El estudio realizado anteriormente, está hecho suponiendo conocidos los parámetros poblacionales, sin embargo en las aplicaciones prácticas no vamos a conocer dichos parámetros, por lo que tendremos que recurrir a estimaciones de los mismos, mediante la muestra que obtengamos. Las estimaciones generalmente las haremos por el método de máxima verosimilitud, por aquello de que es el método de estimación con mejores propiedades asintóticas.

Por tanto si vamos a utilizar estimaciones de los parámetros poblacionales, tendremos que calcular estimaciones de la función influencia y en virtud del teorema de Zehna (ROHATGI, 1.976), sabemos que la estimación de máxima verosimilitud de la función influencia, será la función influencia de los estimadores de máxima verosimilitud de los parámetros poblacionales.

Con esto tendremos determinada la función influencia para los parámetros muestrales. La distribución de la función influencia muestral se distribuirá asintóticamente como la de la función influencia poblacional, de la que conocemos su distribución y parámetros.

Hemos determinado la distribución de la función influencia poblacional bajo la hipótesis de que la variable $W \in N_p(\mu_1, \Sigma_1)$ es decir bajo la hipótesis de que la variable no es outliers. Por esta razón tendremos que la distribución de la función influencia muestral estará calculada bajo la hipótesis de que las n observaciones muestrales obtenidas no son outliers, debiendo exigirse que el tamaño de la muestra sea suficientemente grande pues estamos trabajando con procedimientos asintóticos.

Señalemos que el calculo de la función influencia, se ha realizado suponiendo que existe algún punto W en el cual hay definida una distribución degenerada. Se han extraído n observaciones por lo que tendremos una muestra de tamaño n de la variable función influencia, obtenida al ir suponiendo cada una de las observaciones como aquel punto en el que existe dicha distribución degenerada. Luego si se cumple la hipótesis nula de que las n observaciones pertenecen a una $N_p(\mu_1, \Sigma_1)$ tendremos que dichas observaciones se ajustarán a la distribución muestral de la función influencia.

Como se ha comprobado en los apartados anteriores la distribución de las funciones influencia son de tipo gamma. Por lo que proponemos para cada función un método gráfico para la detección de outliers basados en las representaciones gráficas de distribuciones gammas obtenidas mediante procedimientos computacionales por WILK, GNANADESIKAN y HUYETT (1.962)

CAPITULO TERCERO

DETECCION DE OUTLIERS MEDIANTE EL

COCIENTE DE VEROSIMILITUDES

CONTENIDO.

3.1 - Introducción

3.2 - Contraste de cociente de verosimilitudes

3.3 - Detección de outliers

3.4 - Aplicación

3.1. INTRODUCCION

En el capitulo anterior hemos dado un metodo para la deteccion de outliers en poblaciones normales, cuando se supone que estas proceden de una poblacion con funcion de distribucion dada por $(1-\varepsilon)F(x) + \varepsilon \delta(x)$, y además para una técnica determinada como es el analisis discriminante con vectores medias y matrices de covarianza distintas, para funciones discriminantes lineales.

En este capitulo lo que estudiamos es un criterio para detectar outliers cuando se ha extraido una muestra procedente de una poblacion $N_p(\mu, \Sigma)$ no singular.

Para ello se supone que dada una muestra de tamaño n siempre es posible afirmar que en dicha muestra existen k elementos con $k \geq \left[\frac{n}{2} \right]$, que no son outliers es decir que son observaciones de la poblacion que estamos muestreando, pues en caso contrario la muestra no seria representativa.

El criterio que se propone para la deteccion de outliers es el siguiente: Mediante un muestreo aleatorio simple extraemos k observaciones de entre las n de la muestra y sobre estas observaciones hacemos la suposición de que todas son de la poblacion $N_p(\mu, \Sigma)$ y con cada una de las $n - k$ observaciones restantes, realizamos un contraste de hipotesis para comprobar si cada una de estas observaciones es o no outlier.

Este contraste obtenido por el procedimiento de razon de verosimilitud, se basa en un estadistico cuya distribucion es una ley F de Snedecor.

Al final del capitulo se incluye un programa de ordenador, para la deteccion de outliers, mediante el metodo dado en este capitulo.

NOTA: Con $\left[\frac{n}{2} \right]$ indicamos parte entera de $\frac{n}{2}$.

3.2, CONTRASTE DE COCIENTE DE VEROSIMILITUDES.

Sea X_1, X_2, \dots, X_k, Y una muestra aleatoria simple procedente de una población normal p -dimensional no singular. Supongamos que X_1, X_2, \dots, X_k es una muestra de tamaño k extraída de una población $N_p(\mu, \Sigma)$ no singular.

Se quiere contrastar la hipótesis de que la observación Y se ha extraído de una población $N_p(\mu, \Sigma)$ no singular frente a la hipótesis alternativa de que Y pertenece a una población $N_p(\mu_1, a^2 \Sigma)$ con a ($a \neq 1$) una constante real no nula y fijada de antemano.

Utilizaremos para contrastar esta hipótesis el procedimiento del cociente de verosimilitudes.

Bajo la hipótesis nula, la distribución conjunta de la muestra aleatoria de tamaño $k+1$, es.

$$f(x_1, x_2, \dots, x_k, y) = \left[\prod_{i=1}^k (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp\left\{-\frac{1}{2}(x_i - \mu)' \Sigma^{-1}(x_i - \mu)\right\} \right] \cdot (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp\left\{-\frac{1}{2}(y - \mu)' \Sigma^{-1}(y - \mu)\right\}$$

Donde $\{X_i\}_{i=1,2,\dots,k}, Y$ y μ son vectores cuyas componentes X_{ji}, Y_j y μ_j con $j = 1 \leq j \leq p$ toman sus valores en el intervalo $(-\infty, \infty)$ y por último Σ es una matriz simétrica definida positiva, que representaremos por $\Sigma > 0$.

Esta distribución conjunta también se puede expresar de la forma (KSHIRSAGAR, 1.972).

$$f(x_1, x_2, \dots, x_k, y) = (2\pi)^{-\frac{(k+1)p}{2}} |\Sigma|^{-\frac{k+1}{2}} \cdot \exp\left\{-\frac{1}{2} \text{tr.} \Sigma^{-1} \left[\{(X|Y) - \mu E_{1,k+1}\} \{(X|Y) - \mu E_{1,k+1}\}' \right] \right\} \quad (1)$$

donde

$$(X|Y) = \begin{pmatrix} X_{11} & X_{12} & \dots & X_{1k} & | & Y_1 \\ X_{21} & X_{22} & \dots & X_{2k} & | & Y_2 \\ \dots & \dots & \dots & \dots & | & \dots \\ X_{p1} & X_{p2} & \dots & X_{pk} & | & Y_p \end{pmatrix} \quad \text{y} \quad E_{1,k+1} = (1, 1, \dots, 1)$$



Y los parámetros μ y Σ toman sus valores en el espacio paramétrico Ω ,

$$\Omega = \{ \mu, \Sigma \mid -\infty < \mu < \infty, \Sigma > 0 \}$$

siendo $p + \frac{p(p+1)}{2}$ el número de parámetros desconocidos y que estimaremos mediante el método de máxima verosimilitud.

De (1) se deduce que el logaritmo de la función de verosimilitud es

$$\ln L(\Omega) = \ln \int (x_1, x_2, \dots, x_k, y) = -\frac{(k+1)p}{2} \ln 2\pi - \frac{k+1}{2} \ln |\Sigma| - \frac{1}{2} \text{tr.} \Sigma^{-1} S_X - \frac{k}{2} \text{tr.} \Sigma^{-1} (\bar{X} - \mu)(\bar{X} - \mu)' - \frac{1}{2} \text{tr.} \Sigma^{-1} (y - \mu)(y - \mu)' \quad (2)$$

donde

$$\bar{X} = \frac{1}{k} X' E_{k,1}$$

y S_X es la matriz de suma de cuadrados y suma de productos de la matriz de datos X y viene dada por

$$S_X = X(I - \frac{1}{k} E_{kk})X'$$

Las ecuaciones de verosimilitud vendrían dadas por.

$$\frac{\partial \ln L}{\partial \mu} = 0 \quad (3)$$

$$\frac{\partial \ln L}{\partial \sigma_{ij}} = 0, \quad i=1,2,\dots,p; j \geq i \quad (4)$$

donde σ_{ij} es el elemento que está en la i -ésima fila y j -ésima columna de la matriz Σ .

Escribiendo $\hat{\mu}$ por μ y $\hat{\Sigma}$ por Σ en estas ecuaciones para distinguir los estimadores de los parámetros originales, de (3) se obtiene

$$\frac{\partial \ln L}{\partial \mu} = \mu \hat{\Sigma}^{-1} (\bar{X} - \hat{\mu}) + \hat{\Sigma}^{-1} (y - \hat{\mu}) = 0 \quad (5)$$

y de (4) se obtiene

$$\begin{aligned} \frac{\partial \ln L}{\partial \sigma_{ij}} &= -\frac{k+1}{2} (2 - \delta_{ij}) \hat{\sigma}^{ij} + \frac{1}{2} \text{tr.} \frac{\partial \hat{\Sigma}}{\partial \hat{\sigma}_{ij}} \hat{\Sigma}^{-1} S_X \hat{\Sigma}^{-1} + \\ &+ \frac{n}{2} \text{tr.} \frac{\partial \hat{\Sigma}}{\partial \hat{\sigma}_{ij}} \hat{\Sigma}^{-1} (\bar{X} - \hat{\mu})(\bar{X} - \hat{\mu})' \hat{\Sigma}^{-1} + \frac{1}{2} \text{tr.} \frac{\partial \hat{\Sigma}}{\partial \hat{\sigma}_{ij}} (y - \hat{\mu})(y - \hat{\mu})' \hat{\Sigma}^{-1} \\ &= 0 \quad (6) \end{aligned}$$

habiéndose utilizado los siguientes resultados de análisis matricial (BELLMAN, 1.965)

$$\frac{\partial \log |\Sigma|}{\partial \sigma_{ij}} = (2 - \delta_{ij}) \sigma^{-ij} \quad i, j = 1, 2, \dots, p$$

donde σ^{-ij} es el elemento que está en la i -ésima fila y j -ésima columna de Σ^{-1} y δ_{ij} es la delta de Kronecker. Y

$$\frac{\partial \Sigma^{-1}}{\partial \sigma_{ij}} = -\Sigma^{-1} \frac{\partial \Sigma}{\partial \sigma_{ij}} \Sigma^{-1}$$

Multiplicando a la izquierda por $\hat{\Sigma}$, de la expresión (5) se deduce.

$$\hat{\mu} = \frac{k \bar{x} + y}{k+1} \quad (7)$$

que no es más que el vector media muestral de las $k+1$ observaciones.

Y despejando en (6) $\hat{\sigma}^{-ij}$ se obtiene

$$\hat{\sigma}^{-ij} = \frac{1}{(k+1)(2 - \delta_{ij})} \left[\frac{1}{2} + r \frac{\partial \hat{\Sigma}}{\partial \hat{\sigma}_{ij}} \hat{\Sigma}^{-1} \left[S_X + k(\bar{x} - \hat{\mu})(\bar{x} - \hat{\mu})' + (y - \hat{\mu})(y - \hat{\mu})' \right] \hat{\Sigma}^{-1} \right]$$

luego

$$\hat{\Sigma}^{-1} = \frac{1}{k+1} \hat{\Sigma}^{-1} \left[S_X + k(\bar{x} - \hat{\mu})(\bar{x} - \hat{\mu})' + (y - \hat{\mu})(y - \hat{\mu})' \right] \hat{\Sigma}^{-1}$$

y multiplicando a izquierda y derecha por $\hat{\Sigma}$,

$$\hat{\Sigma} = \frac{1}{k+1} \left[S_X + k(\bar{x} - \hat{\mu})(\bar{x} - \hat{\mu})' + (y - \hat{\mu})(y - \hat{\mu})' \right]$$

y sustituyendo $\hat{\mu}$ por su valor dado en (7) se obtiene finalmente

$$\hat{\Sigma} = \frac{1}{k+1} S_{(k+1)} \quad (8)$$

donde

$$S_{(k+1)} = (X|Y) \left(I - \frac{1}{k+1} E_{k+1, k+1} \right) (X|Y)'$$

que es la matriz de suma de cuadrados y suma de productos de las $k+1$ observaciones p -dimensionales X_1, X_2, \dots, X_k, Y .

Para comprobar que estos estimadores maximizan la función de verosimilitud, bastaría aplicar el procedimiento de WATSON (1.964).

Bajo la hipótesis alternativa la función de densidad conjunta de la muestra aleatoria viene dada por (GIRI , 1.977).

$$f(x_1, x_2, \dots, x_k, y) = \left[\prod_{i=1}^k (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp\left\{-\frac{1}{2}(x_i - \mu)' \Sigma^{-1}(x_i - \mu)\right\} \right] \times (2\pi)^{-p/2} |a^2 \Sigma|^{-1/2} \exp\left\{-\frac{1}{2}(y - \mu_1)' a^2 \Sigma^{-1}(y - \mu_1)\right\} \quad (9)$$

y los parámetros μ , μ_1 y Σ toman sus valores en el espacio paramétrico Ω_1 .

$$\Omega_1 = \left\{ \mu, \mu_1, \Sigma \mid -\infty < \mu < \infty, -\infty < \mu_1 < \infty, \Sigma > 0 \right\}$$

siendo $2p + \frac{p(p+1)}{2}$ el número de parámetros desconocidos y que estimaremos mediante el método de máxima verosimilitud.

Tomando logaritmos en (9) se deduce.

$$\ln h(\Omega_1) = \ln_{\Omega_1} f(x_1, \dots, x_k, y) = -\frac{(k+1)p}{2} \ln 2\pi - \frac{k+1}{2} \ln |\Sigma| - \frac{1}{2} \text{tr.} \Sigma^{-1} S_x - \frac{k}{2} \text{tr.} \Sigma^{-1} (\bar{x} - \mu)(\bar{x} - \mu)' - \frac{1}{2} \ln a^{2p} - \frac{1}{2} \text{tr.} a^2 \Sigma^{-1} (y - \mu_1)(y - \mu_1)' \quad (10)$$

y las ecuaciones de verosimilitud vendrían dadas por el sistema

$$\frac{\partial \ln h}{\partial \mu} = k \hat{\Sigma}^{-1} (\bar{x} - \hat{\mu}) = 0 \quad (11)$$

$$\frac{\partial \ln h}{\partial \mu_1} = a^2 \hat{\Sigma}^{-1} (y - \hat{\mu}_1) = 0 \quad (12)$$

$$\frac{\partial \ln h}{\partial \hat{\sigma}_{ij}} = -\frac{k+1}{2} \frac{\partial}{\partial \hat{\sigma}_{ij}} \ln |\hat{\Sigma}| - \frac{1}{2} \text{tr.} \frac{\partial \hat{\Sigma}^{-1}}{\partial \hat{\sigma}_{ij}} S_x - \frac{k}{2} \text{tr.} \frac{\partial \hat{\Sigma}^{-1}}{\partial \hat{\sigma}_{ij}} (\bar{x} - \hat{\mu})(\bar{x} - \hat{\mu})' - \frac{1}{2} \text{tr.} a^2 \frac{\partial \hat{\Sigma}^{-1}}{\partial \hat{\sigma}_{ij}} (y - \hat{\mu}_1)(y - \hat{\mu}_1)' = 0 \quad (13)$$

De (11) multiplicando por $\hat{\Sigma}$ a la izquierda se tiene.

$$\hat{\mu} = \bar{x} \quad (14)$$

De (12) multiplicando por $a^2 \hat{\Sigma}$ a la izquierda se tiene

$$\hat{\mu}_1 = y \quad (15)$$

y sustituyendo en (13) obtenemos

$$-\frac{k+1}{2} \frac{\partial}{\partial \hat{\sigma}_{ij}} \ln |\hat{\Sigma}| - \frac{1}{2} \text{tr} \frac{\partial \hat{\Sigma}^{-1}}{\partial \hat{\sigma}_{ij}} S_X = 0$$

y por ultimo

$$\hat{\Sigma} = \frac{1}{k+1} S_X \quad (16)$$

Veamos ahora que estos estimadores maximizan la función de verosimilitud

$$\begin{aligned} \ln L(\hat{\Omega}_1) - \ln L(\Omega_1) &= -\frac{k+1}{2} \ln \left| \frac{1}{k+1} S_X \right| - \frac{(k+1)p}{2} + \frac{k+1}{2} \ln |\Sigma| + \\ &+ \frac{1}{2} \text{tr} \cdot \Sigma^{-1} S_X + \frac{k}{2} \text{tr} \cdot \Sigma^{-1} (\bar{X} - \mu)(\bar{X} - \mu)' + \frac{1}{2} \text{tr} \cdot \Sigma^{-1} (Y - \mu_1)(Y - \mu_1)' \end{aligned}$$

y como

$$\text{tr} \cdot \Sigma^{-1} (\bar{X} - \mu)(\bar{X} - \mu)' = (\bar{X} - \mu)' \Sigma^{-1} (\bar{X} - \mu)$$

y

$$\text{tr} \cdot \Sigma^{-1} (Y - \mu_1)(Y - \mu_1)' = (Y - \mu_1)' \Sigma^{-1} (Y - \mu_1)$$

son definidas positivas, se tiene

$$\begin{aligned} \ln L(\hat{\Omega}_1) - \ln L(\Omega_1) &\geq -\frac{k+1}{2} \ln \left| \frac{1}{k+1} S_X \right| - \frac{(k+1)p}{2} + \frac{k+1}{2} \ln |\Sigma| + \\ &+ \frac{1}{2} \text{tr} \cdot \Sigma^{-1} S_X = -\frac{k+1}{2} \ln \left| \frac{1}{k+1} S_X \right| - \frac{(k+1)p}{2} + \frac{1}{2} \text{tr} \cdot \Sigma^{-1} S_X = \\ &= -\frac{k+1}{2} \ln \left| \frac{1}{k+1} S_X \Sigma^{-1} \right| - \frac{(k+1)p}{2} + \frac{1}{2} \text{tr} \cdot \Sigma^{-1} S_X \quad (17) \end{aligned}$$

Notemos por $\lambda_1, \lambda_2, \dots, \lambda_p$ los autovalores de la matriz $\frac{1}{k+1} S_X \Sigma^{-1}$ entonces,

$$\left| \frac{1}{k+1} S_X \Sigma^{-1} \right| = \prod_{i=1}^p \lambda_i \quad \text{y} \quad \frac{1}{k+1} \text{tr} \cdot \Sigma^{-1} S_X = \sum_{i=1}^p \lambda_i$$

y sustituyendo en (17)

$$\ln L(\hat{\Omega}_1) - \ln L(\Omega_1) \geq -\frac{k+1}{2} \sum_{i=1}^p \ln \lambda_i - \frac{(k+1)p}{2} + \frac{1}{2} (k+1) \sum_{i=1}^p \lambda_i =$$

$$= \frac{k+1}{2} \left[\sum_{i=1}^p (-\ln \lambda_i - 1 + \lambda_i) \right] \quad (18)$$

ya que para cualquier t real no negativo se tiene

$$t - \ln t \geq 1$$

de (18) deducimos

$$\ln h(\hat{\Omega}_j) - \ln h(\Omega_j) \geq 0$$

luego

$$\ln h(\hat{\Omega}_j) \geq \ln h(\Omega_j)$$

y esto es cierto para cualquier punto $(\hat{\mu}, \hat{\mu}_j, \hat{\Sigma}) \in \Omega_j$.

Sustituyendo ahora los valores obtenidos para $\hat{\mu}, \hat{\mu}_j, \hat{\Sigma}$, en (7) y (8) en la expresión (2) tendremos.

$$\begin{aligned} \max_{\Omega} \ln h(\Omega) &= -\frac{(k+1)p}{2} \ln 2\pi - \frac{k+1}{2} \ln \left| \frac{1}{k+1} S_{(k+1)} \right| - \frac{1}{2} \text{tr}[(k+1) S_{(k+1)}^{-1} S_{(k+1)}] \\ &= -\frac{(k+1)p}{2} \ln 2\pi - \frac{(k+1)}{2} \ln \left| \frac{1}{k+1} S_{(k+1)} \right| - \frac{(k+1)p}{2} \end{aligned} \quad (19)$$

Analogamente sustituyendo los valores obtenidos para $\hat{\mu}, \hat{\mu}_j$, y $\hat{\Sigma}$ en (14), (15), y (16) en la expresión (10), obtenemos.

$$\max_{\Omega_j} \ln h(\Omega_j) = -\frac{(k+1)p}{2} \ln 2\pi - \frac{k+1}{2} \ln \left| \frac{1}{k+1} S_X \right| - \frac{(k+1)p}{2} - \frac{1}{2} \ln a^{2p} \quad (20)$$

El criterio de razón de verosimilitud para contrastar la hipótesis nula es.

$$\lambda = \frac{\max_{\Omega} h(\Omega)}{\max_{\Omega_j} h(\Omega_j)}$$

de (19) y (20) se deduce

$$\lambda = \frac{|S_X|^{\frac{k+1}{2}} a^{p/2}}{|S_{(k+1)}|^{\frac{k+1}{2}}}$$

por lo que

$$\lambda_j = \frac{|S_X|}{|S_{(k+1)}|}$$

Ahora bien

$$S_{(k+1)} = S_X + \frac{k}{k+1} (\bar{x} - y)(\bar{x} - y)'$$

por tanto

$$\lambda_1 = \frac{|S_k|}{|S_X + \frac{k}{k+1} (\bar{x} - y)(\bar{x} - y)'|}$$

y basandonos en la igualdad (AITKEN, 1.965)

$$|A + qq'| = |A| (1 + q'A^{-1}q)$$

$$\lambda_1 = \frac{1}{1 + \frac{k}{k+1} (y - \bar{x})' S_X^{-1} (y - \bar{x})}$$

Por lo que la región crítica para este contraste viene dada por:

$$\frac{1}{1 + \frac{k}{k+1} (y - \bar{x})' S_X^{-1} (y - \bar{x})} \leq \lambda_0$$

que también puede expresarse en la forma

$$\frac{1 - \lambda_0}{\lambda_0} \leq \frac{k}{k+1} (y - \bar{x})' S_X^{-1} (y - \bar{x})$$

Calculemos ahora la distribución del estadístico

$$\frac{k}{k+1} (y - \bar{x})' S_X^{-1} (y - \bar{x})$$

bajo la hipótesis nula de que X_1, X_2, \dots, X_k, Y es una muestra aleatoria simple de tamaño $k+1$ extraída de una población que se distribuye según una ley $N_p(\mu, \Sigma)$ no singular.

Luego

$$Y \in N_p(\mu, \Sigma), \quad X_i \in N_p(\mu, \Sigma) \quad \forall i \quad i=1, 2, \dots, k$$

y por tanto (ANDERSON, 1.958)

$$\bar{X} \in N_p\left(\mu, \frac{1}{k} \Sigma\right)$$

de donde

$$Y - \bar{X} \in N_p\left(0, \frac{k+1}{k} \Sigma\right)$$

y

$$\sqrt{\frac{k}{k+1}} (Y - \bar{X}) \in N_p(0, \Sigma)$$

Por otro lado para $K > p$ (KSHIRSAGAR, 1.972)

$$S_X \in W_p(k-1, \Sigma)$$

y como $(y - \bar{X})$ y S_X son independientes (WILKS, 1.962), el estadístico

$$\frac{k}{k+1} (y - \bar{X})' S_X^{-1} (y - \bar{X}) \equiv \frac{T_p^2}{k-1}$$

donde T_p^2 es una variable que sigue una distribución de Hotelling centrada p - dimensional con $k - 1$ grados de libertad.

Por otro lado según (KSHIRSAGAR, 1.972)

" Si $U \in N_p(0, \Sigma)$ y D es independiente de U y distribuida según $W_p(f, \Sigma)$ la distribución de la variable

$$T_p^2 = U' D^{-1} U$$

es una distribución de Hotelling centrada p - dimensional con f grados de libertad y

$$\frac{f-p+1}{p} \cdot \frac{T_p^2}{f}$$

se distribuye según una ley F de Snedecor con p y $f - p + 1$ grados de libertad "

Por tanto el estadístico

$$\frac{k-p}{p} \cdot \frac{k}{k+1} (y - \bar{X})' S_X^{-1} (y - \bar{X})$$

se distribuye según una F de Snedecor con p y $k - p$ grados de libertad

Por lo para un nivel de significación α prefijado la región crítica será:

$$\frac{k-p}{p} \cdot \frac{k}{k+1} (y - \bar{X})' S_X^{-1} (y - \bar{X}) \geq F_{1-\alpha, p, k-p}$$

siendo $F_{1-\alpha, p, k-p}$ el cuantil de orden $1 - \alpha$, es decir

$$P\left[F_{p, k-p} \geq F_{1-\alpha, p, k-p} \right] = \alpha$$

3.3. DETECCION DE OUTLIERS.

Consideremos una población que se distribuye según una ley $N_p(\mu, \Sigma)$ no singular de la que hemos extraído una muestra de tamaño n . Se trata de detectar cuáles de estas observaciones son outliers.

Como ya vimos en la introducción, dada una muestra de tamaño n siempre es posible afirmar que en dicha muestra existen k elementos con $k > \lfloor \frac{n}{2} \rfloor$ que no son outliers, es decir son observaciones de la población $N_p(\mu, \Sigma)$ que estamos muestreando.

El método que proponemos para la detección de outliers es el siguiente:

Mediante un muestreo aleatorio simple, extraemos k observaciones de entre las n de la muestra. Sobre estas k observaciones hacemos la hipótesis de que todas ellas pertenecen a la población $N_p(\mu, \Sigma)$ no singular. A continuación calculamos el vector media de estas k observaciones \bar{X} y la matriz de suma de cuadrados y suma de productos S_X de estas observaciones e imponemos a k la condición $k > p$ para que la distribución de S_X sea no degenerada.

Para cada una de las $n - k$ observaciones restantes Y , calculamos el valor del estadístico

$$\frac{k-p}{p} \cdot \frac{k}{k+1} (Y - \bar{X})' S_X^{-1} (Y - \bar{X})$$

Fijado un nivel de significación α , se busca el valor del cuantil $F_{1-\alpha, p, k-p}$ tal que

$$P[F > F_{1-\alpha, p, k-p}] = \alpha$$

donde F se distribuye según una ley F de Snedecor con p y $k - p$ grados de libertad.

Para cada una de las $n - k$ observaciones Y la regla de decisión



a seguir seria.

Si

$$F_{1-\alpha, p, k-p} < \frac{k-p}{p} \cdot \frac{k}{k+1} (y - \bar{X}) S_X^{-1} (y - \bar{X})$$

se rechaza la hipótesis nula de que $Y \in N_p(\mu, \Sigma)$ aceptandose por tanto que dicha observación es outlier.

Midamos ahora el riesgo de error que cometemos, al suponer que las k observaciones obtenidas mediante el muestreo aleatorio simple, X_1, X_2, \dots, X_k no son outliers.

Sea Z la variable aleatoria que representa el número de outliers en una muestra de tamaño n , Z podrá tomar los valores $0, 1, \dots, \lfloor \frac{n}{2} \rfloor$

Si notamos por θ la probabilidad de ser outlier

$$P[Z=r] = \frac{\binom{n}{r} \theta^r (1-\theta)^{n-r}}{\sum_{j=0}^{\lfloor \frac{n}{2} \rfloor} \binom{n}{j} \theta^j (1-\theta)^{n-j}}$$

Notemos por A el suceso de que ninguna de las k observaciones extraídas de la muestra es outlier. Entonces

$$\begin{aligned} P[A] &= \sum_{r=0}^{\lfloor \frac{n}{2} \rfloor} P[A/Z=r] \cdot P[Z=r] = \\ &= \sum_{r=0}^{\lfloor \frac{n}{2} \rfloor} \frac{\binom{r}{0} \binom{n-r}{k}}{\binom{n}{k}} \times \frac{\binom{n}{r} \theta^r (1-\theta)^{n-r}}{\sum_{j=0}^{\lfloor \frac{n}{2} \rfloor} \binom{n}{j} \theta^j (1-\theta)^{n-j}} = \\ &= \frac{1}{\sum_{j=0}^{\lfloor \frac{n}{2} \rfloor} \binom{n}{j} \theta^j (1-\theta)^{n-j}} \sum_{r=0}^{\lfloor \frac{n}{2} \rfloor} \binom{n-k}{r} \theta^r (1-\theta)^{n-r} = \\ &= (1-\theta)^k \frac{\sum_{r=0}^{\lfloor \frac{n}{2} \rfloor} \binom{n-k}{r} \theta^r (1-\theta)^{n-k-r}}{\sum_{j=0}^{\lfloor \frac{n}{2} \rfloor} \binom{n}{j} \theta^j (1-\theta)^{n-j}} \end{aligned}$$

Solo quedaria dar una estimación de la probabilidad θ de ser outlier, para ello pueden seguirse los siguientes criterios.

1ª .- Considerar la probabilidad θ de ser outlier como la probabilidad de que el estadístico caiga en la región crítica. Es decir tomar θ como el nivel de significación del contraste

2ª .- Estimar θ mediante la proporción entre el número de observaciones que detectamos como outliers y el número de total de observaciones.

3ª .- La estimación de θ dada por SHAH (1.966)

$$\hat{\theta} = \frac{a_3 - \left(\left[\frac{n}{2}\right] + 1\right)a_2}{(n-2)a_2 + \frac{1}{2}(n-1)\left(n - \left[\frac{n}{2}\right]\right) - n(n-2)a_1 - n\left[\frac{n}{2}\right]}$$

donde

$$a_i = \frac{\sum_{j=1}^n f_j^i}{n}$$

4ª .- Elegir como probabilidad de ser outlier la que se desprende de los trabajos de DEVLIN y otros (1.975), GNANADESIKAN (1.977), BARNETT (1.978) etc., que consideran esta probabilidad del orden de $1/n$, siendo n el tamaño de la muestra.

3.4. APLICACION.

Aplicamos este metodo al siguiente ejemplo dado por Barnett (1.978) donde se representan las edades y salarios estos ultimos expresados en libras esterlinas , de 55 de los 374 miembros de la seccion de estudiantes del Instituto de Ingenieria Electrica.

27.67	2930	25.33	2300
23.42	2330	25.25	2200
24.67	2480	27.92	3500
27.92	4100	29.50	3600
26.92	2500	21.17	1470
28.92	3380	25.67	2690
26.08	2720	28.42	2860
29.92	4930	26.00	3000
29.50	3020	27.42	3100
22.42	1970	26.58	2600
23.42	1700	25.50	2250
28.00	3100	29.25	3600
23.00	1950	26.00	2750
25.17	2320	29.33	3500
26.58	2750	26.25	3400
22.75	1960	26.83	4500
25.00	2300	27.92	2800
30.00	4120	27.08	3610
23.50	3900	28.33	3100
26.58	5200	28.33	2900
28.25	3200	30.00	3600



28.25	3600	29.92	3610
24.67	2030	30.58	4200
25.42	3520	23.83	3050
22.67	1900	26.33	2760
25.92	3230	26.83	4000
25.25	2500	28.25	3100
25.58	3020		

Para un nivel de significación $\alpha = 0.05$ y para $k = 10$, mediante el siguiente programa de ordenador obtenemos los elementos OUTLIERS.

```

1*      PARAMETER N=5,K=10,J=5,M=2,NK=N-K
2*      EXCL=1,4,5
3*
4*      C
5*      C NEN. DE DATOS, NEN. DE EXCLUSIONES, J= N. INICIAL DE ALEATORIZACION
6*      C M= DIMENSION
7*
8*      C
9*      C DIMENSION DATOS(M*N), EXCL (M*K), IEXCL(K), C(K,K), PR(M,K), S(M,M),
10*      C Y(M,N), X(M,K), YK(M,NK), REM, SINV(M,M), V(6), I(1)
11*
12*      C LECTURA DE UN DATO M-DIMENSIONAL EN CADA TARJETA E IMPRESION
13*      C
14*      READ 1,(DATOS(I,J),I=1,M),J=1,N)
15*      1 FORMAT(2F30.2)
16*      PRINT 1,((DATOS(I,J),I=1,M),J=1,N)
17*      PRINT 2
18*      2 FORMAT(/,/)
19*
20*      C
21*      C ELECCION DE K DATOS ALEATORIOS DE LA MUESTRA SIN REEMPLAZAMIENTO E IMPRESION
22*      C
23*      J=J1
24*      DO 2 I=1,K
25*      30 IEXCL(I)=RANDI(J)
26*      IEXCL(I)=IEXCL(I)+N*(IEXCL(I)/N)
27*      IF(I.EQ.1) GOTO 2
28*      I=I-1
29*      DO 22 J=1,I
30*      IF(IEXCL(I).EQ.IEXCL(J)) GOTO 30
31*      22 CONTINUE
32*      2 CONTINUE
33*      PRINT 3,IEXCL
34*      3 FORMAT(' INDICES ALEATORIOS *7X,10(1X,19),/')
35*      DO 4 I=1,K
36*      I=IEXCL(I)
37*      DO 4 J=1,M
38*      4 EXCL(I,J)=DATOS(I,J)
39*      PRINT 4,EXCL
40*      PRINT 2
41*
42*      C
43*      C CALCULO DE LA MATRIZ DE SUMA DE CUADRADOS Y PRODUCTOS DE K OBSERVACIONES SK
44*      C (EN EL PROGRAMA, LA MATRIZ SK ES S)
45*      C
46*      DIM S
47*      DIM J1
48*      DO 5 I=1,K
49*      DO 6 J=1,K
50*      5 S(I,J)=1./K
51*      DO 6 I=1,K
52*      DO 7 I=1,I+1
53*      DO 7 I=1,M
54*      DO 7 J=1,K
55*      S(I,J)=0
56*      DO 7 L=1,K
57*      7 S(I,J)=S(I,J)+ EXCL(L,I)*C(L,J)
58*      DO 8 I=1,M
59*      DO 8 J=1,M
60*      S(I,J)=S(I,J)+S(I,I)* EXCL(I,L)
61*
62*      C
63*      C IMPRESION DE LA MATRIZ SK
64*      C
65*      PRINT 3(1)
66*      3(1) FORMAT(' MATRIZ SK',/)
67*      PRINT 9,S
68*      9 FORMAT('M(1X,E18.9)')
69*      PRINT 2(1)
70*
71*      C
72*      C CALCULO DE LA MATRIZ SK(-1), INVERSA DE SK, E IMPRESION
73*      C (EN EL PROGRAMA, LA MATRIZ SK(-1) ES SI(1))
74*      C
75*      PRINT 3(2)
76*      3(2) FORMAT(' MATRIZ SK(-1)',/)
77*      DET=S(1,1)*S(2,2)-S(1,2)*S(2,1)
78*      IF((S(1,2).EQ.S(2,1))) GOTO 2(0)
79*      PRINT 2(1), S(1,2),S(2,1)
80*      2(1) FORMAT(' S(2,1)=S(1,2),S(2,1)=S(1,2),/')
81*      2(0) IF(DET.EQ.0) STOP
82*      SINVI(1,1)=S(2,2)/DET
83*      SINVI(2,2)=S(1,1)/DET
84*      SINVI(1,2)=-S(1,2)/DET
85*      SINVI(2,1)=-S(2,1)/DET
86*      PRINT 9,SINV

```

```

83*      PRINT 23
84*
85*      C
86*      C FORMACION DEL VECTOR DE DATOS NO ELEGIDOS E IMPRESION
87*      C (EN EL PROGRAMA, EL VECTOR DE DATOS NO ELEGIDOS ES Y)
88*      C
89*      T0=0
90*      DO 10 I=1,N
91*      A=1
92*      DO 11 J=1,K
93*      A=A*(1-TEXCL(J))
94*      IF(A.C0.0) GOTO 12
95*      DO 13 L=1,M
96*      Y(L,I-TCI)=DATOS(L,I)
97*      GOTO 10
98*      I 12 T0=TCI+1
99*      I 10 CONTINUE
100*      PRINT 303
101*      I 13 FORMAT( ' N-K DATOS NO ELEGIDOS',/)
102*      PRINT 310(I,(Y(J,I),J=1,K),I=1,NK)
103*      I 310 FORMAT(I5,2F20.0)
104*      PRINT 23
105*
106*      C
107*      C
108*      C
109*      C CALCULO DEL ESTADISTICO
110*      C
111*      C
112*      C
113*      C CALCULO DEL VECTOR MEDIA E IMPRESION
114*      C (EN EL PROGRAMA, EL VECTOR MEDIA ES XMEDK)
115*      C
116*      DO 14 I=1,M
117*      XMEDK(I)=0
118*      DO 15 J=1,K
119*      DO 16 I=1,M
120*      XMEDK(I)=XMEDK(I)+EXCL(I,J)
121*      DO 16 I=1,M
122*      XMEDK(I)=XMEDK(I)/K
123*      PRINT 210
124*      I 210 FORMAT(IX,'VECTOR MEDIA',/)
125*      PRINT 1,XMEDK
126*      PRINT 21
127*
128*      C
129*      C
130*      C
131*      C CALCULO DEL VECTOR DIFERENCIA DEL VECTOR DE DATOS NO ELEGIDOS Y EL VECTOR
132*      C MEDIA, E IMPRESION
133*      C (EN EL PROGRAMA, ES YK)
134*      C
135*      DO 17 I=1,M
136*      DO 17 J=1,NK
137*      I 17 YK(I,J)=Y(I,J)-XMEDK(I)
138*      PRINT 211
139*      I 211 FORMAT(IX,'VECTOR DIFERENCIA DEL VECTOR DE DATOS NO ELEGIDOS Y EL
140*      * VECTOR MEDIA',/)
141*      PRINT 1,YK
142*      PRINT 23
143*
144*      C
145*      C
146*      C
147*      C CALCULO DEL F EXPERIMENTAL PARA CADA DATO NO ELEGIDO, E IMPRESION
148*      C (EN EL PROGRAMA ES FEXP)
149*      C
150*      PRINT 212
151*      I 212 FORMAT(IX,'CALCULO DEL ESTADISTICO PARA LOS N-K DATOS NO ELEGIDOS',
152*      **/)
153*      DO 18 J=1,NK
154*      DO 19 I=1,M
155*      P(I)=0
156*      DO 19 L=1,M

```

```

146*      10 P(I)=P(I)+YK(L,J)+SINV(L,I)
147*      FEXP=C
148*      DO 20 I=1,M
149*      20 FEXP=FEXP+P(I)*YK(I,J)
150*      FEXP=FEXP*(K-M)*K/(M*(K+1))
151*      IF(FEXP.GT.87500) GOTO 500
152*      V1=' '
153*      V2=' '
154*      GOTO 501
155*      500 V1='OUTLIE'
156*      V2='R'
157*      501 PRINT 21,J,FEXP,V1,V2
158*      21 FORMAT(IX,'DATO NO ELEGIDO NUMERO =',I4.9X,'F.EXPERIMENTAL=',F10.6,
159*      'BY',Z4C)
160*      PRINT 23
161*      10 CONTINUE
162*
163*      C -----
164*      C GRAFICA DE LOS DATOS
165*      C (EN EL PROGRAMA, LOS DATOS NO ELEGIDOS VIENEN REPRESENTADOS POR '*', Y LOS
166*      C DATOS ELEGIDOS, POR '0')
167*      C -----
168*      DO 503 I=1,61
169*      DO 503 J=1,121
170*      503 V(I,J)=' '
171*      DO 505 I=1,121
172*      505 V(1,I)='-'
173*      DO 506 I=1,61
174*      506 V(I,1)='I'
175*      DO 502 I=1,NK
176*      II=(60*(Y(1,I)-20))/11+1
177*      JJ=(120*(Y(2,I)-1400))/3800+1
178*      502 V(II,JJ)='*'
179*      DO 504 I=1,K
180*      TI=(60*(EXCL(1,I)-20))/11+1
181*      JJ=(120*(EXCL(2,I)-1400))/3800+1
182*      504 V(TI,JJ)='0'
183*      PRINT 507
184*      507 FORMAT(3H1)
185*      DO 508 I=1,61
186*      PRINT 508,(V(I,J),J=1,121)
187*      508 FORMAT(IX,121A1)
188*      508 CONTINUE
189*      STOP
190*      END

```

NO OF COMPI LATION: NO DIAGNOSTICS.

27.57	2376.00
23.42	2370.00
24.57	2430.00
27.92	4100.00
28.92	2500.00
28.92	3320.00
28.02	2720.00
29.92	4970.00
23.50	3020.00
22.42	1370.00
23.42	1700.00
28.00	3100.00
23.00	1950.00
25.17	2320.00
26.58	2750.00
22.75	1360.00
28.00	2300.00
30.00	4120.00
23.50	2900.00
26.58	5200.00

28.25	: 2200.00
28.33	: 2300.00
25.75	: 2200.00
27.92	: 2500.00
29.50	: 2600.00
21.17	: 1470.00
25.67	: 2050.00
20.42	: 2060.00
26.00	: 3000.00
27.42	: 3100.00
26.56	: 2600.00
25.50	: 2250.00
29.75	: 2600.00
27.00	: 2750.00
29.33	: 3500.00
26.25	: 3400.00
26.83	: 4500.00
27.92	: 2800.00
27.09	: 2510.00
20.32	: 3100.00
28.32	: 2900.00
30.00	: 3000.00
28.25	: 3600.00
24.67	: 2020.00
25.42	: 3520.00
22.67	: 1200.00
28.92	: 2230.00
25.25	: 2500.00
25.58	: 2020.00
29.92	: 3510.00
30.50	: 4200.00
23.83	: 3050.00
26.33	: 2760.00
26.82	: 4000.00
28.25	: 3100.00

INDICE ALFABETICO

28 53 5 37 41 33 25 1 45 14

25.42	: 2860.00
26.33	: 2760.00
20.92	: 2500.00
26.83	: 4500.00
27.33	: 2900.00
29.25	: 3000.00
28.50	: 3500.00
27.67	: 2020.00
25.42	: 3520.00
25.17	: 2320.00

MATRI SK

.70401815 4+02 .22577271+04
.72677528 3+04 .382286394+07

MATRI SK(-1)

83.05 .826775283+04 8215 .22577271+04

- 2 47 93 23 9-01 - 3 11 28 62 54-04
 - 11 12 26 48 9-04 - 2 90 04 77 08-06

K-K DATOS NO ELECTUOS

1	23.	2330.
2	25.	2480.
3	28.	4100.
4	29.	3380.
5	26.	2720.
6	30.	4570.
7	25.	3020.
8	22.	1870.
9	21.	1700.
10	28.	3100.
11	23.	1950.
12	27.	2750.
13	23.	1960.
14	25.	2700.
15	30.	4120.
16	23.	3500.
17	27.	5200.
18	28.	3200.
19	25.	2300.
20	26.	2200.
21	28.	3500.
22	21.	1470.
23	26.	2690.
24	26.	3000.
25	27.	3100.
26	27.	2600.
27	25.	2250.
28	26.	2750.
29	28.	3500.
30	26.	3400.
31	28.	2800.
32	27.	3610.
33	28.	3100.
34	30.	3800.
35	28.	3600.
36	25.	2030.
37	21.	1500.
38	26.	3230.
39	28.	2500.
40	25.	3020.
41	30.	3610.
42	31.	4200.
43	24.	3050.
44	27.	4000.
45	28.	3100.

VECTOR MEDIA

27.32

3149.00

VECTOR DIFERENCIA DEL VECTOR DE DATOS NO ELEGIDOS Y EL VECTOR MEDIA

-3.25	-819.00
-2.71	-869.00
.54	951.00
1.54	211.00
-1.30	-429.00
2.54	1781.00
2.12	-179.00
-4.96	-1179.00
-3.96	-1949.00
.63	-49.00
-4.39	-1199.00
-.80	-799.00
-4.01	-1189.00
-2.38	-849.00
2.82	971.00
-3.88	751.00
-.80	2051.00
.87	51.00
-2.05	-899.00
-2.12	-909.00
.54	351.00
-6.21	-1679.00
-1.71	-499.00
-1.39	-199.00
-.04	-49.00
-.80	-549.00
-1.88	-899.00
-1.38	-399.00
1.35	351.00
-1.12	251.00
.54	-249.00
-.30	451.00
.95	-49.00
2.67	451.00
.97	451.00
-2.71	-1119.00
-4.71	-1249.00
-1.46	91.00
-2.12	-549.00
-1.90	-179.00
2.54	451.00
3.20	1051.00
-3.55	-99.00
-.55	951.00
-.87	-49.00

CALCULO DEL ESTADISTICO PARA LOS N-K DATOS NO ELEGIDOS

DATO NO ELEGIDO NUMERO = 1 F.EXPERIMENTAL= 2.345495

DATO NO ELEGIDO NUMERO = 2 F.EXPERIMENTAL= 1.450264

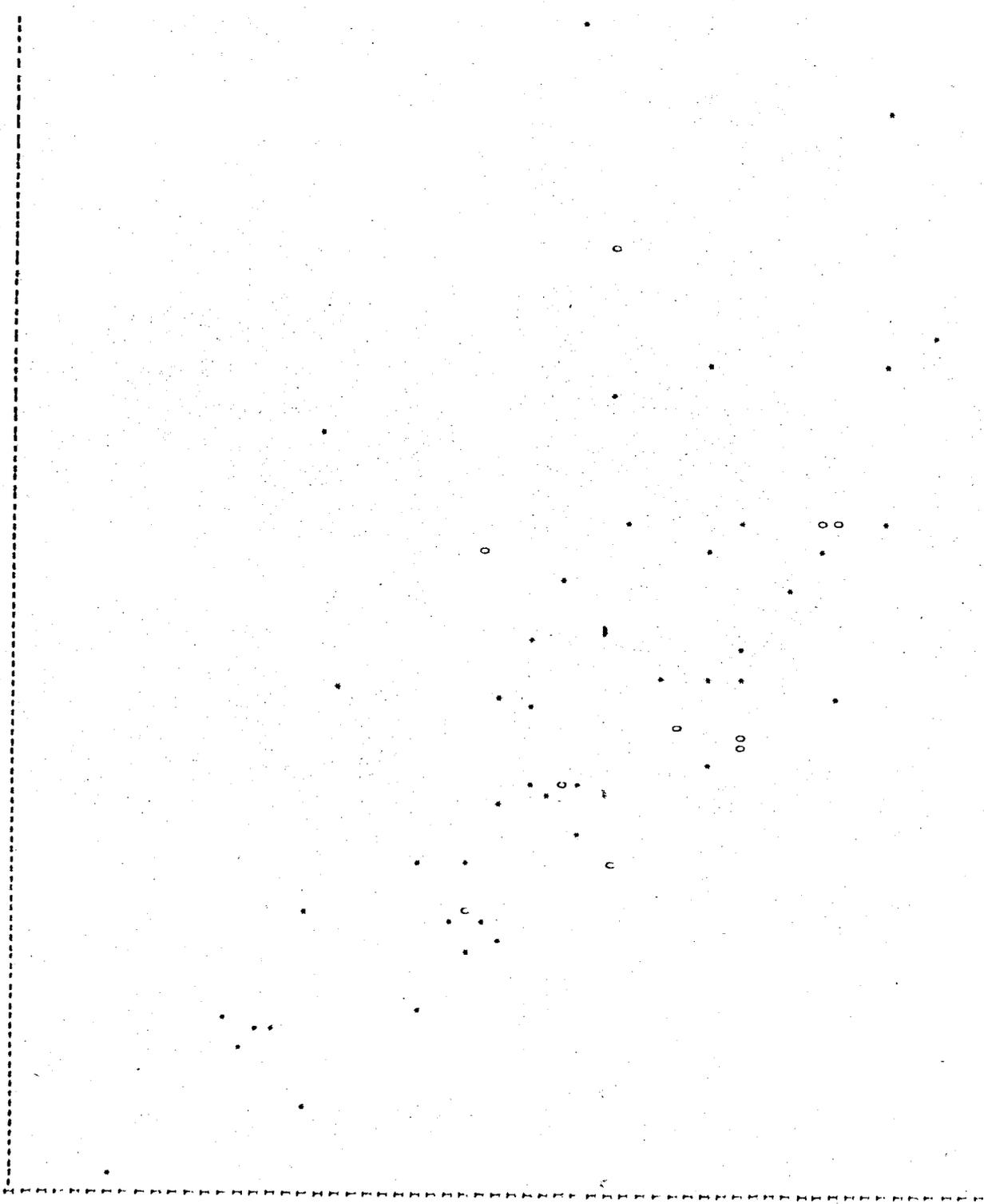
DATO NO ELEGIDO NUMERO = 3 F.EXPERIMENTAL= .870424



DATO NO. ELECTO NUMERO =	4	F.EXPERIMENTAL=	.404213	
DATO NO. ELECTO NUMERO =	5	F.EXPERIMENTAL=	.385245	
DATO NO. ELECTO NUMERO =	6	F.EXPERIMENTAL=	3.474880	
DATO NO. ELECTO NUMERO =	7	F.EXPERIMENTAL=	.973179	
DATO NO. ELECTO NUMERO =	8	F.EXPERIMENTAL=	4.792631	OUTLIER
DATO NO. ELECTO NUMERO =	9	F.EXPERIMENTAL=	3.836191	
DATO NO. ELECTO NUMERO =	10	F.EXPERIMENTAL=	.081686	
DATO NO. ELECTO NUMERO =	11	F.EXPERIMENTAL=	3.941431	
DATO NO. ELECTO NUMERO =	12	F.EXPERIMENTAL=	.212847	
DATO NO. ELECTO NUMERO =	13	F.EXPERIMENTAL=	4.283946	
DATO NO. ELECTO NUMERO =	14	F.EXPERIMENTAL=	1.300327	
DATO NO. ELECTO NUMERO =	15	F.EXPERIMENTAL=	1.690949	
DATO NO. ELECTO NUMERO =	16	F.EXPERIMENTAL=	4.113308	
DATO NO. ELECTO NUMERO =	17	F.EXPERIMENTAL=	4.780483	OUTLIER
DATO NO. ELECTO NUMERO =	18	F.EXPERIMENTAL=	.135750	
DATO NO. ELECTO NUMERO =	19	F.EXPERIMENTAL=	1.144202	

DATO NO. ELECCION NUMERO = 20	F.EXPERIMENTAL= 1.327276	
DATO NO. ELECCION NUMERO = 21	F.EXPERIMENTAL= .137532	
DATO NO. ELECCION NUMERO = 22	F.EXPERIMENTAL= 7.277046	OUTLIER
DATO NO. ELECCION NUMERO = 23	F.EXPERIMENTAL= .537029	
DATO NO. ELECCION NUMERO = 24	F.EXPERIMENTAL= .341429	
DATO NO. ELECCION NUMERO = 25	F.EXPERIMENTAL= .003032	
DATO NO. ELECCION NUMERO = 26	F.EXPERIMENTAL= .370354	
DATO NO. ELECCION NUMERO = 27	F.EXPERIMENTAL= 1.116299	
DATO NO. ELECCION NUMERO = 28	F.EXPERIMENTAL= .402513	
DATO NO. ELECCION NUMERO = 29	F.EXPERIMENTAL= .527444	
DATO NO. ELECCION NUMERO = 30	F.EXPERIMENTAL= .377331	
DATO NO. ELECCION NUMERO = 31	F.EXPERIMENTAL= .221208	
DATO NO. ELECCION NUMERO = 32	F.EXPERIMENTAL= .255784	
DATO NO. ELECCION NUMERO = 33	F.EXPERIMENTAL= .183707	
DATO NO. ELECCION NUMERO = 34	F.EXPERIMENTAL= 1.245900	
DATO NO. ELECCION NUMERO = 35	F.EXPERIMENTAL= .261200	

DATO NO ELEGIDO NUMERO = 36	F.EXPERIMENTAL=	1.897141	
DATO NO ELEGIDO NUMERO = 37	F.EXPERIMENTAL=	4.495087	OUTLIER
DATO NO ELEGIDO NUMERO = 38	F.EXPERIMENTAL=	.442510	
DATO NO ELEGIDO NUMERO = 39	F.EXPERIMENTAL=	.904372	
DATO NO ELEGIDO NUMERO = 40	F.EXPERIMENTAL=	.505268	
DATO NO ELEGIDO NUMERO = 41	F.EXPERIMENTAL=	1.178966	
DATO NO ELEGIDO NUMERO = 42	F.EXPERIMENTAL=	2.317540	
DATO NO ELEGIDO NUMERO = 43	F.EXPERIMENTAL=	2.340552	
DATO NO ELEGIDO NUMERO = 44	F.EXPERIMENTAL=	.902792	
DATO NO ELEGIDO NUMERO = 45	F.EXPERIMENTAL=	.155158	



CAPITULO CUARTO

DETECCION DE OUTLIERS BASADO EN LA
DISTANCIA ENTRE MATRICES SIMETRICAS Y
DEFINIDAS POSITIVAS



CONTENIDO.

- 4.1 - Introducción
- 4.2 - Distancia entre matrices de sumas de cuadrados y sumas de productos
 - Distancia geodésica entre matrices simétricas definidas positivas
 - Distribución del estadístico distancia entre matrices de sumas de cuadrados y sumas de productos
- 4.3 - Acotación
- 4.4 - Detección de outliers
- 4.5 - Aplicación

4.1. INTRODUCCION.

En este capítulo damos un procedimiento para la detección de outliers en muestras procedentes de poblaciones $N_2(\mu, \Sigma)$, utilizando como estadístico el máximo del cuadrado de la distancia entre matrices de sumas de cuadrados y sumas de productos de observaciones muestrales.

Se comienza calculando la distancia entre matrices simétricas y definidas positivas, comprobándose que dicha distancia se puede definir entre matrices de sumas de cuadrados y suma de productos de observaciones muestrales.

A continuación obtenemos la función de densidad del cuadrado de la distancia, acotándose dicha función por otra que puede ser tabulada.

Finalizamos el capítulo exponiendo un criterio para la detección de outliers que se basa en el estadístico ordenado

$$\max_i d^2(S_{(n)}, S_{(n-k)}^{(i)}) \quad i = 1, 2, \dots, \binom{n}{k}$$

4.2. DISTANCIA ENTRE MATRICES DE SUMAS DE CUADRADOS Y SUMAS DE PRODUCTOS.

En este apartado vamos a calcular la distancia entre matrices simétricas y definidas positivas, si se utiliza como forma métrica diferencial la dada por MAAS (1.955). A continuación se comprueba que dicha distancia se puede definir entre las matrices de sumas de cuadrados y sumas de productos (s. c. y s. p.) de observaciones muestrales. Encontrando finalmente la distribución del estadístico distancia.

4.2.1. Distancia geodésica entre matrices simétricas definidas positivas.

4.2.1.1. Definición. Sea H un espacio conexo sobre el que se ha

definido una forma métrica diferencial $(ds)^2$; la distancia entre dos puntos $x, y \in H$ viene dada por

$$d(x, y) = \inf_{\mathcal{X}} \int_x^y ds$$

donde \mathcal{X} es el conjunto de todos los posibles caminos que unen x con y .

La función $d(x, y)$ verifica las siguientes propiedades (KOBAYASHI, NOMIZU, 1.963).

$$d(x, y) \geq 0$$

$$d(x, y) = 0 \iff x = y$$

$$d(x, y) = d(y, x)$$

$$d(x, z) \leq d(x, y) + d(y, z)$$

La topología definida por d es equivalente a la topología inicial del espacio H .

Las matrices simétricas definidas positivas de dimensión $p \times p$, forman un cono convexo en el espacio euclideo $R^{\frac{1}{2}p(p+1)}$. Utilizando ahora la forma métrica diferencial para matrices simétricas definidas positivas dada por Maas, podemos enunciar el siguiente teorema.

4.2.1.1. Teorema .- Sea A y B dos matrices simétricas definidas positivas. Tomando como forma diferencial la dada por Maas, la distancia geodésica entre estas matrices viene dada por

$$d(A, B) = \left(\sum_{i=1}^p (\log \lambda_i)^2 \right)^{1/2}$$

donde λ_i son las raíces de la ecuación determinante $|B - \lambda A| = 0$.

En efecto :

La distancia entre dos matrices A y B vendrá dada por

$$d(A, B) = \inf_{\mathcal{X}} \int_A^B ds$$

donde \mathcal{X} es el conjunto de todos los posibles caminos que unen A con B y ds es la forma métrica diferencial.

Por otro lado se sabe por Analisis Matricial que (BELLMAN , 1965).

" Dadas dos matrices simétricas reales A y B , con A definida positiva, existe una matriz no singular T tal que

$$TAT' = I \quad \text{y} \quad TBT' = \Lambda$$

donde Λ es una matriz diagonal y si B es definida positiva los elementos que componen la matriz Λ , $\lambda_i > 0 \forall i$.

Como $(ds)^2$ es invariante por trasformaciones de congruencias podemos escribir,

$$d(A, B) = \inf_{\mathcal{X}} \int_{I_p} ds$$

donde

$$(ds)^2 = \sum_{i=1}^p \left(\frac{d\lambda_i}{\lambda_i} \right)^2 + \sum_{i < j}^p \frac{(\lambda_i - \lambda_j)^2}{\lambda_i \lambda_j} (dz_{ij})^2$$

Haciendo el cambio $\eta_i = \log \lambda_i \quad i=1, 2, \dots, p$ se tiene

$$d(A, B) = \inf_{\mathcal{X}} \int_0^\eta ds$$

donde ahora

$$\begin{aligned} (ds)^2 &= \sum_{i=1}^p (d\eta_i)^2 + \sum_{i < j}^p \left(e^{\eta_i - \eta_j} + e^{-(\eta_i - \eta_j)} - 2e^{\frac{\eta_i - \eta_j}{2}} e^{-\frac{\eta_i - \eta_j}{2}} \right) (dz_{ij})^2 = \\ &= \sum_{i=1}^p (d\eta_i)^2 + \sum_{i < j}^p \left(\operatorname{sen} h. \left(\frac{1}{2} (\eta_i - \eta_j) \right) \right)^2 (dz_{ij})^2 \end{aligned}$$

Por tanto

$$d(A, B) = \inf_{\mathcal{X}} \int_0^\eta ds = \inf_{\mathcal{X}} \int_0^\eta \sqrt{\sum_{i=0}^p (d\eta_i)^2 + \sum_{i < j}^p \left(\operatorname{sen} h. \left(\frac{1}{2} (\eta_i - \eta_j) \right) \right)^2 (dz_{ij})^2} ds$$

$$\geq \inf_{\alpha} \int_0^{\eta} \left(\sum_{i=0}^p (d\eta_i)^2 \right)^{1/2}$$

Y esto ultimo se hace minimo (1) si como distancia tomamos la euclidea entre el punto $(0,0,\dots,0)$ y $(\eta_1, \eta_2, \dots, \eta_p)$.

En definitiva

$$d(A,B) = \left(\sum_{i=1}^p \eta_i^2 \right)^{1/2}$$

4.2.2. Distribución del estadístico distancia entre matrices de suma de cuadrados y suma de productos.

Aplicaremos el resultado anterior a las matrices de s. c. y s. p. de observaciones muestrales.

Así si

$$X = \begin{pmatrix} X_{11} & X_{12} & \dots & X_{1n} \\ X_{21} & X_{22} & \dots & X_{2n} \\ \dots & \dots & \dots & \dots \\ X_{p1} & X_{p2} & \dots & X_{pn} \end{pmatrix}$$

es una muestra aleatoria simple procedente de una población $N_p(\mu, \Sigma)$ no singular, entonces la matriz de s. c. y s. p. de observaciones muestrales para esta muestra de tamaño n viene dada por

$$S_{(n)} = X \left(I - \frac{1}{n} E_{nn} \right) X'$$

y sabemos que esta matriz es simétrica por su propia construcción, además si no conocemos los parametros poblacionales, $S_{(n)}$ está distribuida según una ley de Wishart p - dimensional con $n - 1$ grados de

(1) . El minimo se puede alcanzar de la misma forma si como trayectorias tomamos

$$\eta_i(t) = t \eta_i \quad 0 \leq t \leq 1 \quad i = 1, 2, \dots, p$$

$$Z_{ij}(t) = 0$$



libertad y matriz asociada $\sum (W_p(n-1, \Sigma))$ (ANDERSON, 1.958) y si se tiene que $n > p$ entonces $S_{(n)}$ es definida positiva con probabilidad uno, esto se puede representar como

$$\int_{S > 0} W_p(n-1, \Sigma) dS = 1$$

Luego podemos afirmar que las matrices de s. c. y s. p. muestrales, cumplen las condiciones exigidas para poder definir la función distancia $d(S_{(n-k)}, S_{(n)})$, salvo conjuntos de medida nula, en la forma.

$$d(S_{(n-k)}, S_{(n)}) = \left[\sum_{i=1}^p (\log \lambda_i)^2 \right]^{1/2} \quad (2)$$

donde $S_{(n-k)}$ es la matriz de s. c. y s. p. de $n-k$ ($n-k > p$) observaciones, y los λ_i son las soluciones de la ecuación determinante.

$$|S_{(n-k)} - \lambda S_{(n)}| = 0$$

Como estamos trabajando con matrices de tipo aleatorio la distancia definida en (2), podemos considerarla como una variable aleatoria y por tanto podemos hablar de la distribución de la variable distancia y para calcular dicha distribución vamos a necesitar la distribución conjunta de las raíces $\lambda_1, \lambda_2, \dots, \lambda_p$ que determinaremos basándonos en los siguientes resultados.

4.2.2.1. Lema. Sea X una muestra aleatoria simple de tamaño n ($n > p$) procedente de una población $N_p(\mu, \Sigma)$ no singular. La matriz de s. c. y s. p. de la muestra X , $S_{(n)}$ puede descomponerse en la forma

$$S_{(n)} = S_{(n-k)} + S_{(k)} + \frac{k(n-k)}{n} (\bar{X}_{(n-k)} - \bar{X}_{(k)}) (\bar{X}_{(n-k)} - \bar{X}_{(k)})'$$

donde $S_{(n-k)}$ es la matriz de s. c. y s. p. de $n-k$ ($n-k > p$) observaciones y $\bar{X}_{(n-k)}$ su vector media, siendo $S_{(k)}$ la matriz de s. c. y s. p. de la k ($k > p$) observaciones restantes y $\bar{X}_{(k)}$ su vector media.

Además

$$S_{(k)} + \frac{k(n-k)}{n} (\bar{X}_{(n-k)} - \bar{X}_{(k)}) (\bar{X}_{(n-k)} - \bar{X}_{(k)})' \in W_p'(k, \Sigma)$$

En efecto

$$\begin{aligned}
 S_{(n)} &= X \left(I - \frac{1}{n} E_{nn} \right) X' = (X_{(n-k)} | X_{(k)}) \left(I - \frac{1}{n} E_{nn} \right) (X_{(n-k)} | X_{(k)})' \\
 &= X_{(n-k)} X_{(n-k)}' + X_{(k)} X_{(k)}' - \frac{1}{n} (X_{(n-k)} | X_{(k)}) \begin{pmatrix} E_{n-k, n} \\ E_{k, n} \end{pmatrix} \begin{pmatrix} X_{(n-k)}' \\ X_{(k)}' \end{pmatrix} \\
 &= X_{(n-k)} X_{(n-k)}' + X_{(k)} X_{(k)}' - \frac{1}{n} \left[X_{(n-k)} E_{n-k, n-k} X_{(n-k)}' + \right. \\
 &\quad \left. + X_{(n-k)} E_{(n-k), k} X_{(k)}' + X_{(k)} E_{k, (n-k)} X_{(n-k)}' + X_{(k)} E_{k, k} X_{(k)}' \right]
 \end{aligned}$$

Sumando y restando

$$\frac{1}{n-k} X_{(n-k)} E_{(n-k), (n-k)} X_{(n-k)}' \quad \text{y} \quad \frac{1}{k} X_{(k)} E_{k, k} X_{(k)}'$$

se obtiene

$$\begin{aligned}
 S_{(n)} &= S_{(n-k)} + S_{(k)} + (n-k) \bar{X}_{(n-k)} \bar{X}_{(n-k)}' - \frac{(n-k)^2}{n} \bar{X}_{(n-k)} \bar{X}_{(n-k)}' + \\
 &\quad + k \bar{X}_{(k)} \bar{X}_{(k)}' - \frac{k^2}{n} \bar{X}_{(k)} \bar{X}_{(k)}' - \frac{(n-k) \cdot k}{n} \bar{X}_{(k)} \bar{X}_{(n-k)}' - \frac{(n-k)k}{n} \bar{X}_{(n-k)} \bar{X}_{(k)}' = \\
 &= S_{(n-k)} + S_{(k)} + \frac{k(n-k)}{n} (\bar{X}_{(n-k)} - \bar{X}_{(k)}) (\bar{X}_{(n-k)} - \bar{X}_{(k)})'
 \end{aligned}$$

Ahora bien, la matriz de s. c. y s. p. $S_{(k)}$ se distribuye segun una ley $\mathcal{W}_p(k-1, \Sigma)$. Por otro lado como:

$$\bar{X}_{(n-k)} \in \mathcal{N}_p \left(\mu, \frac{1}{n-k} \Sigma \right) \quad \text{y} \quad \bar{X}_{(k)} \in \mathcal{N}_p \left(\mu, \frac{1}{k} \Sigma \right)$$

y ambas son independientes se deduce

$$\sqrt{\frac{k(n-k)}{n}} (\bar{X}_{(n-k)} - \bar{X}_{(k)}) \in \mathcal{N}_p(0, \Sigma)$$

Por tanto

$$\left(\sqrt{\frac{k(n-k)}{n}} \right)^2 (\bar{X}_{(n-k)} - \bar{X}_{(k)}) (\bar{X}_{(n-k)} - \bar{X}_{(k)})'$$

se distribuye segun una ley pseudo Wishart p - dimensional con un grado de libertad y matriz asociada Σ .

Por último como

$$S_{(k)} \quad \text{y} \quad \frac{k(n-k)}{n} (\bar{X}_{(n-k)} - \bar{X}_{(k)}) (\bar{X}_{(n-k)} - \bar{X}_{(k)})'$$

son independientes (WILKS , 1.962)

$$S' = S_{(k)} + \frac{k(n-k)}{n} (\bar{X}_{(n-k)} - \bar{X}_{(k)}) (\bar{X}_{(n-k)} - \bar{X}_{(k)})'$$

se distribuye segun una ley $W_p(k, \Sigma)$ en virtud de la reproducti-
vidad de esta ley.

Otro resultado que utilizaremos para encontrar la funcion de den-
sidad del estadistico distancia es el dado por ; HSU (1.939) , FISHER
(1.939) , ROY (1.939) y KHRISNAIAH (1.978).

4.2.2.2. Lema. Sea S' una matriz $p \times p$ que se distribuye segun una
ley $W_p(k, \Sigma)$ y $S_{(n-k)}$ una matriz $p \times p$ que se distribuye segun una ley
 $W_p(n-k-1, \Sigma)$ independiente de la anterior y sea $S_{(n)} = S_{(n-k)} + S'$

Entonces la distribuci3n conjunta de las raices de la ecuaci3n deter-
minante.

$$|S_{(n-k)} - \lambda S_{(n)}| = 0$$

es de la forma

$$f(\lambda_1, \lambda_2, \dots, \lambda_p) = A \prod_{i=1}^p \prod_{j=i+1}^p (\lambda_j - \lambda_i) \left[\prod_{i=1}^p \lambda_i \right]^{\frac{1}{2}(n-k-p-2)} \\ \times \left[\prod_{i=1}^p (1 - \lambda_i) \right]^{\frac{1}{2}(k-p-1)}$$

donde

$$0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_p \leq 1$$

y

$$A = \pi^{p/2} \prod_{i=1}^p \frac{\Gamma(\frac{1}{2}(n-i))}{\Gamma(\frac{1}{2}(n-k-i)) \Gamma(\frac{1}{2}(k-i+1)) \Gamma(\frac{1}{2}(p-i+1))}$$

4.2.2.1. Teorema. Sea X una muestra aleatoria simple de tama3o n
($n > 2$) extraida de una poblaci3n con distribuci3n $N_2(\mu, \Sigma)$ no
singular. Sea $S_{(n)}$ la matriz de s. c. y s. p. de las observaciones

muestrales y $S_{(n-k)}$ la matriz de s. c. y s. p. de $n - k$ observaciones ($k > 2$).

La función de densidad del estadístico cuadrado de la distancia

$$U = d^2(S_{(n-k)}, S_{(n)}) = (\log \lambda_1)^2 + (\log \lambda_2)^2$$

donde λ_1 y λ_2 son las soluciones de la ecuación determinante

$$|S_{(n-k)} - \lambda S_{(n)}| = 0$$

viene dada por

$$\begin{cases} A e^{-\sqrt{u} \left(\frac{n-k-2}{2}\right)} \int_{\sqrt{2}-1}^1 \left(e^{-\sqrt{u} \frac{1-w^2}{1+w^2}} - e^{-\sqrt{u} \frac{2w}{1+w^2}} \right) e^{-\sqrt{u} (n-k-2) w \frac{1-w}{1+w^2}} \frac{1}{1+w^2} \times \\ \times \left[\left(1 - e^{-\sqrt{u} \frac{2w}{1+w^2}} \right) \left(1 - e^{-\sqrt{u} \frac{1-w^2}{1+w^2}} \right) \right]^{\frac{k-3}{2}} dw, \quad 0 \leq u < \infty. \\ 0 \end{cases} \quad \text{en el resto}$$

siendo

$$A = \frac{\sqrt{\pi} \Gamma\left(\frac{1}{2}(n-1)\right) \Gamma\left(\frac{1}{2}(n-2)\right)}{\Gamma\left(\frac{1}{2}(n-k-1)\right) \Gamma\left(\frac{1}{2}(n-k-2)\right) \Gamma\left(\frac{k}{2}\right) \Gamma\left(\frac{k-1}{2}\right)}$$

En efecto

En virtud del Lema 4.2.2.2.

$$f(\lambda_1, \lambda_2) = A (\lambda_2 - \lambda_1) (\lambda_1 \lambda_2)^{\frac{1}{2}(n-k-4)} \left[(1-\lambda_1)(1-\lambda_2) \right]^{\frac{k-3}{2}} \quad 0 \leq \lambda_1 \leq \lambda_2 \leq 1$$

y

$$A = \frac{\sqrt{\pi} \Gamma\left(\frac{1}{2}(n-1)\right) \Gamma\left(\frac{1}{2}(n-2)\right)}{\Gamma\left(\frac{1}{2}(n-k-1)\right) \Gamma\left(\frac{1}{2}(n-k-2)\right) \Gamma\left(\frac{k}{2}\right) \Gamma\left(\frac{k-1}{2}\right)}$$

Haciendo el cambio de variables

$$X = -\ln \lambda_1$$

$$Y = -\ln \lambda_2$$

se obtiene

$$f(x, y) = A (e^{-y} - e^{-x}) e^{-(x+y) \frac{1}{2}(n-k-2)} \left[(1-e^{-y})(1-e^{-x}) \right]^{\frac{k-3}{2}} \quad 0 < y < x < \infty$$

y si hacemos

$$\theta = X^2$$

$$\varphi = Y^2$$

obtendremos:

$$f(\theta, \varphi) = \frac{A}{4} (e^{-\sqrt{\varphi}} - e^{-\sqrt{\theta}}) e^{-(\sqrt{\theta} + \sqrt{\varphi}) \frac{1}{2}(n-k-2)} \left[(1 - e^{-\sqrt{\varphi}})(1 - e^{-\sqrt{\theta}}) \right]^{\frac{k-3}{2}} (\theta\varphi)^{-\frac{1}{2}}$$

$0 < \varphi < \theta < \infty$

y por ultimo realizando la transformación

$$U = \theta + \varphi$$

$$t = \theta$$

Se obtiene

$$f(u, t) = \frac{A}{4} (e^{-\sqrt{u-t}} - e^{-\sqrt{t}}) e^{-(\sqrt{t} + \sqrt{u-t}) \frac{1}{2}(n-k-2)} \left[(1 - e^{-\sqrt{u-t}})(1 - e^{-\sqrt{t}}) \right]^{\frac{k-3}{2}} \times (t(u-t))^{-\frac{1}{2}}$$

$0 < U < \infty, \frac{U}{2} < t < U$

Por lo que la función de densidad marginal de U, viene dada por

$$f(u) = \frac{A}{4} \int_{U/2}^u (e^{-\sqrt{u-t}} - e^{-\sqrt{t}}) e^{-(\sqrt{t} + \sqrt{u-t}) \frac{1}{2}(n-k-2)} \left[(1 - e^{-\sqrt{u-t}})(1 - e^{-\sqrt{t}}) \right]^{\frac{k-3}{2}} [t(u-t)]^{-\frac{1}{2}} dt$$

$0 < U < \infty$

Efectuando ahora el cambio de variable

$$z = \sqrt{\frac{2(u-t)}{u}}$$

se obtiene

$$f(u) = \frac{A}{4} \int_0^1 (e^{-\frac{z}{\sqrt{2}} \sqrt{u}} - e^{-\sqrt{u} \sqrt{1 - \frac{z^2}{2}}}) e^{-\sqrt{u} \frac{n-k-2}{2} \left(\frac{z}{\sqrt{2}} + \sqrt{1 - \frac{z^2}{2}} \right)} \left[(1 - e^{-\sqrt{u} \sqrt{1 - \frac{z^2}{2}}}) (1 - e^{-\frac{z}{\sqrt{2}} \sqrt{u}}) \right]^{\frac{k-3}{2}} \times \sqrt{2} \left(1 - \frac{z^2}{2}\right)^{-\frac{1}{2}} dz$$

$0 < U < \infty$

y realizando la transformación

$$\sqrt{1 - \frac{z^2}{2}} = \left(1 + \frac{z}{\sqrt{2}}\right) w$$

se deduce finalmente

$$f(u) = A e^{-\sqrt{u} \frac{(n-k-2)}{2}} \int_{\sqrt{2}-1}^1 (e^{-\sqrt{u} \frac{1-w^2}{1+w^2}} - e^{-\sqrt{u} \frac{2w}{1+w^2}}) e^{-\sqrt{u} (n-k-2) w \frac{1-w}{1+w^2}} \left[(1 - e^{-\sqrt{u} \frac{2w}{1+w^2}}) (1 - e^{-\sqrt{u} \frac{1-w^2}{1+w^2}}) \right]^{\frac{k-3}{2}} \times \frac{dw}{1+w^2}$$

$0 < U < \infty$

4.3. ACOTACION.

4.3.1. Teorema. La función

$$f(w) = A e^{-\sqrt{u} \left(\frac{n-k-2}{2}\right)} \int_{\sqrt{2}-1}^1 \left(e^{-\sqrt{u} \frac{1-w^2}{1+w^2}} - e^{-\sqrt{u} \frac{2w}{1+w^2}} \right) e^{-\sqrt{u} (n-k-2) w \frac{1-w}{1+w^2}} \frac{1}{1+w^2} \times \\ \times \left[\left(1 - e^{-\sqrt{u} \frac{2w}{1+w^2}} \right) \left(1 - e^{-\sqrt{u} \frac{1-w^2}{1+w^2}} \right) \right]^{\frac{k-3}{2}} dw \quad 0 < u < \infty$$

siendo

$$A = \frac{\sqrt{\pi} \Gamma\left(\frac{1}{2}(n-1)\right) \Gamma\left(\frac{1}{2}(n-2)\right)}{\Gamma\left(\frac{1}{2}(n-k-1)\right) \Gamma\left(\frac{1}{2}(n-k-2)\right) \Gamma\left(\frac{k}{2}\right) \Gamma\left(\frac{k-1}{2}\right)}$$

está acotada por la función

$$g(u) = A e^{-\sqrt{u} \left(\frac{n-k-2}{2}\right)} (1 - e^{-\sqrt{u}})^{k-2}$$

En efecto:

Sea

$$y_1 = \left(e^{-\sqrt{u} \frac{1-w^2}{1+w^2}} - e^{-\sqrt{u} \frac{2w}{1+w^2}} \right)$$

$$y_1' = e^{-\sqrt{u} \frac{1-w^2}{1+w^2}} \sqrt{u} \frac{4w}{(1+w^2)^2} + e^{-\sqrt{u} \frac{2w}{1+w^2}} \sqrt{u} \frac{2-2w^2}{(1+w^2)^2} > 0 \quad \forall w \in I$$

luego y_1 es una función creciente en el intervalo de integración

$I = (\sqrt{2}-1, 1)$. Por tanto

$$\sup_I y_1 = (1 - e^{-\sqrt{u}})$$

Sea

$$y_2 = e^{-\sqrt{u} (n-k-2) w \frac{1-w}{1+w^2}}$$

$$y_2' = - e^{-\sqrt{u} (n-k-2) w \frac{1-w}{1+w^2}} \sqrt{u} (n-k-2) \left[\frac{1-2w-w^2}{(1+w^2)^2} \right]$$

esta derivada primera se anula en los puntos

$$w_1 = -1 - \sqrt{2}$$

$$w_2 = -1 + \sqrt{2}$$

dándose

$$y_2''(w_1) < 0$$

$$y_2''(w_2) > 0$$

Luego

$$\sup_I y_2 = 1$$

Tomemos

$$y_3 = \left(1 - e^{-\sqrt{u} \frac{2w}{1+w^2}}\right)^{\frac{k-3}{2}}$$

$$y_3' = \frac{k-3}{2} \left(1 - e^{-\sqrt{u} \frac{2w}{1+w^2}}\right)^{\frac{k-5}{2}} e^{-\sqrt{u} \frac{2w}{1+w^2}} \sqrt{u} \frac{2-2w^2}{(1+w^2)^2}$$

e $y_3' > 0 \quad \forall w \in (\sqrt{2}-1, 1)$, por lo que alcanza el supremo en $w=1$
así

$$\sup_I y_3 = \left(1 - e^{-\sqrt{u}}\right)^{\frac{k-3}{2}}$$

Sea

$$y_4 = \left(1 - e^{-\sqrt{u} \frac{1-w^2}{1+w^2}}\right)^{\frac{k-3}{2}}$$

$$y_4' = \frac{k-3}{2} \left(1 - e^{-\sqrt{u} \frac{1-w^2}{1+w^2}}\right)^{\frac{k-5}{2}} e^{-\sqrt{u} \frac{1-w^2}{1+w^2}} \sqrt{u} \left[\frac{-4w}{(1+w^2)^2} \right]$$

e $y_4' < 0$, por tanto el supremo en el intervalo de integración se alcanza en el extremo inferior del mismo así:

$$\sup_I y_4 = \left(1 - e^{-\sqrt{u} \cdot 0.7071}\right)^{\frac{k-3}{2}}$$

Y por ultimo

$$\int_{\sqrt{2}-1}^1 \frac{dw}{1+w^2} = \arctg w \Big|_{\sqrt{2}-1}^1 \leq 0.5$$

Por tanto

$$f(w) = A \int_{\sqrt{2}-1}^1 y_1 \cdot y_2 \cdot y_3 \cdot y_4 \frac{dw}{1+w^2} \leq A e^{-\sqrt{u} \left(\frac{n-k-2}{2}\right)} (1 - e^{-\sqrt{u}}) (1 - e^{-\sqrt{u}})^{\frac{k-3}{2}} \\ \times \left(1 - e^{-\sqrt{u} \cdot 0.7071}\right)^{\frac{k-3}{2}} \times \frac{1}{2}$$

así

$$f(w) \leq A e^{-\sqrt{u} \left(\frac{n-k-2}{2}\right)} (1 - e^{-\sqrt{u}})^{k-2} = g(u) \quad \forall u \quad 0 < u < \infty$$



4.3.2. Teorema. La función $q(x)$ definida mediante la expresión

$$q(x) = \frac{1}{2 \sum_{j=0}^{k-2} \left[\binom{k-2}{j} (-1)^j \frac{1}{\left(\frac{n-k-2}{2} + j\right)^2} \right]} e^{-\sqrt{x} \left(\frac{n-k-2}{2}\right)} (1 - e^{-\sqrt{x}})^{k-2} \quad 0 \leq x < \infty$$

con $k > 2$ y $(n-k) > 2$ es función de densidad.

En efecto:

Para que $q(x)$ sea función de densidad ha de ocurrir que

$$q(x) \geq 0 \quad \forall x$$

$$\int_0^{\infty} q(x) dx = 1$$

La primera propiedad es evidente que se cumple. Veamos que es cierta la segunda.

$$\begin{aligned} \int_0^{\infty} q(x) dx &= \frac{1}{2 \sum_{j=0}^{k-2} \left[\binom{k-2}{j} (-1)^j \frac{1}{\left(\frac{n-k-2}{2} + j\right)^2} \right]} \int_0^{\infty} e^{-\sqrt{x} \left(\frac{n-k-2}{2}\right)} (1 - e^{-\sqrt{x}})^{k-2} dx = \\ &= \frac{1}{2 \sum_{j=0}^{k-2} \left[\binom{k-2}{j} (-1)^j \frac{1}{\left(\frac{n-k-2}{2} + j\right)^2} \right]} \sum_{j=0}^{k-2} (-1)^j \binom{k-2}{j} \int_0^{\infty} e^{-\sqrt{x} \left(\frac{n-k-2}{2} + j\right)} dx \end{aligned}$$

Realizando el cambio de variable

$$\sqrt{x} = y$$

se obtiene

$$\int_0^{\infty} q(x) dx = \frac{1}{2 \sum_{j=0}^{k-2} \left[\binom{k-2}{j} (-1)^j \frac{1}{\left(\frac{n-k-2}{2} + j\right)^2} \right]} \sum_{j=0}^{k-2} (-1)^j \binom{k-2}{j} \int_0^{\infty} y e^{-y \left(\frac{n-k-2}{2} + j\right)} dy$$

La función de densidad $g(x)$ presenta la ventaja sobre $f(u)$ de que su función de distribución se puede tabular, aplicando para ello el método de Simpson. Y además $F(u)$ domina estocásticamente a $G(x)$.

En las páginas 83 a 87 se incluye un programa de ordenador para las representaciones gráficas de $f(u)$ y $g(x)$, determinando dichas representaciones en el caso particular $n = 55$ y $k = 10$.

A continuación en las páginas 88 y 89 se incluye un programa para la tabulación de $G(x)$ en la situación práctica $n = 15$ y $k = 6$.

Representación gráfica de las funciones $f(u)$ y $g(x)$.

```

1* C
2* C COMPARACION Y REPRESENTACION DE LAS FUNCIONES DE DENSIDAD DE LA VARIABLE
3* C DISTANCIA AL CUADRADO F(U) Y LA DADA EN EL TEOREMA 4.3.2.
4* C (EN EL PROGRAMA, F(U) ES F, Y LA DEL TEOREMA ES G)
5* C
6* -----
7* IMPLICIT DOUBLE PRECISION(A-H,O-U,W,Z)
8* COMMON N,K,U
9* DIMENSION A(2),V(1:81),U1(1),F2(1),G1(1)
10* EXTERNAL F1
11* LOGICAL REL
12* N=5
13* U=0
14* K=10
15* DO 10 I=1,F1
16* DO 10 J=1,F1
17* 10 V(I,J)=*
18* DO 11 I=1,F1
19* 11 V(I,1)=*
20* DO 12 I=1,81
21* 12 V(1,I)=*
22* G1(1)=0
23* F2(1)=0
24* U1(1)=0
25* C
26* -----
27* C CALCULO DE LAS CONSTANTES B Y C
28* C
29* B=1/(N-K-2)**2
30* K2=K-2
31* DO 1 I=1,K2
32* 1 B=0+(-1)**I*(FACT(K2)/FACT(IR)*FACT(K2-IR))
33* *(4./(N-K-2+2*IR)**2)
34* C=FACT(N-3)/(4*FACT(K2)*FACT(N-K-1))
35* C
36* -----
37* C CALCULO DE LAS FUNCIONES F Y G
38* C
39* A1=DSQRT(2.0)-1.0
40* A2=1.0
41* E=1.0-E
42* MAXIT=200
43* REL=TRUE
44* PRINT 9
45* 9 FORMAT(IX,/,1X,*,VALORES COMPARATIVOS DE LAS FUNCIONES F Y G,/)
46* DO 2 I=1,80
47* U=*.05U/80
48* AINT=SIN(PI*(F1+A*E*REL*MAXIT+.1*D.17))
49* F=C*AINT*DEXP(-DSQRT(U)*(N-K-2)/2.)
50* G=(DEXP(-DSQRT(U))*((N-K-2)/2.))*(1.-DEXP(-DSQRT(U)))*(K-2)
51* G=1/G
52* I=1
53* J=100*F/10+1
54* J=100*G/10+1
55* IF(J.NE.J0) GOTO22
56* V(I,J)=*
57* GOTO 22
58* 22 V(I,J)=*
59* V(I,J0)=*
60* 23 G1(I)=F
61* F2(I)=F
62* U1(I)=U
63* PRINT6,U,F,G
64* 6 FORMAT(IX,*,U=*,09.3,10X,*,F(U)=*,D12.6,10X,*,G(X)=*,D12.6,/)
65* GOTO 2
66* 3 PRINT 7,U
67* 7 FORMAT(IX,*,ERROR EN U=*,D9.7)
68* 2 CONTINUE
69* C
70* -----
71* C GRAFICA DE LAS FUNCIONES DE DENSIDAD F Y G
72* C
73* PRINT 31
74* 31 FORMAT(1H1)
75* PRINT 30,(U1(I),F2(I),G1(I),(V(I,J),J=1,81)),I=1,F1)
76* 30 FORMAT(1X,D1.3,1X,D12.6,1X,D17.6,10X,81A1)
77* STOP
78* END

```

END OF COMPILATION: NO DIAGNOSTICS.

```

1*      DOUBLE PRECISION FUNCTION SIM3NI(FX,A,E,REL,MAXIT,FK,S)
2*      IMPLICIT DOUBLE PRECISION(A-H,O-Z)
3*      DIMENSION A(2)
4*      LOGICAL REL
5*      PREV=0.
6*      H=(A(2)-A(1))/3.
7*      X=A(1)
8*      N=0
9*      M=3
10*     S=0.
11*     DO 3 J=1,MAXIT
12*     DO 1 I=N,M
13*     R=3.
14*     IF(MOD(I,3).EQ.2*N) R=N+1.
15*     S=S+FX(X,FK)*R
16*     1 X=X+H
17*     SIM3NI=S*H+.37500/(N+1.00)
18*     IF(N.EQ.0) GO TO 2
19*     H=H*.500
20*     M=2*M
21*     R=SIM3NI-PREV
22*     IF(ABS(R).GT.1/SIM3NI)
23*     IF(ABS(R).LT.E) GO TO 4
24*     2 PREV=SIM3NI
25*     N=1
26*     3 X=A(1)+.500*H
27*     RETURN 7
28*     4 RETURN
29*     END

```

SIM3NT
SIM3NI

SIM3NI

END OF COMPILATION: NO DIAGNOSTICS.

```

1*      DOUBLE PRECISION FUNCTION FACTM
2*      COMMON N,K
3*      IF(M) 1,3,4
4*      1 PRINT 2,N,K
5*      2 FORMAT(' ARGUMENTO NEGATIVO PARA UN FACTORIAL *N=*I4,* K=*I4)
6*      STOP
7*      3 FACT=1
8*      RETURN
9*      4 FACT=1
10*     DO 5 I=1,M
11*     5 FACT=FACT*I
12*     RETURN
13*     END

```

END OF COMPILATION: NO DIAGNOSTICS.

```

1*      DOUBLE PRECISION FUNCTION FJ(T,FK)
2*      IMPLICIT DOUBLE PRECISION(A-H,O-U,W-Z)
3*      COMMON N,K,U
4*      R=EXP(-DSQRT(U)*(1.-T**2)/(1.+T**2))
5*      S=EXP(-DSQRT(U)*2.*T/(1.+T**2))
6*      Q=EXP(-DSQRT(U)*(N-K-2)*T*(1.-T)/(1.+T**2))
7*      FJ=(R-S)*Q*(1.-S)*(1.-R)**((K-3)/2.)/(1.+T**2)
8*      RETURN
9*      END

```

END OF COMPILATION: NO DIAGNOSTICS.

VALORES COMPARATIVOS DE LAS FUNCIONES F Y G

U= .108 -0.01	F(U)= .053437+0.00	G(X)= .147171+0.00
U= .217 -0.01	F(U)= .267517+0.01	G(X)= .787149+0.00
U= .325 -0.01	F(U)= .477508+0.01	G(X)= .172039+0.01
U= .433 -0.01	F(U)= .628437+0.01	G(X)= .267974+0.01
U= .542 -0.01	F(U)= .710423+0.01	G(X)= .350505+0.01
U= .650 -0.01	F(U)= .730630+0.01	G(X)= .414475+0.01
U= .758 -0.01	F(U)= .720672+0.01	G(X)= .457726+0.01
U= .867 -0.01	F(U)= .660306+0.01	G(X)= .482580+0.01
U= .975 -0.01	F(U)= .625784+0.01	G(X)= .492021+0.01
U= .108 +0.00	F(U)= .565043+0.01	G(X)= .489245+0.01
U= .217 +0.00	F(U)= .503272+0.01	G(X)= .477255+0.01
U= .325 +0.00	F(U)= .443713+0.01	G(X)= .458846+0.01
U= .433 +0.00	F(U)= .389215+0.01	G(X)= .435578+0.01
U= .542 +0.00	F(U)= .337689+0.01	G(X)= .409697+0.01
U= .650 +0.00	F(U)= .292446+0.01	G(X)= .382388+0.01
U= .758 +0.00	F(U)= .252420+0.01	G(X)= .354630+0.01
U= .867 +0.00	F(U)= .217327+0.01	G(X)= .327175+0.01
U= .975 +0.00	F(U)= .186766+0.01	G(X)= .300488+0.01
U= .108 +0.01	F(U)= .160289+0.01	G(X)= .274984+0.01
U= .217 +0.01	F(U)= .137438+0.01	G(X)= .250871+0.01
U= .325 +0.01	F(U)= .117777+0.01	G(X)= .228281+0.01
U= .433 +0.01	F(U)= .100836+0.01	G(X)= .207272+0.01
U= .542 +0.01	F(U)= .864275+0.00	G(X)= .187851+0.01
U= .650 +0.01	F(U)= .740399+0.00	G(X)= .169986+0.01
U= .758 +0.01	F(U)= .634428+0.00	G(X)= .153619+0.01
U= .867 +0.01	F(U)= .543819+0.00	G(X)= .138677+0.01
U= .975 +0.01	F(U)= .466365+0.00	G(X)= .125073+0.01
U= .108 +0.02	F(U)= .400159+0.00	G(X)= .112718+0.01
U= .217 +0.02	F(U)= .343560+0.00	G(X)= .101520+0.01
U= .325 +0.02	F(U)= .295162+0.00	G(X)= .917881+0.00
U= .433 +0.02	F(U)= .253763+0.00	G(X)= .822337+0.00
U= .542 +0.02	F(U)= .218332+0.00	G(X)= .739726+0.00
U= .650 +0.02	F(U)= .187995+0.00	G(X)= .665273+0.00
U= .758 +0.02	F(U)= .162002+0.00	G(X)= .598174+0.00
U= .867 +0.02	F(U)= .139717+0.00	G(X)= .537796+0.00

U= .390 +0.00	F(U) = .120598 +0.00	G(X) = .487463 +0.00
U= .401 +0.00	F(U) = .104183 +0.00	G(X) = .434649 +0.00
U= .412 +0.00	F(U) = .900789 -0.01	G(X) = .790757 +0.00
U= .422 +0.00	F(U) = .779507 -0.01	G(X) = .751319 +0.00
U= .433 +0.00	F(U) = .675133 -0.01	G(X) = .715890 +0.00
U= .444 +0.00	F(U) = .585236 -0.01	G(X) = .784070 +0.00
U= .455 +0.00	F(U) = .507744 -0.01	G(X) = .755492 +0.00
U= .466 +0.00	F(U) = .440890 -0.01	G(X) = .729828 +0.00
U= .477 +0.00	F(U) = .383164 -0.01	G(X) = .706781 +0.00
U= .487 +0.00	F(U) = .333278 -0.01	G(X) = .686087 +0.00
U= .498 +0.00	F(U) = .290130 -0.01	G(X) = .667495 +0.00
U= .509 +0.00	F(U) = .252783 -0.01	G(X) = .650799 +0.00
U= .520 +0.00	F(U) = .220420 -0.01	G(X) = .635801 +0.00
U= .531 +0.00	F(U) = .192360 -0.01	G(X) = .622326 +0.00
U= .542 +0.00	F(U) = .168009 -0.01	G(X) = .610217 +0.00
U= .552 +0.00	F(U) = .146859 -0.01	G(X) = .599349 -0.01
U= .553 +0.00	F(U) = .128474 -0.01	G(X) = .589552 -0.01
U= .574 +0.00	F(U) = .112478 -0.01	G(X) = .580756 -0.01
U= .585 +0.00	F(U) = .985514 -0.02	G(X) = .572846 -0.01
U= .596 +0.00	F(U) = .864354 -0.02	G(X) = .565730 -0.01
U= .607 +0.00	F(U) = .758315 -0.02	G(X) = .559277 -0.01
U= .617 +0.00	F(U) = .665940 -0.02	G(X) = .553585 -0.01
U= .628 +0.00	F(U) = .585253 -0.02	G(X) = .548372 -0.01
U= .639 +0.00	F(U) = .514730 -0.02	G(X) = .543705 -0.01
U= .650 +0.00	F(U) = .453015 -0.02	G(X) = .539497 -0.01

Tabulación de la función de distribución $G(x)$

```

1*      DOUBLE PRECISION FUNCTION SIM3NI(FX,A,E,REL,MAXIT,FK,S)      SIM3NI
2*      IMPLICIT DOUBLE PRECISION(A-H,O-Z)                          SIM3NI
3*      C-----
4*      C      INTEGRACION NUMERICA UTILIZANDO LA REGLA 3/8 DE SIMPSON      SIM3NI
5*      C-----
6*      DIMENSION A(2)
7*      LOGICAL REL
8*      PREV=0.
9*      C-----
10*     C      INICIALIZAR H,X,N,M,S
11*     C-----
12*     H=(A(2)-A(1))/3.
13*     X=A(1)
14*     N=0
15*     M=3
16*     S=0.
17*     C-----
18*     C      CICLO HASTA EL MAXIMO NUMERO DE EVALUACIONES
19*     C-----
20*     DO 3 J=1,MAXIT
21*     C-----
22*     C      CICLO PARA EL NUMERO DE EVALUACIONES
23*     C-----
24*     DO 1 I=N,M
25*     R=3.
26*     C-----
27*     C      DETERMINACION DEL COEFICIENTE R
28*     C-----
29*     IF(MOD(I,3).EQ.2*N) R=N+1.
30*     C-----
31*     C      SUMA DE LAS EVALUACIONES DE LA FUNCION
32*     C-----
33*     S=S+FX(X,FK)*R
34*     C-----
35*     C      INCREMENTAR X
36*     C-----
37*     1 X=X+H
38*     C-----
39*     C      OBTENCION DEL NUEVO VALOR DE INTEGRACION
40*     C-----
41*     SIM3NI=S*H*.37500/(N+1.001)
42*     C-----
43*     C      TEST PARA EL PRIMER CICLO
44*     C-----
45*     IF(N.EQ.0) GO TO 2
46*     C-----
47*     C      PARA LOS CICLOS DISTINTOS DEL PRIMERO, MITAD DE H Y DOBLE DE H
48*     C-----
49*     H=H*.500
50*     M=2*M
51*     C-----
52*     C      REVISION DEL ERROR DE CONTROL
53*     C-----
54*     R=SIM3NI-PREV
55*     IF(REL) P=R/SIM3NI
56*     C-----
57*     C      SI EL ERROR ESTA DENTRO DE LOS LIMITES DADOS, FIN
58*     C-----
59*     IF(DABS(R).LT.F) GO TO 4
60*     C-----
61*     C      FIJAR UN NUEVO VALOR DE INTEGRACION
62*     C-----
63*     2 PREV=SIM3NI
64*     N=1
65*     C-----
66*     C      OBTENER EL NUEVO LIMITE INFERIOR PARA LA EVALUACION DE LA FUNCION
67*     C-----
68*     3 X=A(1)+.500*H
69*     RETURN 7
70*     4 RETURN
71*     END
SIM3NI

```

END OF COMPILATION: NO DIAGNOSTICS.

```

1*      IMPLICIT DOUBLE PRECISION(A-H,O-Z)
2*      EXTERNAL F1
3*      LOGICAL REL
4*      COMMON N,K
5*      N=15
6*      K=6
7*      DIMENSION A(2)
8*      B=4./(N-K-2.)**2
9*      K2=K-2
10*     DO 1 IR=1,K2
11*     1 B=B*(-1)**IR*(FACT(K2)/(FACT(IR)*FACT(K2-IR)))
12*     *I4./(N-K-2.*2*IR)**2)
13*     A(1)=0
14*     A(2)=9
15*     E=1.D-6
16*     MAXIT=200
17*     REL=.TRUE.
18*     DO 3 I=1,200
19*     A(2)=A(2)+1.D-1
20*     AINT=SIMSNI(F1*A*E*REL*MAXIT*1.D*52)
21*     F=AINT/(2*B)
22*     PRINT 5,A(2)*F
23*     6 FORMAT(10X,'F(',D9.4,') =',D20.10)
24*     GO TO 3
25*     2 PRINT 4,N,K
26*     4 FORMAT(' ERROR EN LA INTEGRAL PARA N=',I4,' K=',I4)
27*     3 CONTINUE
28*     STOP
29*     END

```

END OF COMPILATION: NO DIAGNOSTICS.

```

1*      DOUBLE PRECISION FUNCTION FACT(M)
2*      COMMON N,K
3*      IF(M) 1,3,4
4*      1 PRINT 2,N,K
5*      2 FORMAT(' ARGUMENTO NEGATIVO PARA UN FACTORIAL .N=',I4,' K=',I4)
6*      STOP
7*      3 FACT=1
8*      RETURN
9*      4 FACT=1
10*     DO 5 I=1,M
11*     5 FACT=FACT*I
12*     RETURN
13*     END

```

END OF COMPILATION: NO DIAGNOSTICS.

```

1*      DOUBLE PRECISION FUNCTION F1(U,K)
2*      IMPLICIT DOUBLE PRECISION(A-H,O-Z)
3*      COMMON N,K
4*      F1=(DEXP(-DSORT(U))*((N-K-2.)/2.))*(1-DEXP(-DSORT(U)))**K-2)
5*      RETURN
6*      END

```

END OF COMPILATION: NO DIAGNOSTICS.



4.4 DETECCION DE OUTLIERS.

En el apartado 4.2 se ha obtenido la función de densidad de el cuadrado de la distancia entre la matriz de s.c. y s.p. de una muestra aleatoria de tamaño n extraída de una población $N_2(\mu, \Sigma)$ no singular y la matriz de s. c. y s. p. de $n - k$ de estas observaciones.

Debido a que se pueden extraer de $\binom{n}{k}$ formas diferentes estas $n - k$ observaciones de entre las n de la muestra para cada k fijo variando dentro de las condiciones del teorema 4.2.2.1, obtendríamos matrices de s. c. y s. p. $S_{(n-k)}$. Por otro lado bajo la hipótesis de que los elementos de la muestra de tamaño n pertenecen a la misma población $N_2(\mu, \Sigma)$, se tendrá que $\exists q \in \mathbb{R}, q > 0$ tal que:

$$d^2[S_{(n-k)}^{(i)}, S_{(n)}] \leq q \quad \forall i \quad i = 1, 2, \dots, \binom{n}{k}$$

y por tanto

$$\max_i d^2[S_{(n-k)}^{(i)}, S_{(n)}] \leq q$$

Evidentemente si $\max_i d^2[S_{(n-k)}^{(i)}, S_{(n)}] > q$ esto se interpretaría afirmando que los n elementos que componen la muestra no pertenecen a la misma población. Por otro lado los k elementos que se han suprimido de la muestra para obtener $S_{(n-k)}$ serían los elementos que no pertenecen a la población $N_2(\mu, \Sigma)$ ya que al añadir estos k elementos a la muestra de tamaño $n - k$ se produce una dispersión que hace que la distancia entre dichas matrices sobrepase a esa cota q . Estos k elementos serán los posibles outliers para esta distribución $N_2(\mu, \Sigma)$

Ya que los estadísticos $\{d^2(S_{(n-k)}^{(i)}, S_{(n)}), i = 1, 2, \dots, \binom{n}{k}\}$ no son independientes no podemos calcular la distribución del estadístico.

ordenado $\max_i d^2(S_{(n-k)}^{(i)}, S_{(n)})$. Por ello recurriremos a hallar una cota superior para la cola de dicha distribución, es decir acotaremos

$$P[\max_i d^2(S_{(n-k)}^{(i)}, S_{(n)}) > q]$$

Si definimos E_i como el suceso $\{d^2(S_{(n-k)}^{(i)}, S_{(n)}) > q\}$ para $i=1, 2, \dots, \binom{n}{k}$.

entonces

$$\begin{aligned} P[\max_i d^2(S_{(n-k)}^{(i)}, S_{(n)}) > q] &= P\left[\bigcup_{i=1}^{\binom{n}{k}} E_i\right] \leq \sum_{i=1}^{\binom{n}{k}} P(E_i) = \\ &= \binom{n}{k} P(E_i) = \binom{n}{k} P[d^2(S_{(n-k)}^{(i)}, S_{(n)}) > q] \end{aligned}$$

Sin embargo nos interesa trabajar con la variable aleatoria X con función de densidad la dada en el teorema 4.3.2., cuya distribución que es dominada por la del estadístico $d^2(S_{(n-k)}, S_{(n)})$, hemos tabulado en el apartado 4.3.

Ya que

$$P[X \leq q] \leq P[d^2(S_{(n-k)}, S_{(n)}) \leq q] \quad \forall q \quad 0 \leq q \leq \infty$$

se deduce

$$P[X > q] \geq P[d^2(S_{(n-k)}, S_{(n)}) > q]$$

y en definitiva

$$\begin{aligned} P[\max_i d^2(S_{(n-k)}^{(i)}, S_{(n)}) > q_\alpha] &\leq \binom{n}{k} P[d^2(S_{(n-k)}, S_{(n)}) > q_\alpha] \leq \\ &\leq \binom{n}{k} P[X > q_\alpha] = \alpha. \end{aligned}$$

En definitiva el proceso a seguir sería el siguiente:

Dado un nivel de significación α se determinaría el cuantil q_α mediante la tabulación que se ha hecho de la distribución de la variable aleatoria X .

A continuación calcularíamos el valor del estadístico $\max_i d^2(S_{(n-k)}^{(i)}, S_{(n)})$ y la regla de decisión sería la siguiente.

Si

$$\max_i d^2(S_{(n-k)}^{(i)}, S_{(n)}) > q_\alpha \quad i = 1, 2, \dots, \binom{n}{k}$$

los n elementos que componen la muestra no pertenecen a la misma población, siendo considerados como outliers los k elementos que se han suprimido de la muestra para obtener $S_{(n-k)}$.

En las paginas siguientes incluimos un programa de ordenador para calcular el valor del estadístico ordenado $\max_i d^2(S_{(n-k)}^{(i)}, S_{(n)})$ y obtener los k elementos de la muestra excluidos en el calculo de dicho estadístico y que determinaríamos si son outliers mediante la regla de decisión dada anteriormente.



Calculo del estadístico máxima distancia y detección de los posibles outliers.

```

1*      PARAMETER N=15,N4=N-4,N3=N-3,N2=N-2,N1=N-1
2*      DIMENSION DATOS(2,N),CIN,N),PR(2,N),S(2,2),SINV(2,2),VI(2,2),
3*      *X1(2,N4),CN4(N4,N4),PR2(2,N4),S2(2,2),A(2,2)
4*      C -----
5*      C LECTURA DE UN DATO BIDIMENSIONAL EN CADA TARJETA E IMPRESION
6*      C -----
7*      PFAO 1, ((DATOS(I,J),I=1,2),J=1,N)
8*      PRINT1, ((DATOS(I,J),I=1,2),J=1,N)
9*      1 FORMAT(2F20.2)
10*     PRINT 23
11*     23 FORMAT(/,/)
12*     DI=0
13*     JI=0
14*     C -----
15*     C CALCULO DE LA MATRIZ DE SUMA DE CUADRADOS Y PRODUCTOS DE N OBSERVACIONES SN
16*     C (EN EL PROGRAMA, LA MATRIZ SN ES S)
17*     C -----
18*     DO 2 I=1,N
19*     DO 2 J=1,N
20*     2 C(I,J)=-1./N
21*     DO 3 I=1,N
22*     3 C(I,I)=C(I,I)+1
23*     DO 4 I=1,2
24*     DO 4 J=1,N
25*     PR(I,J)=0
26*     DO 4 K=1,N
27*     4 PR(I,J)=PR(I,J)+DATOS(I,K)*C(K,J)
28*     DO 5 I=1,2
29*     DO 5 J=1,2
30*     S(I,J)=0
31*     DO 5 K=1,N
32*     5 S(I,J)=S(I,J)+PR(I,K)*DATOS(J,K)
33*     C -----
34*     C IMPRESION DE LA MATRIZ SN
35*     C -----
36*     PRINT 301
37*     301 FORMAT(' MATRIZ SN',/)
38*     PRINT 300, ((S(I,J),I=1,2),J=1,2)
39*     300 FORMAT(2E10.9)
40*     PRINT 23
41*     C -----
42*     C COMPROBACION DE QUE LA MATRIZ SN ES INVERTIBLE Y SIMETRICA, E IMPRESION ERROR
43*     C -----
44*     DET=S(1,1)*S(2,2)-S(1,2)*S(2,1)
45*     IF((S(1,2)).EQ.(S(2,1))) GOTO 200
46*     PRINT 201, S(1,2),S(2,1)
47*     201 FORMAT(' S12='E16.9,' S21='E18.9)
48*     200 IF(DET.EQ.0) STOP1
49*     C -----
50*     C CALCULO DE LA MATRIZ SN(-1), INVERSA DE SN, E IMPRESION
51*     C (EN EL PROGRAMA, SN(-1) ES SINV)
52*     C -----
53*     SINV(1,1)=S(2,2)/DET
54*     SINV(2,2)=S(1,1)/DET
55*     SINV(1,2)=-S(1,2)/DET
56*     SINV(2,1)=-S(2,1)/DET
57*     PRINT 302
58*     302 FORMAT(' MATRIZ SN(-1)',/)
59*     PRINT 300, ((SINV(I,J),I=1,2),J=1,2)
60*     PRINT 23
61*     C -----
62*     C ENUMERACION UNO A UNO DE TODAS LAS POSIBLES CUATERNAS
63*     C -----
64*     DO 6 I1=1,N3
65*     I11=I1+1
66*     DO 6 I2=I11,N2
67*     I21=I2+1
68*     DO 6 I3=I21,N1
69*     I31=I3+1
70*     DO 6 I4=I31,N
71*     IO=0
72*     DO 9 J=1,N
73*     IF((J-I1)*(J-I2)*(J-I3)*(J-I4) .EQ.0) GOTO7
74*     DO 8 M=1,2
75*     8 X1(M,J-IO)=DATOS(M,J)
76*     GOTO 9
77*     7 IO=IO+1
78*     9 CONTINUE
79*     C -----

```

```

80* C CALCULO DE LA MATRIZ DE SUMA DE CUADRADOS Y PRODUCTOS DE N=4 OBSERVACIONES SN4
81* C (EN EL PROGRAMA LA MATRIZ SN4 ES S2)
82* -----
83* DO 11 I=1,N4
84* DO 11 J=1,N4
85* 11 CN4(I,J)=-1./N4
86* DO 12 I=1,N4
87* 12 CN4(I,I)=CN4(I,I)+1
88* DO 13 I=1,2
89* DO 13 J=1,N4
90* PR2(I,J)=0
91* DO 13 K=1,N4
92* 13 PR2(I,J)=PR2(I,J)+X1(I,K)*CN4(K,J)
93* DO 14 I=1,2
94* DO 14 J=1,2
95* S2(I,J)=0
96* DO 14 K=1,N4
97* 14 S2(I,J)=S2(I,J)+PR2(I,K)*X1(J,K)
98* -----
99* C CALCULO DE LA MATRIZ SN4*SN(-1)
100* C (EN EL PROGRAMA LA MATRIZ PRODUCTO SN4*SN(-1) ES A)
101* -----
102* DO 15 I=1,2
103* DO 15 J=1,2
104* A(I,J)=0
105* DO 15 K=1,2
106* 15 A(I,J)=A(I,J)+S2(I,K)*SIN(K,J)
107* -----
108* C CALCULO DE LOS AUTOVALORES
109* -----
110* RAD=SQRT((A(1,1)-A(2,2))**2+.4.*A(2,1)*A(1,2))
111* Y1=(A(1,1)+A(2,2)+RAD)/2.
112* Y2=(A(1,1)+A(2,2)-RAD)/2.
113* -----
114* C COMPROBACION DE QUE LOS AUTOVALORES ESTAN EN (0,1)
115* C
116* IF((0.LT.Y1).AND.(Y1.LT.1)) GOTO 16
117* PRINT 17,Y1,I1,I2,I3,I4
118* 17 FORMAT(' Y1 NO ESTA EN (0,1). Y1=*F15.9*X.4I9./)
119* GOTO 63
120* 16 IF((0.LT.Y2).AND.(Y2.LT.1)) GOTO 18
121* PRINT 19,Y2,I1,I2,I3,I4
122* 19 FORMAT(' Y2 NO ESTA EN (0,1). Y2=*F15.9*X.4I9./)
123* GOTO 60
124* -----
125* C CALCULO DEL ESTADISTICO MAXIMA DISTANCIA E IMPRESION DE LOS DATOS EXCLUIDOS
126* C (POSIBLES OUTLIERS)
127* C J1 ES EL CONTADOR DEL NUMERO DE COMPARACIONES
128* C
129* 18 D2=((ALOG(Y1))**2+(ALOG(Y2))**2)
130* IF(D2.LT.D1) GOTO 60
131* I11=I1
132* I22=I2
133* I33=I3
134* I44=I4
135* D1=D2
136* 60 J1=J1+1
137* IF(((J1/1000.)-INT(J1/1000)).GT.(0.000001)) GOTO 6
138* PRINT 20,J1,I1,I2,I3,I4, D1,I11,I22,I33,I44
139* 20 FORMAT(' IYEPACION*I7.* EN* 4I7.20X.* MAX. DISTANCIA PARCIAL =*
140* *F15.9*X EN*4I7./)
141* E CONTINUE
142* PRINT 20,D1

```

```

14 3*      20 FORMAT( ' MAXIMA DISTANCIA =*F15.9. /)
14 4*      PRINT 21, DATOS(1,I111),DATOS(1,I222),DATOS(1,I333),DATOS(1,I444)
14 5*      PRINT 22, DATOS(2,I111),DATOS(2,I222),DATOS(2,I333),DATOS(2,I444)
14 6*      21 FORMAT( ' DATOS EXCLUIDOS EN EL MAXIMO* 10X*4F15.9)
14 7*      22 FORMAT(39X*4F15.9)

```

```

C -----
C GRAFICA DE LOS PUNTOS
C -----

```

```

15 1*      PRINT 103
15 2*      103 FORMAT(1H1)
15 3*      DO 101 I=1,61
15 4*      DO 105 J=1,121
15 5*      105 V(J)=* *
15 6*      V(61)=*I*
15 7*      IF(I.NE.41) GO TO 106
15 8*      DO 107 J=1,121
15 9*      107 V(J)=*- *
16 0*      106 DO 102 J=1,N
16 1*      IF((DATOS(2,J).GE.(20.25-.5*I)).AND.(DATOS(2,J).LT.(20.75-.5*I)))
16 2*      * GO TO 104
16 3*      GO TO 102
16 4*      104 J1=(10+DATOS(1,J))*121/20
16 5*      V(J)=* *
16 6*      102 CONTINUE
16 7*      PRINT 108,V
16 8*      108 FORMAT(1X,121A1)
16 9*      101 CONTINUE
17 0*      STOP
17 1*      END

```

END OF COMPILATION: NO DIAGNOSTICS.

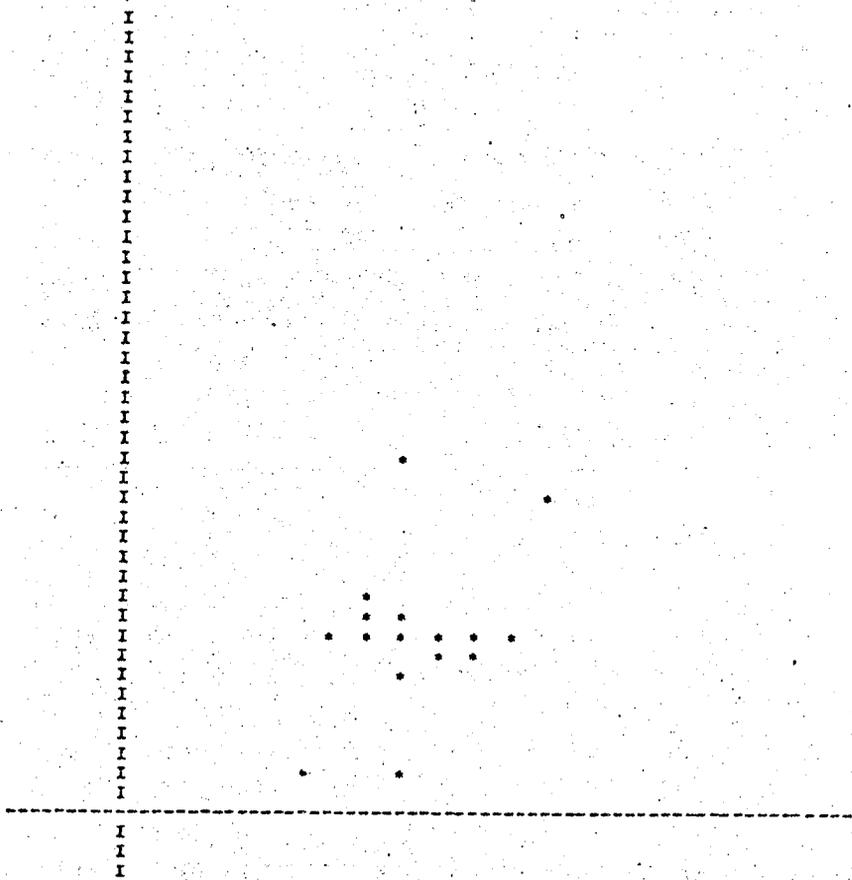
4.5. APLICACION.

Se ha extraído una muestra de tamaño quince de las calificaciones obtenidas en teoría y problemas en una asignatura de tercer curso de la licenciatura de Matemáticas, con los siguientes resultados.

5.00	4.00
6.00	8.00
3.00	4.50
4.00	3.50
3.50	5.50
4.00	9.00
4.50	4.50
4.00	5.00
5.50	4.50
4.00	4.50
4.50	4.00
4.00	1.00
3.50	4.50
5.00	4.50
3.50	5.00

Representando en el eje de abscisas las calificaciones de teoría y en el de ordenadas las de problemas, la nube de puntos que se obtiene es de la forma





Se quieren detectar las calificaciones que destacan de entre las de la muestra, es decir las observaciones outliers aplicando para ello el metodo de la distancia entre las matrices de sumas de cuadrados y sumas de productos de observaciones muestrales desarrollado en el apartado anterior.

Mediante el programa dado en las pags. 93 a 95 , para $k = 3$ se obtiene:

MAXIMA DISTANCIA = 8.6990

DATOS EXCLUIDOS EN EL MAXIMO:

6.00	8.00
4.00	9.00
4.00	1.00

Para un nivel de significación $\alpha = 0.01$ el percentil de la distribución del estadístico $\max_{\lambda} d^2(S_{(n)}, S_{(n-k)}^{(\lambda)})$, que viene dado mediante el programa de ordenador de las págs. 88 y 89 es 8.5.

Como la máxima distancia experimental es mayor que este valor crítico, se tendrá que para los datos excluidos en el máximo, se rechaza la hipótesis de que las tres observaciones pertenecen a la misma población y por tanto serían consideradas outliers.

A continuación comprobamos si pueden existir cuatro observaciones outliers. Mediante el programa de ordenador dado en las páginas 93, 94 y 95, obtenemos.

MAXIMA DISTANCIA ≈ 12.6026

DATOS EXCLUIDOS EN EL MAXIMO

6.00	8.00
4.00	9.00
4.00	1.00
4.00	3.50

Para un nivel de significación $\alpha = 0.01$ el valor teórico del estadístico es 13.6 y como el experimental es menor rechazamos la hipótesis de que esas cuatro observaciones son outliers.

REFERENCIAS BIBLIOGRAFICAS.

- ABRAMOWITZ M. and STEGUN I. A. (1.972). Handbook of mathematical functions. Dover.
- AITKEN A. C. (1.965). Determinantes y Matrices. Dossat.
- ANDERSON T. W. (1.958). An Introduction to Multivariate Statistical Analysis. Wiley.
- ANDERSON T. W. and BAHADUR R.R. (1.962). "Classification into two multivariate normal distributions with different covariance matrices". Annals of Mathematical Statistics (33) .
- ANDREWS D. F. (1.971). "Significance tests based on residuals" Biometrika (58).
- ANDREWS D. F. (1.972). "Plots of high-dimensional data". Biometric (28) .
- ANSCOMBE F. J. (1.960). "Rejection of outliers". Technometrics (2).
- ANSCOMBE F. J. and TUKEY J. W. (1.963). "The examination and analysis of residuals". Technometrics (5).
- BARNETT V. and LEWIS T. (1.978). Outliers in Statistical Data. Wiley.

- BELLMAN R. (1.965). *Introduccion al Analisis Matricial*. Reverté.
- BROWN B. M. (1.975). "A short-cut test for outliers using residuals". *Biometrika* (62).
- CHERNOFF H. (1.972). "The selection of effective attributes for deciding between hypotheses using linear discriminant functions". In *Frontiers of Pattern Recognition*, Watanabe (Ed.). Academic Press.
- CHERNOFF H. (1.973). "Some measure for discriminating between normal multivariate distributions with unequal covariance matrices". In *Multivariate Analysis III*, Krishnaiah (Ed.). Academic Press.
- COLLET D. (1.980). "Outliers in circular data". *Applied Statistique* (29).
- COX D. R. and SNELL E. J. (1.968). "A general definition of residuals". *Journal of the Royal Statistical Society B* (30).
- DAVID H. A., HARTLEY H. O. and PEARSON E. S. (1.954). "The Distribution of the ratio, in a single normal sample of range of standard deviation". *Biometrika* (41)
- DE RAULY D. (1.966). *L'estimation Statistique*. Gauthier Villars.
- DEVLIN S. J.; GNANADESIKAN R. and KETTERING J. R. (1.975). "Robust estimation and outlier detection with correlation coefficients". *Biometrika* (62)

- DIXON W. J. (1.950). "Analysis of extreme values". Annals of Mathematical Statistics (21).
- FERGUSON T. S. (1.961). "On the rejection of outliers". Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability, Vol I.
- FISHER R. A. (1.936). "The use of multiple measurements in taxonomic problems". Annals of Eugenics (7).
- FISHER R. A. (1.939). "The sampling distribution of some statistics obtained from non-linear equations". Annals of Eugenics (9)
- GIRI N. C. (1.977). Multivariate Statistical Inference. Academic Press.
- GNANADESIKAN R. (1.973). "Graphical methods for informal inference in multivariate data analysis". Bull. Int. Statist. Inst. (45).
- GNANADESIKAN R. (1.977). Methods for Statistical Data Analysis of Multivariate Observations. Wiley.
- GNANADESIKAN R. and KETTERING J. R. (1.972). "Robust estimates residuals, and outlier detection with multiresponse data". Biometrics (28)
- GRUBBS F. E. (1.950). "Sample criteria for testing outlying observations". Annals of Mathematical Statistics (21)

- GRUBBS F. E. (1.969). "Procedures for detecting outlying observations in samples". *Technometrics* (11).
- GUMBELL E. J. (1.960). "Discussion on rejection of outliers by Anscombe". *Technometrics* (2).
- GUPTA S. S. (1.960). "Order Statistics from the Gamma distribution". *Technometrics* (2).
- HAMPEL F. R. (1.974). "The influence curve and its role in robust estimation". *Journal of the the American Statistical Association* (69).
- HAWKINS D. M. (1.974). "The detection of errors in multivariate data using principal components". *Journal of the American Statistical Association* (69).
- HSU P. L. (1.939). "On the distribution of roots of certain determinantal equations". *Annals of Eugenics* (9).
- IRWIN J. O. (1.925). "On a criterion for the rejection of outlying observations". *Biometrika* (17).
- JOHNSON N. L. and KOTZ S. (1.969). *Distributions in Statistics; Discrete Distributions*. Wiley.
- JOHNSON N. L. and KOTZ S. (1.970). *Distributions in Statistics; Continuous univariate distributions, I y II*. Wiley
- JOHNSON N. L. and KOTZ S. (1.972). *Distributions in Statistics; Continuous Multivariate Distributions*. Wiley.

- KENDALL M. G. and BUCKLAND W. R. (1.957). A Dictionary of Statistical Terms. Longman.
- KOBAYASHI S. and NOMIZU K. (1.963). Foundation of Differential Geometry, Vol. I. Wiley Interscience.
- KRISHNAIAH P. R. (1.978). "Some recent developments on real multivariate distributions". In Developments in Statistics, Krishnaiah (Ed.). Academic Press.
- KSHIRSAGAR A. M. (1.972). Multivariate Analysis. Marcel Dekker.
- KUDO A. (1.956). "On the testing of outlying observations". Sankhya (17).
- MAAS H. (1.955). "Die bestimmung der dirichletreihen mit grossencharakteren zu den modulformen n-ten grades". J. Indian Math. Soc. (19).
- NAIR K. R. (1.948). "The distribution of the extreme deviate from the sample mean and its studentized form". Biometrika (35).
- PATEL J. K.; KAPADIA C. H. and OWEN D. B. (1.976). Handbook of Statistical Distributions. Marcel Dekker.
- PATNAIK P. B. (1.949). "The non-central χ^2 and F-distributions and their applications". Biometrika (36).
- PEARSON E. S. and CHANDRASEKAR C. (1.936). "The efficiency of statistical tools and a criterion for the rejection of outlying observations". Biometrika (28).



- PRESS S. J. (1.966). "Linear combinations of non central chi square variates". Annals of Mathematical Statistics (37).
- ROHATGI V. K. (1.976): An Introduction to Probability Theory and Mathematical Statistics. Wiley.
- ROSNER B. (1.975). "On the detection of many outliers". Technometrics (11).
- ROY S. N. (1.939). "P-statistics or some generalizations in an analysis of variance appropriate to multivariate problems". Shankya (4).
- SHAH S. M. (1.966). "On estimating the parameter of a doubly truncated binomial distribution". Journal of the American Statistical Association (61).
- SHAPIRO S. S. and WILK M. B. (1.972). "An analysis of variance test for the exponential distribution". Technometrics (14).
- SIOTANI M. (1.959). "The extreme value of the generalised distances of the individual points in the multivariate normal sample". Ann. Inst. Statist. Math. Tokyo (10).
- STUDENT (1.927). "Errors of routine analysis". Biometrika (19).

- TIAO G. C. and GUTTMAN I. (1.967). "Analysis of outliers with adjusted residuals". Technometrics (9).
- TIETJEN G. L. and MOORE R. H. (1.972). "Some Grubbs-Type Statistics for the detection of several outliers". Technometrics (14).
- THOMPSON W.R. (1.935). "On a criterion for the rejection of observations and the distribution of the ratio of the deviation to the sample standard deviation". Annals of Mathematical Statistics (6).
- WATSON G. W. (1.964). "A note on maximum likelihood". Shankya (26)
- WILK M. B. GNANADESIKAN R. and HUYETT J. (1.962). "Probability plots for the gamma distribution". Technometrics (4)
- WILKS S. S. (1.962). Mathematical Statistics. Wiley.
- WILKS S. S. (1.963). "Multivariate Statistical outliers". Shankya A, (25).

