

03061
1ej.
1
MAY 23 1983

UNIVERSIDAD NACIONAL
AUTONOMA DE MEXICO

UNIDAD ACADÉMICA DE LOS CICLOS PROFESIONAL
Y DE POSGRADO DEL C C H

"MEZCLAS DE DISTRIBUCIONES NORMALES"
MULTIVARIADAS

T E S I S
MAESTRIA EN ESTADISTICA E INVESTIGACION DE
OPERACIONES

JAVIER ALAGON CANO

TESIS CON
TALLA DE ORIGEN

JULIO 1983



Universidad Nacional
Autónoma de México



UNAM – Dirección General de Bibliotecas
Tesis Digitales
Restricciones de uso

DERECHOS RESERVADOS ©
PROHIBIDA SU REPRODUCCIÓN TOTAL O PARCIAL

Todo el material contenido en esta tesis esta protegido por la Ley Federal del Derecho de Autor (LFDA) de los Estados Unidos Mexicanos (México).

El uso de imágenes, fragmentos de videos, y demás material que sea objeto de protección de los derechos de autor, será exclusivamente para fines educativos e informativos y deberá citar la fuente donde la obtuvo mencionando el autor o autores. Cualquier uso distinto como el lucro, reproducción, edición o modificación, será perseguido y sancionado por el respectivo titular de los Derechos de Autor.

INDICE

CAPITULO I. INTRODUCCION	1
CAPITULO II. MAXIMA VEROSIMILITUD EN MEZCLAS DE DISTRIBUCIONES MULTIVARIADAS	6
1. Definición y cálculo de estimadores	6
2. Métodos Iterativos	12
CAPITULO III. MEZCLAS DE DISTRIBUCIONES NORMALES	17
1. Caso General	18
2. Matrices de Dispersion Iguales	24
3. Caso	29
4. Matrices de Correlación Iguales	34
CAPITULO IV. PRUEBAS SOBRE EL NUMERO DE GRUPOS EN UNA MEZCLA	41
1. Revisión de algunos métodos	41
2. Prueba del cociente de verosimilitud	44
CAPITULO V. USO DE NORMIX	47
1. Características generales y manual de usuario	47
2. Agrupamiento de los datos de Iris de Fisher	
BIBLIOGRAFIA	

CAPITULO I

INTRODUCCION

El término más usado para referirse a las técnicas estadísticas que tienen por objetivo separar un conjunto de datos y formar distintos grupos es el de "Análisis de Clasificación". Muchos autores usan este término para referirse también a técnicas que buscan agrupar variables.

Dentro del Análisis de Clasificación se encuentran el Análisis Discriminante y el Análisis de Conglomerados*. En el primero de éstos se conoce de antemano el tipo de grupos que se formarán, así como el número de éstos; en el segundo se desconocen a priori ambas características. El Análisis Discriminante es usado también para estudiar las características principales que hacen diferente a un grupo de individuos de otro; en base a esto, se puede asignar un nuevo individuo en alguno de los grupos. Por su parte, el Análisis de Conglomerados tiene como principal objetivo, simplemente el formar grupos de individuos de acuerdo a un conjunto de características medidas en éstos.

Este trabajo trata con una de las técnicas de Análisis de Conglomerados, específicamente la técnica que consis-

* En Inglés: Cluster Analysis.

te en ajustar mezclas de distribuciones multivariadas a un conjunto de datos multivariados. El objetivo de esta técnica es resolver el siguiente problema: dada una muestra de N objetos o individuos a cada uno de los cuales se han medido p variables, agrupar estos objetos o individuos en g grupos o clases distintas. El número de estos grupos no se conoce de antemano. Esta técnica de análisis de conglomerados es de tipo no jerárquico.

Clasificar gente en distintos tipos ha sido para el ser humano todo un pasatiempo desde épocas muy remotas. Los antiguos hindúes usaban características físicas y de comportamiento así como el sexo de las personas para clasificar a la gente en seis tipos distintos que ellos designaban con nombres de animales. Los antiguos fisiólogos griegos y romanos desarrollaron varias tipologías basadas en variaciones de características físicas resultantes de una mezcla de cuatro humores. La más destacada de estas clasificaciones fue la de Galeno (129-199 d.C.), que definió nueve tipos de personas de acuerdo a su temperamento y que se suponía se relacionaba con la susceptibilidad de la persona a distintas enfermedades y con las diferencias individuales en el comportamiento. En el siglo XVIII, Linceo hizo la primera clasificación científica dentro de los reinos animal y vegetal.

La mayoría del trabajo de clasificación se hizo

primero en los campos de la Botánica y la Zoología en donde se le conoció por Taxonomía. Esta fue en un principio, más un arte que una técnica científica, pero gradualmente se fueron desarrollando técnicas más objetivas hasta llegar a los métodos de la Taxonomía Numérica. Posteriormente, son particularmente importantes los intentos de Zubin en 1938 y de Thorndike en 1953 al usar técnicas de clasificación numérica en campos muy alejados de los de las ciencias naturales. Sin embargo, estas técnicas tuvieron que esperar hasta el desarrollo de las computadoras electrónicas puesto que incluyen una gran cantidad de cálculos numéricos.

Desde hace dos décadas, tanto estadísticos como matemáticos han estado trabajando en un enfoque más formal al análisis de conglomerados y se han hecho muchos intentos para formular modelos estadísticos y matemáticos acordes con los problemas de clasificación.

La mayoría de los procedimientos de análisis de conglomerados tratan de medir la similitud entre cualesquiera dos objetos y entonces agrupan los objetos maximizando la similitud dentro de grupos. Desafortunadamente, no existe una "medida apropiada" de similitud, pues esta depende del problema bajo investigación, por lo que forzosamente se cae en esquemas subjetivos de clasificación. Para cierto tipo de datos puede ser muy conveniente una técnica de análisis de conglomerados sin suposiciones arbitrarias de similitud. La

técnica de ajustar mezclas de distribuciones a un conjunto de datos, es de este tipo, puesto que no incluye medidas de similitud subjetivas.

El principal objetivo de esta tesis consistió en resolver un problema particular en la técnica de mezclas. Específicamente, el problema consistía en encontrar los estimadores máximo verosímiles de los parámetros involucrados en una mezcla de distribuciones normales multivariadas (mediante los cuales podemos identificar los distintos grupos existentes en los datos) en el caso particular en que las matrices de dispersión de los grupos son del tipo $\Sigma_s = a_s \Sigma$ donde Σ es una matriz positiva definida y a_s es un número real positivo, $s=1, 2, \dots, g$ y g es el número de grupos a formar. Tal y como se demuestra en el Capítulo III de esta tesis, este caso es un caso particular de la situación en donde se tienen matrices de correlación iguales en los grupos. Este caso se presenta frecuentemente en problemas de Biología, en donde se pueden tener grupos en donde las varianzas y covarianzas de las variables medidas cambian proporcionalmente en los distintos grupos; por ejemplo, en zoología podemos encontrarnos con una situación en donde al clasificar diferentes especies animales de acuerdo a un conjunto de medidas sobre varios órganos del cuerpo, estas difieren en medias y varianzas para distintos grupos aunque se tengan las mismas correlaciones en los distintos grupos. Desde luego

que el caso en que se tienen matrices de correlación iguales ya estaba resuelto [Ver #1]. Sin embargo, simplificar el problema tiene muchas ventajas estadísticas y numéricas, puesto que se tiene un menor número de parámetros por estimar, lo cual no hace necesario un tamaño de muestra muy grande para calcular los estimadores de las varianzas y covarianzas en los grupos; además se simplifican los cálculos numéricos dentro de la computadora y la convergencia de los métodos iterativos usados es más rápida.

La simplificación anterior es semejante a la que se hace cuando en el caso general se suponen matrices de dispersión iguales. Estos fueron resueltos y presentados por J.H. Wolfe en 1967 [Ver 18]. El mismo Wolfe implantó estos casos en computadora en un programa al que tituló NORMIX (del Inglés NORmal MIXtures). Este programa se encuentra en ITHAS. En el último capítulo de este trabajo presentamos la manera de usar NORMIX.

CAPITULO II

MAXIMA VEROSIMILITUD EN MEZCLAS
DE DISTRIBUCIONES MULTIVARIADAS

Tal y como lo señalamos en el capítulo anterior, este trabajo tratará con una técnica para el análisis de datos multivariados, cuyo objetivo es resolver el siguiente problema: Dada una muestra de N objetos o individuos a cada uno de los cuales se le han medido p variables, agrupar estos objetos o individuos en g grupos o clases distintas. El número de estos grupos no se conoce de antemano.

El objetivo central de este capítulo es presentar la técnica de máxima verosimilitud en mezclas de distribuciones multivariadas para resolver el problema anterior. Presentamos también algunos métodos iterativos para la resolución de las ecuaciones máximo verosímiles.

§ 1. DEFINICION Y CALCULO DE ESTIMADORES.

Dado que los objetos o individuos dentro de un grupo difieren uno de otro, es razonable suponer la existencia de una distribución de probabilidad de las características de una población que pertenece a un grupo. Así, los individuos

de grupos diferentes tendrán distintas distribuciones de probabilidad de sus características. La población combinada tomada de todos los grupos tendrá una distribución de probabilidad conocida como una mezcla de distribuciones. El problema de análisis de conglomerados consiste en identificar y describir las distribuciones componentes de una mezcla a partir de la muestra. Usualmente se supone que las distribuciones componentes son distribuciones estadísticas estándar con parámetros desconocidos. El problema de agrupación se resuelve entonces con técnicas estadísticas estándar de estimación paramétrica y asignando a cada individuo al grupo al grupo al cual tenga mayor probabilidad de pertenecer. Este trabajo trata única y exclusivamente, estimación por máxima verosimilitud.

Para ello, sean $\alpha_1(x, \Theta_1), \dots, \alpha_g(x, \Theta_g)$, g distribuciones de probabilidad definidas en un espacio de dimensión p de vectores aleatorios

$$x = (x_1, \dots, x_p)$$

Supongamos además, que α_s es una función dos veces derivable en sus parámetros

$$\Theta_s = (\theta_{s1}, \dots, \theta_{sq}), \quad s = 1, \dots, g; \quad q \in \mathbb{N}$$

Formamos una mezcla de distribuciones tomando proporciones λ_s de la población del grupo s , es decir, la distribución de probabilidad de la mezcla está dada por

$$(2.1) \quad f(\underline{x}) = \sum_{s=1}^g \lambda_s \alpha_s(\underline{x}, \Theta_s)$$

donde

$$(2.2) \quad \sum_{s=1}^g \lambda_s = 1.$$

La probabilidad de que el vector \underline{x} pertenezca al grupo s está dada por

$$(2.3) \quad P(s/\underline{x}) = \frac{P(s) P(\underline{x}/s)}{P(\underline{x})} = \frac{\lambda_s \alpha_s(\underline{x}, \Theta_s)}{f(\underline{x})}$$

Supongamos que obtenemos una muestra de N vectores aleatorios. El vector k -ésimo lo representaremos por

$$\underline{x}_k = (x_{1k}, \dots, x_{pk}), \quad k = 1, \dots, N.$$

Los estimadores máximo verosímiles de los parámetros involucrados en la mezcla, λ_s, Θ_s , $s = 1, \dots, g$, son aquellos valores de λ_s, Θ_s que maximizan la verosimilitud de la muestra

$$\log L = \sum_{k=1}^N \log f(\underline{x}_k)$$

sujeta a la restricción

$$\sum_{s=1}^g \lambda_s = 1$$

Usando un multiplicador de Lagrange γ , formamos la función objetivo

$$L = \sum_{k=1}^N \log f(x_k) - \gamma \left(\sum_{s=1}^g \lambda_s - 1 \right)$$

Las ecuaciones máximo verosímiles son

$$\frac{\partial L}{\partial \lambda_s} = \sum_{k=1}^N \frac{\kappa_s(x_k, \theta_s)}{f(x_k)} - \gamma = 0$$

$$\frac{\partial L}{\partial \theta_{si}} = \sum_{k=1}^N \frac{\lambda_s}{f(x_k)} \frac{\partial \kappa_s}{\partial \theta_{si}} = 0$$

Si multiplicamos la primera de estas ecuaciones por λ_s y sustituimos en (2.3) obtenemos que

$$\sum_{k=1}^N \frac{\lambda_s \kappa_s(x_k, \theta_s)}{f(x_k)} - \gamma \lambda_s = 0$$

es decir

$$\sum_{k=1}^N P(s/x_k) - \gamma \lambda_s = 0$$

Si ahora sumamos sobre s concluimos que

$$\sum_{s=1}^g \sum_{k=1}^N P(s/x_k) - \gamma \sum_{s=1}^g \lambda_s = 0$$

i.e.,

$$\gamma = \sum_{s=1}^g \sum_{k=1}^N P(s/x_k) = \sum_{k=1}^N \sum_{s=1}^g \frac{\lambda_s \kappa_s(x_k, \theta_s)}{f(x_k)} = \sum_{k=1}^N \frac{f(x_k)}{f(x_k)} = N$$

Podemos obtener más información de estas ecuaciones máximo verosímiles, como presentamos en los siguientes teoremas.

Teorema 1. La estimación por máxima verosimilitud de las proporciones de una mezcla es igual a la media muestral de las probabilidades de pertenencia al grupo dado, es decir

$$(2.4) \quad \hat{\lambda}_s = \frac{1}{N} \sum_{k=1}^N P(s/Y_k)$$

Demostración. Tenemos que
$$\sum_{k=1}^N P(s/Y_k) - \gamma \lambda_s = 0$$

Dado que $\gamma = N$, concluimos que

$$\hat{\lambda}_s = \sum_{k=1}^N \hat{P}(s/Y_k). \quad \square$$

Teorema 2. Las ecuaciones máximo verosímiles para estimar los parámetros de las distribuciones de una mezcla están dadas por

$$(2.5) \quad \frac{\partial L}{\partial \theta_{si}} = \sum_{k=1}^N P(s/Y_k) \frac{\partial \log \kappa_s(Z_k, \Theta_s)}{\partial \theta_{si}} = 0$$

Demostración. Tenemos que
$$P(s/Y_k) = \frac{\lambda_s \kappa_s(Z_k, \Theta_s)}{f(Z_k)}$$

Por lo tanto,

$$\begin{aligned} \frac{\partial L}{\partial \theta_{si}} &= \sum_{k=1}^N \frac{\lambda_s}{f(Z_k)} \frac{\partial \kappa_s(Z_k, \Theta_s)}{\partial \theta_{si}} = \sum_{k=1}^N \frac{P(s/Y_k)}{\kappa_s(Z_k, \Theta_s)} \frac{\partial \kappa_s(Z_k, \Theta_s)}{\partial \theta_{si}} = \\ &= \sum_{k=1}^N P(s/Y_k) \frac{\partial \log \kappa_s(Z_k, \Theta_s)}{\partial \theta_{si}} = 0. \quad \square \end{aligned}$$

Notemos que si toda la población fuera obtenida de un solo grupo, las ecuaciones máximo verosímiles para θ_s serían

$$\sum_{k=1}^N \frac{\partial \log \kappa_s}{\partial \theta_{s_i}} = 0$$

Por lo tanto, las ecuaciones máximo verosímiles para las estimaciones de los parámetros de una mezcla, son los promedios ponderados de las expresiones usadas al obtener las estimaciones máximo verosímiles para grupos puros, donde las ponderaciones son las probabilidades de pertenencia a los grupos.

Usualmente, el número de grupos g , es una hipótesis que debe probarse contra la alternativa de g' y esto se hace encontrando las estimaciones por máxima verosimilitud bajo ambas hipótesis y probando el cociente de verosimilitud con la fórmula

$$\chi^2 = -2 \log \left[\frac{k_g}{k_{g'}} \right]$$

con grados de libertad igual a la diferencia de los números de parámetros a estimar bajo las dos hipótesis. Por supuesto que la distribución del logaritmo del cociente de verosimilitud es aproximadamente χ^2 para muestras grandes solamente. Esto se verá con más detalle en el capítulo IV.

Notemos que las ecuaciones (2.4) y (2.5) no están en forma cerrada, por lo que tienen que ser resueltas numéricamente. Para ello, veremos algunos métodos iterativos de resolución de estas ecuaciones.

2. METODOS ITERATIVOS

En la mayoría de los casos, las ecuaciones máximo verosímiles

$$(2.4) \quad \frac{1}{\lambda_s} \sum_{k=1}^N \hat{p}(s/x_k) - N = 0$$

y

$$(2.5) \quad \frac{\partial L}{\partial \theta_{sj}} = \sum_{k=1}^N \hat{p}(s/x_k) \frac{\partial \log \alpha(s/x_k, \theta_{sj})}{\partial \theta_{sj}} = 0$$

tienen que ser resueltas numéricamente. La manera estándar de hacer esto es mediante un proceso Newton-Raphson. Kale [13], en 1962, estudió la convergencia de este método así como la de otro método de resolución al que llama el 'método de puntajes' [*]. Estos métodos son parecidos. El primero usa la matriz J^{-1} donde

$$J = \begin{pmatrix} J_{\lambda\lambda} & J_{\lambda\theta} \\ J_{\theta\lambda} & J_{\theta\theta} \end{pmatrix}$$

y J se encuentra particionada en las submatrices $J_{\lambda\lambda}$, $J_{\lambda\theta}$, $J_{\theta\lambda}$, $J_{\theta\theta}$, las cuales están definidas por

$$J_{\lambda\lambda} = \left\{ \frac{\partial^2 \log L}{\partial \lambda_s \partial \lambda_p} \right\}$$

$$J_{\lambda\theta} = \left\{ \frac{\partial^2 \log L}{\partial \lambda_s \partial \theta_{pj}} \right\}$$

* En inglés: 'Method of scoring'

$$J_{\theta\theta} = \left\{ \frac{\partial^2 \log L}{\partial \theta_{si} \partial \theta_{pj}} \right\}$$

y $J_{\theta\lambda} = J_{\lambda\theta}^t$. Las ecuaciones máximo verosímiles pueden ser resueltas iterativamente usando las siguientes ecuaciones:

$$(12.6) \quad \begin{pmatrix} T_{r+1} \\ S_{r+1} \end{pmatrix} = \begin{pmatrix} T_r \\ S_r \end{pmatrix} - J^{-1} \begin{pmatrix} \left[\frac{1}{\lambda_s} \sum_{k=1}^N P(s|x_k) - n \right]_{\lambda} T_r \\ \left[\sum_{k=1}^N P(s|x_k) \frac{\partial \log ds(x_k, \theta_k)}{\partial \theta_{si}} \right]_{\theta} S_r \end{pmatrix}$$

donde $\lambda = (\lambda_1, \dots, \lambda_g)$, $\theta = (\theta_1, \dots, \theta_q, \dots, \theta_{g1}, \dots, \theta_{gq})$, son los vectores de parámetros, T_r y S_r son las estimaciones r -ésimas para los λ_s y para los θ_s respectivamente. A su vez, el método de puntajes utiliza la inversa de la matriz de información:

$$I = \begin{pmatrix} I_{\lambda\lambda} & I_{\lambda\theta} \\ I_{\theta\lambda} & I_{\theta\theta} \end{pmatrix}$$

donde I se encuentra particionada en las submatrices $I_{\lambda\lambda}$, $I_{\lambda\theta}$, $I_{\theta\lambda}$, $I_{\theta\theta}$, que se encuentran definidas por

$$(12.7) \quad I_{\lambda\lambda} = \left\{ E \left(\frac{\partial \log L}{\partial \lambda_s} \frac{\partial \log L}{\partial \lambda_p} \right) \right\} = N \left\{ \frac{1}{\lambda_s \lambda_p} E [P(s|x) P(p|x)] \right\}$$

$$(12.8) \quad I_{\theta\theta} = \left\{ E \left(\frac{\partial \log L}{\partial \theta_{si}} \frac{\partial \log L}{\partial \theta_{pj}} \right) \right\} = N \left\{ \frac{1}{\lambda_s} E [P(s|x) P(p|x) \frac{\partial \log ds}{\partial \theta_{pj}}] \right\}$$

$$(2.9) \quad T_{00} = \left\{ E \left(\frac{\partial \log L}{\partial \theta_{si}} \frac{\partial \log L}{\partial \theta_{pj}} \right) \right\} = N \left\{ E \left[P(s/x) P(p/z) \frac{\partial \log ds}{\partial \theta_{si}} \frac{\partial \log dp}{\partial \theta_{pj}} \right] \right\}$$

La submatriz T_{0j} , es la transpuesta de T_{j0} . Si utilizamos la misma notación que en (2.6), las ecuaciones máximo verosimilables pueden ser resueltas iterativamente con el método de puntajes usando las siguientes ecuaciones:

$$(2.10) \quad \begin{pmatrix} T_{r+1} \\ S_{r+1} \end{pmatrix} = \begin{pmatrix} T_r \\ S_r \end{pmatrix} + \frac{1}{N} I^{-1} \begin{pmatrix} \left[\frac{1}{ds} \sum_{k=1}^N P(s/z_k) - N \right]_{s=T_r} \\ \left[\sum_{k=1}^N P(s/z_k) \frac{\partial \log ds(z_k, \theta_s)}{\partial \theta_{si}} \right]_{s=S_r} \end{pmatrix}$$

Tanto la matriz I como la matriz J se evalúan en (T_1, S_1) que es la estimación inicial de los parámetros. Para que exista convergencia en ambos métodos es necesario que (T_1, S_1) sea un estimador consistente de (θ, θ) . En la mayoría de los casos, las esperanzas en la matriz de información contienen integrales imposibles de evaluar en forma cerrada y difícilmente aproximables por series. El enfoque usual ha sido estimar la matriz de información con la muestra reemplazando el símbolo E en (2.7), (2.8) y (2.10) por $\frac{1}{N} \sum_{k=1}^N$.

Varios investigadores, al llevar a cabo Análisis de conglomerados (por ejemplo Friedman y Rubin en 1967 [2] y Ball en 1967 [4]), han reportado que los métodos iterativos pueden producir una multiplicidad de soluciones diferentes,

dependiendo de los estimadores iniciales que se utilicen al comenzar el proceso iterativo. Los métodos aquí presentados también tienen esta inconveniencia. Desde un punto de vista puramente matemático, no es sorprendente que un conjunto de ecuaciones no lineales tenga varias raíces diferentes. Sin embargo, puede sorprender el hecho de que surjan soluciones múltiples en el contexto de estimación estadística.

Veamos por medio de un ejemplo, lo que intuitivamente está sucediendo. Supongamos que se tiene una colección heterogénea de personas y los queremos dividir exactamente en dos grupos. Desde luego que existe un gran número de formas en que podemos hacer esto: hombres y mujeres, niños y adultos, liberales y conservadores, ricos y pobres, etc. La división apropiada depende de las variables usadas en el análisis. Si la única variable usada en el análisis es la altura, cualquier método razonable de agrupamiento asignará a la mayoría de los niños en un grupo y a la mayoría de los adultos en el otro. Si la longitud del pelo es usada como única variable de agrupamiento, una buena separación de hombres y mujeres se obtendría para personas adultas, pero quizás esta separación no sería muy buena para niños y jóvenes. Supongamos ahora que en el análisis se incluyen una multitud de variables. Un método que diera iguales pesos o ponderaciones a todas las variables, no sería un buen método. La ventaja de los métodos iterativos en la formación de los

grupos, el que el mismo procedimiento va determinando los pesos aplicados a las variables. En cada iteración, estos pesos se revisan de tal forma que se produzca una mejor discriminación entre grupos. Precisamente esta virtud del método puede originar soluciones múltiples.

Así pues, supongamos que en este ejemplo se hace una partición aleatoria del conjunto de individuos para determinar estimadores iniciales. Supongamos que debido al azar existe una pequeña mayoría de liberales en un grupo y una mayoría de conservadores en el otro. Entonces, las variables relacionadas con la actitud política discriminan mejor entre los grupos que cualquier otro conjunto de variables. Las actitudes políticas tendrán mayor peso en la siguiente iteración. Los liberales tendrán mayor probabilidad de pertenecer a un grupo que al otro. En la siguiente iteración, esta composición revisada de grupos generará pesos aún mayores para las variables de actitud política. De esta manera, el proceso convergerá a una tipología liberal-conservadora pura.

CAPITULO III

MEZCLAS DE DISTRIBUCIONES NORMALES

Este capítulo constituye la parte más importante de este trabajo. Aquí presentamos la estimación máxima verosímil de los parámetros involucrados en una mezcla, en el caso particular en que las distribuciones que forman la mezcla son normales multivariadas. Así pues, obtenemos los estimadores máximo verosímiles de las medias, varianzas y covarianzas en cada uno de los grupos, así como las proporciones de estos grupos en la mezcla.

El capítulo consta de cuatro secciones: en la primera se estudia el caso general en que no hay restricción alguna sobre las matrices de dispersión en los distintos grupos; en la segunda se supone que estas matrices son iguales para todos los grupos; en la tercera, que las matrices de dispersión son de la forma $\Sigma_i = \alpha_i \Sigma$, donde α_i es un escalar, y Σ es una matriz común; finalmente, en la cuarta sección se estudia el caso en que se tenga la misma matriz de correlación en los distintos grupos.

1. CASO GENERAL.

En este caso tenemos que

$$(3.1) \quad \alpha_s(\underline{x}, \underline{\mu}_s) = \alpha_s(\underline{x}, \underline{\mu}_s, \Sigma_s) \\ = (2\pi)^{-p/2} |\Sigma_s|^{-1/2} \exp \left\{ -\frac{1}{2} (\underline{x} - \underline{\mu}_s)' \Sigma_s^{-1} (\underline{x} - \underline{\mu}_s) \right\}$$

donde $\Sigma_s = \{ \sigma_{ij}^s \}$ es la matriz $p \times p$ de dispersión de las p variables x_1, x_2, \dots, x_p en la s -ésima distribución de la mezcla y $\Sigma_s^{-1} = \{ \tau_{ij}^s \} = \Sigma_s^{-1}$

Para obtener los estimadores máximo verosímiles de $\underline{\mu}_s = (\mu_{s1}, \dots, \mu_{sp})$ y de $\Sigma_s = \{ \sigma_{ij}^s \}$ utilizaremos la ecuación (2.5) y por lo tanto necesitamos calcular explícitamente $\frac{\partial \log \alpha_s}{\partial \mu_{si}}$ y $\frac{\partial \log \alpha_s}{\partial \tau_{ij}^s}$.

Para ello notemos que

$$\log \alpha_s(\underline{x}, \underline{\mu}_s, \Sigma_s) = -\frac{p}{2} \log 2\pi + \frac{1}{2} \log |\Sigma_s|^{-1} - \frac{1}{2} (\underline{x} - \underline{\mu}_s)' \Sigma_s^{-1} (\underline{x} - \underline{\mu}_s)$$

de donde

$$\frac{\partial \log \alpha_s(\underline{x}_k, \underline{\mu}_s, \Sigma_s)}{\partial \mu_{si}} = -\frac{1}{2} \frac{\partial}{\partial \mu_{si}} \left[(\underline{x}_k - \underline{\mu}_s)' \Sigma_s^{-1} (\underline{x}_k - \underline{\mu}_s) \right]$$

Es decir, tenemos que derivar la forma cuadrática

$$(\underline{x}_k - \underline{\mu}_s)' \Sigma_s^{-1} (\underline{x}_k - \underline{\mu}_s).$$

Pero,

$$\frac{\partial}{\partial \mu_{si}} (\underline{x}_k - \underline{\mu}_s)' \Sigma_s^{-1} (\underline{x}_k - \underline{\mu}_s) = -2 \Sigma_s^{-1} (\underline{x}_k - \underline{\mu}_s),$$

de donde

$$\frac{\partial}{\partial \mu_{si}} (x_k - \mu_s)' \Sigma^s (x_k - \mu_s) = -2 \sum_{j=1}^p \bar{v}_{ij}^s (x_{jk} - \mu_{sj})$$

Entonces, obtenemos finalmente que

$$(3.2) \quad \frac{\partial \log \alpha_s(x_k, \mu_s, \Sigma_s)}{\partial \mu_{si}} = \sum_{j=1}^p \bar{v}_{ij}^s (x_{jk} - \mu_{sj})$$

Ahora obtendremos
$$\frac{\partial \log \alpha_s(x_k, \mu_s, \Sigma_s)}{\partial \bar{v}_{ij}^s}$$

Si $i=j$, entonces

$$\frac{\partial \log \alpha_s(x_k, \mu_s, \Sigma_s)}{\partial \bar{v}_{ii}^s} = \frac{1}{2} \frac{\partial}{\partial \bar{v}_{ii}^s} \log |\Sigma^s| - \frac{1}{2} \frac{\partial}{\partial \bar{v}_{ii}^s} (x_k - \mu_s)' \Sigma^s (x_k - \mu_s)$$

Ahora bien, como

$$|\Sigma^s| = \bar{v}_{ii}^s [\text{cofactor de } \bar{v}_{ii}^s] + [\text{términos sin } \bar{v}_{ii}^s]$$

y

$$(x_k - \mu_s)' \Sigma^s (x_k - \mu_s) = (x_{ki} - \mu_{si}) \bar{v}_{ii}^s + [\text{términos sin } \bar{v}_{ii}^s]$$

entonces

$$\frac{\partial \log \alpha_s(x_k, \mu_s, \Sigma_s)}{\partial \bar{v}_{ii}^s} = \frac{1}{2} \left\{ \frac{[\text{cofactor de } \bar{v}_{ii}^s]}{|\Sigma^s|} - (x_{ki} - \mu_{si})^2 \right\}$$

Pero,

$$\frac{[\text{cofactor de } \bar{v}_{ii}^s]}{|\Sigma^s|} = \bar{v}_s^{ii}$$

Por lo tanto

$$\frac{\partial \log \alpha_s(\underline{x}_k, \mu_s, \Sigma_s)}{\partial \bar{v}_{ij}^s} = \frac{1}{2} \left\{ \bar{v}_s^{-ij} - (x_{ik} - \mu_{si})^2 \right\}$$

Si $i \neq j$, entonces

$$\frac{\partial \log \alpha_s(\underline{x}_k, \mu_s, \Sigma_s)}{\partial \bar{v}_{ij}^s} = \frac{1}{2} \frac{\partial}{\partial \bar{v}_{ij}^s} \log |\Sigma^s| - \frac{1}{2} \frac{\partial}{\partial \bar{v}_{ij}^s} (\underline{x}_k - \mu_s)' \Sigma^s (\underline{x}_k - \mu_s)$$

Dado que

$$|\Sigma^s| = 2 \bar{v}_{ij}^s [\text{cofactor de } \bar{v}_{ij}^s] + [\text{términos sin } \bar{v}_{ij}^s],$$

y

$$(\underline{x}_k - \mu_s)' \Sigma^s (\underline{x}_k - \mu_s) = 2(x_{ik} - \mu_{si}) \bar{v}_{ij}^s (x_{jk} - \mu_{sj}) + [\text{términos sin } \bar{v}_{ij}^s],$$

tenemos que

$$\frac{\partial \log \alpha_s(\underline{x}_k, \mu_s, \Sigma_s)}{\partial \bar{v}_{ij}^s} = \frac{[\text{cofactor de } \bar{v}_{ij}^s]}{|\Sigma^s|} - (x_{ik} - \mu_{si})(x_{jk} - \mu_{sj})$$

y nuevamente, usando que

$$\frac{[\text{cofactor de } \bar{v}_{ij}^s]}{|\Sigma^s|} = \bar{v}_s^{-ij},$$

obtenemos

$$\frac{\partial \log \alpha_s(\underline{x}_k, \mu_s, \Sigma_s)}{\partial \bar{v}_{ij}^s} = \bar{v}_s^{-ij} - (x_{ik} - \mu_{si})(x_{jk} - \mu_{sj})$$

Podemos juntar ambos casos ($i \neq j$, $i = j$) para escribir una sola expresión:

$$(3.3) \quad \frac{\partial \log \alpha_s(\underline{x}_k, \mu_s, \Sigma_s)}{\partial \bar{v}_{ij}^s} = \left(1 - \frac{\delta_{ij}}{2}\right) \left\{ \bar{v}_s^{-ij} - (x_{ik} - \mu_{si})(x_{jk} - \mu_{sj}) \right\}$$

en donde $\delta_{ij} = \begin{cases} 1 & \text{si } i=j \\ 0 & \text{si } i \neq j \end{cases}$

Podemos usar toda esta información para obtener el siguiente resultado.

Teorema 1. Los estimadores máximo verosímiles de los parámetros de una mezcla de distribuciones normales multivariadas con distintas matrices de dispersión están dados por

$$(3.4) \quad \hat{\lambda}_s = \frac{1}{N} \sum_{k=1}^N \hat{P}(s | \underline{x}_k)$$

$$(3.5) \quad \hat{M}_{si} = \frac{1}{N \hat{\lambda}_s} \sum_{k=1}^N \hat{P}(s | \underline{x}_k) x_{ik}$$

$$(3.6) \quad \hat{V}_{ij} = \frac{1}{N \hat{\lambda}_s} \sum_{k=1}^N \hat{P}(s | \underline{x}_k) (x_{ik} - \hat{\mu}_{si})(x_{jk} - \hat{\mu}_{sj})$$

Demostración. Si sustituimos (3.2) en (2.5) obtenemos que

$$\sum_{k=1}^N \hat{P}(s | \underline{x}_k) \sum_{j=1}^p \hat{V}_{ij} (x_{jk} - \hat{\mu}_{sj}) = 0$$

Entonces

$$\begin{aligned} & \sum_{j=1}^p \hat{V}_{ij} \left[\sum_{k=1}^N \hat{P}(s | \underline{x}_k) (x_{jk} - \hat{\mu}_{sj}) \right] = \\ & = \sum_{j=1}^p \hat{V}_{ij} \left[\sum_{k=1}^N \hat{P}(s | \underline{x}_k) x_{jk} - \hat{\mu}_{sj} \sum_{k=1}^N \hat{P}(s | \underline{x}_k) \right] = \\ & = \sum_{j=1}^p \hat{V}_{ij} \left[\sum_{k=1}^N \hat{P}(s | \underline{x}_k) x_{jk} - \hat{\mu}_{sj} N \hat{\lambda}_s \right] = 0. \end{aligned}$$

Por lo tanto

$$\sum_{k=1}^N \hat{p}(s/x_k) x_{jk} - \hat{\mu}_{sj} \cdot N \hat{\lambda}_s = 0,$$

de donde

$$\hat{\mu}_{sj} = \frac{1}{N \hat{\lambda}_s} \sum_{k=1}^N \hat{p}(s/x_k) x_{jk},$$

que es la ecuación (3.5).

Para demostrar (3.6) sustituiremos (3.3) en (2.5), con lo que obtenemos que

$$\sum_{k=1}^N \hat{p}(s/x_k) \left(1 - \frac{\delta_{ij}}{2}\right) \left[\hat{v}_s^{ij} - (x_{ik} - \hat{\mu}_{si})(x_{jk} - \hat{\mu}_{sj}) \right] = 0.$$

Entonces

$$\left(1 - \frac{\delta_{ij}}{2}\right) \sum_{k=1}^N \hat{p}(s/x_k) \left[\hat{v}_s^{ij} - (x_{ik} - \hat{\mu}_{si})(x_{jk} - \hat{\mu}_{sj}) \right] = 0.$$

De donde

$$\hat{v}_s^{ij} \sum_{k=1}^N \hat{p}(s/x_k) = \sum_{k=1}^N \hat{p}(s/x_k) (x_{ik} - \hat{\mu}_{si})(x_{jk} - \hat{\mu}_{sj}),$$

y despejando \hat{v}_s^{ij} obtenemos finalmente que

$$\hat{v}_s^{ij} = \frac{1}{N \hat{\lambda}_s} \sum_{k=1}^N \hat{p}(s/x_k) (x_{ik} - \hat{\mu}_{si})(x_{jk} - \hat{\mu}_{sj}).$$

La ecuación (3.4) es simplemente (2.4). \square

Tal y como lo indicamos en el capítulo anterior, las ecuaciones (3.4) a (3.6) no están en forma cerrada, pues no se conocen las probabilidades de pertenencia a los distintos grupos para cada uno de los objetos o individuos en la muestra. En

el capítulo anterior presentamos algunos métodos iterativos de resolución de estas ecuaciones e indicamos que las soluciones dependían del agrupamiento inicial. Usualmente, este agrupamiento inicial se hace con técnicas jerárquicas de análisis de conglomerados o bien minimizando la traza o el determinante de la matriz W , que es la matriz de dispersión dentro de grupos. Esto se verá con mayor detalle en el capítulo V. Cuando ya se tienen los valores numéricos de los parámetros de la mezcla (después del proceso iterativo) se asigna a cada objeto o individuo al grupo al que tenga mayor probabilidad de pertenencia.

Finalmente, notemos que en el caso general de mezclas normales (distintas matrices de dispersión) se tienen que estimar $(g-1) + gp + \frac{1}{2}gp(p+1)$ parámetros (el primer término corresponde a las λ'_i , el segundo a las μ'_i y el tercero a las γ'_i). Si g y p son relativamente grandes, se necesita un tamaño de muestra N muy grande para poder así estimar tantos parámetros -en particular, el número de individuos dentro de cada grupo debe ser suficientemente grande para obtener estimaciones satisfactorias de las matrices de dispersión. Debido a esto, resulta conveniente hacer algunas simplificaciones a las distribuciones componentes de la mezcla, principalmente a las matrices de dispersión de las distribuciones. Estas simplificaciones dependen de la naturaleza del problema que se tenga. Así pues, si se supone que las matrices de dis-

persión son iguales en los distintos grupos, se tienen que estimar $(g-1) + gp + \frac{1}{2} p(p+1)$ parámetros, lográndose una reducción de $\frac{1}{2}(g-1)p(p+1)$ parámetros. Otra posibilidad es suponer matrices de correlación iguales. En este caso se tienen que estimar $(g-1) + gp + \frac{1}{2} p(p-1) + gp$ parámetros, es decir, $\frac{1}{2} p(p-1)(g-1)$ parámetros menos que estimar. Aún otra posibilidad es suponer que las matrices de dispersión en los distintos grupos son iguales excepto por una constante multiplicativa, i.e., $\Sigma_s = a_s \Sigma$ $s = 1, 2, \dots, g$. En este caso se tienen que estimar $(g-1) + gp + \frac{1}{2} p(p+1) + g$ parámetros (el cuarto término corresponde a las a_s). Estas simplificaciones, aparte de reducir el número de parámetros por estimar, hacen que el problema de agrupación sea computacionalmente más simple. Desde luego, las suposiciones deben ser realistas, es decir, deben corresponder a la naturaleza del problema; de lo contrario las estimaciones no serán buenas.

A continuación presentamos las estimaciones máximo verosímiles para cada una de estas simplificaciones.

2. MATRICES DE DISPERSIÓN IGUALES.

En este caso suponemos que $\Sigma_s = \Sigma = \{ \sigma^{-2} \delta \}$.

La ecuación (2.5) no es válida, pues fue derivada bajo la suposición de que los parámetros de distintos grupos no estaban relacionados funcionalmente. La ecuación de verosimilitud correcta es

$$\frac{\partial \mathcal{L}}{\partial \hat{v}_{ij}} = \sum_{s=1}^g \frac{\partial \mathcal{L}}{\partial \hat{v}_{ij}^s} \frac{\partial \hat{v}_{ij}^s}{\partial \hat{v}_{ij}}$$

De hecho tenemos el siguiente resultado.

Teorema 2. Los estimadores máximo verosímiles de los parámetros de una mezcla de distribuciones normales multivariadas con matrices de dispersión iguales están dadas por

$$(3.7) \quad \hat{\lambda}_s = \frac{1}{N} \sum_{k=1}^N \hat{p}(s | \underline{x}_k)$$

$$(3.8) \quad \hat{\mu}_{si} = \frac{1}{N \hat{\lambda}_s} \sum_{k=1}^N \hat{p}(s | \underline{x}_k) x_{ik}$$

$$(3.9) \quad \hat{v}_{ij}^{-1} = \frac{1}{N} \sum_{k=1}^N x_{ik} x_{jk} - \sum_{s=1}^g \hat{\lambda}_s \hat{\mu}_{si} \hat{\mu}_{sj}$$

Demostración. La ecuación (3.7) es simplemente (2.4). La ecuación (3.8) es clara si se observa la demostración de la ecuación (3.5) del Teorema 1 de este capítulo. Para demostrar (3.9) primero obtendremos $\frac{\partial \mathcal{L}}{\partial \hat{v}_{ij}^s}$. Para esto, ya hicimos notar que

$$\frac{\partial \mathcal{L}}{\partial \hat{v}_{ij}} = \sum_{s=1}^g \frac{\partial \mathcal{L}}{\partial \hat{v}_{ij}^s} \frac{\partial \hat{v}_{ij}^s}{\partial \hat{v}_{ij}}$$

y dado que

$$\frac{\partial \hat{v}_{ij}^s}{\partial \hat{v}_{ij}} = 1$$

tenemos que

$$= (1 - \frac{\delta_{ij}}{2}) \left[\sum_{s=1}^g \sum_{k=1}^N \frac{\lambda_s ds}{f(x_k)} \bar{v}_s^{ij} - \sum_{s=1}^g \sum_{k=1}^N \frac{\lambda_s ds}{f(x_k)} (x_{ik} - \mu_{si})(x_{jk} - \mu_{sj}) \right]$$

$$= (1 - \frac{\delta_{ij}}{2}) \left[N \bar{v}^{ij} - \sum_{s=1}^g \sum_{k=1}^N \frac{\lambda_s ds(x_k, \mu_s, \Sigma_s)}{f(x_k)} (x_{ik} - \mu_{si})(x_{jk} - \mu_{sj}) \right]$$

Al igualar esto a cero obtenemos que

$$N \bar{v}^{ij} - \sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s ds(x_k, \mu_s, \Sigma_s)}{f(x_k)} (x_{ik} - \hat{\mu}_{si})(x_{jk} - \hat{\mu}_{sj}) = 0$$

Es decir,

$$N \bar{v}^{ij} = \sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s ds(x_k, \mu_s, \Sigma_s)}{f(x_k)} [x_{ik} x_{jk} - x_{ik} \hat{\mu}_{sj} - x_{jk} \hat{\mu}_{si} + \hat{\mu}_{si} \hat{\mu}_{sj}]$$

$$= \sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s ds(x_k, \mu_s, \Sigma_s)}{f(x_k)} x_{ik} x_{jk} + \sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s ds(x_k, \mu_s, \Sigma_s)}{f(x_k)} \hat{\mu}_{si} \hat{\mu}_{sj}$$

$$- \sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s ds(x_k, \mu_s, \Sigma_s)}{f(x_k)} (x_{ik} \hat{\mu}_{sj} + x_{jk} \hat{\mu}_{si})$$

Ahora bien,

$$\sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s ds(x_k, \mu_s, \Sigma_s)}{f(x_k)} \hat{\mu}_{si} \hat{\mu}_{sj} = \sum_{s=1}^g \hat{\lambda}_s \hat{\mu}_{si} \hat{\mu}_{sj} \sum_{k=1}^N \frac{ds(x_k, \mu_s, \Sigma_s)}{f(x_k)}$$

$$= \sum_{s=1}^g \hat{\lambda}_s \hat{\mu}_{si} \hat{\mu}_{sj} \left[\frac{1}{\lambda_s} \sum_{k=1}^N \hat{P}(s|x_k) \right] = N \sum_{s=1}^g \hat{\lambda}_s \hat{\mu}_{si} \hat{\mu}_{sj}$$

$$\frac{\partial L}{\partial \tau_{ij}^s} = \sum_{s=1}^g \frac{\partial}{\partial \tau_{ij}^s} \left[\sum_{k=1}^N \log f(x_k) \right] = \sum_{s=1}^g \sum_{k=1}^N \frac{\frac{\partial}{\partial \tau_{ij}^s} f(x_k)}{f(x_k)}$$

Así pues, tenemos que encontrar $\frac{\partial}{\partial \tau_{ij}^s} f(x_k)$:

$$\begin{aligned} \frac{\partial}{\partial \tau_{ij}^s} f(x_k) &= \frac{\partial}{\partial \tau_{ij}^s} \sum_{s=1}^g \lambda_s \kappa_s(x_k, \mu_s, \Sigma_s) = \\ &= \frac{\partial}{\partial \tau_{ij}^s} \lambda_s \kappa_s(x_k, \mu_s, \Sigma_s) = \\ &= \frac{\partial}{\partial \tau_{ij}^s} \left[e^{\log \lambda_s \kappa_s(x_k, \mu_s, \Sigma_s)} \right] = \\ &= \lambda_s \kappa_s(x_k, \mu_s, \Sigma_s) \frac{\partial}{\partial \tau_{ij}^s} \log [\lambda_s \kappa_s(x_k, \mu_s, \Sigma_s)] = \\ &= \lambda_s \kappa_s(x_k, \mu_s, \Sigma_s) \frac{\partial}{\partial \tau_{ij}^s} \log \kappa_s(x_k, \mu_s, \Sigma_s). \end{aligned}$$

y usando la expresión (3.3) obtenemos que

$$\frac{\partial}{\partial \tau_{ij}^s} f(x_k) = \lambda_s \kappa_s(x_k, \mu_s, \Sigma_s) \left(1 - \frac{\delta_{ij}}{2}\right) \left[\sqrt{\Sigma_s}^{-1} (x_{ik} - \mu_{si})(x_{jk} - \mu_{sj}) \right].$$

Por lo tanto

$$\frac{\partial L}{\partial \tau_{ij}^s} = \sum_{s=1}^g \sum_{k=1}^N \frac{\left(1 - \frac{\delta_{ij}}{2}\right) \lambda_s \kappa_s(x_k, \mu_s, \Sigma_s)}{f(x_k)} \left[\sqrt{\Sigma_s}^{-1} (x_{ik} - \mu_{si})(x_{jk} - \mu_{sj}) \right] =$$

y también,

$$\begin{aligned}
 & \sum_{s=1}^g \sum_{k=1}^N \frac{\lambda_s ds(x_k, \mu_s, \Sigma_s)}{f(x_k)} (x_{ik} \hat{\mu}_{sj} + x_{jk} \hat{\mu}_{si}) \\
 &= \sum_{s=1}^g \sum_{k=1}^N \left[\hat{P}(s|x_k) x_{ik} \hat{\mu}_{sj} + \hat{P}(s|x_k) x_{jk} \hat{\mu}_{si} \right] \\
 &= \sum_{s=1}^g \left[\hat{\mu}_{sj} \sum_{k=1}^N \hat{P}(s|x_k) x_{ik} + \hat{\mu}_{si} \sum_{k=1}^N \hat{P}(s|x_k) x_{jk} \right] \\
 &= \sum_{s=1}^g \left[\hat{\mu}_{sj} \hat{\mu}_{si} N \hat{\lambda}_s + \hat{\mu}_{si} \hat{\mu}_{sj} N \hat{\lambda}_s \right] \\
 &= 2N \sum_{s=1}^g \hat{\lambda}_s \hat{\mu}_{si} \hat{\mu}_{sj}.
 \end{aligned}$$

Por lo tanto

$$\begin{aligned}
 N \hat{V}^{ij} &= \sum_{s=1}^g \sum_{k=1}^N \frac{\lambda_s ds(x_k, \mu_s, \Sigma_s)}{f(x_k)} x_{ik} x_{jk} - N \sum_{s=1}^g \hat{\lambda}_s \hat{\mu}_{si} \hat{\mu}_{sj} \\
 &= \sum_{k=1}^N x_{ik} x_{jk} - N \sum_{s=1}^g \hat{\lambda}_s \hat{\mu}_{si} \hat{\mu}_{sj}.
 \end{aligned}$$

Es decir

$$\hat{V}^{ij} = \frac{1}{N} \sum_{k=1}^N x_{ik} x_{jk} - \sum_{s=1}^g \hat{\lambda}_s \hat{\mu}_{si} \hat{\mu}_{sj},$$

que es la ecuación (3.9). \square

5. CASO $\Sigma_s = a_s \Sigma$

En ciertas ocasiones, el suponer que las matrices de dispersión son iguales para todos los grupos, puede no ser una suposición muy adecuada. Una suposición un poco más general y con muchos parámetros menos por estimar que en el caso general consiste en pensar que las matrices de dispersión son iguales excepto por una constante multiplicativa, i.e., $\Sigma_s = a_s \Sigma$ donde a_s es un número real positivo y Σ es una matriz positiva definida de tamaño $p \times p$, $s = 1, \dots, g$

En este caso tenemos que

$$h_s(\bar{x}, \mu_s, \Sigma_s) = h_s(\bar{x}, \mu_s, a_s \Sigma) = (2\pi)^{-p/2} |\Sigma_s|^{-1/2} \exp\left\{-\frac{1}{2}(\bar{x} - \mu_s)' \Sigma_s^{-1} (\bar{x} - \mu_s)\right\}$$

donde

$$\Sigma_s = \{\bar{v}_s^{ij}\} = a_s \Sigma = \{a_s v^{ij}\}$$

y

$$\Sigma_s^{-1} = \{\bar{v}_s^{ij}\}^{-1} = \{z_s^{ij}\} = a_s^{-1} \Sigma^{-1} = \{a_s^{-1} v^{ij}\}$$

con

$$a_s^{-1} = (a_s)^{-1}, \quad \Sigma = \{v^{ij}\}, \quad \Sigma^{-1} = \{v^{ij}\}$$

Tenemos que encontrar los estimadores máximo verosímiles de a_s , μ_s , v^{ij} y Σ , $s = 1, 2, \dots, g$.

El siguiente teorema muestra estos resultados.

Teorema 3. Los estimadores máximo verosímiles de los parámetros de una mezcla de distribuciones normales multivariadas con matrices de dispersión de la forma $\Sigma_s = a_s \Sigma$, $s = 1, 2, \dots, g$, están dados por

$$(3.10) \quad \hat{\lambda}_s = \frac{1}{N} \sum_{k=1}^N \hat{p}(s | \underline{x}_k)$$

$$(3.11) \quad \hat{\mu}_{si} = \frac{1}{N \hat{\lambda}_s} \sum_{k=1}^N \hat{p}(s | \underline{x}_k) x_{ik}$$

$$(3.12) \quad \hat{V}^{ij} = \frac{1}{N} \sum_{s=1}^g \sum_{k=1}^N \frac{\hat{p}(s | \underline{x}_k) x_{ik} x_{jk}}{a_s} - \sum_{s=1}^g \frac{\hat{\lambda}_s}{\hat{a}_s} \hat{\mu}_{si} \hat{\mu}_{sj}$$

$$(3.13) \quad \hat{Q}_s = \frac{1}{PN \hat{\lambda}_s} \sum_{k=1}^N \hat{p}(s | \underline{x}_k) (\underline{x}_k - \hat{\mu}_s)' \hat{\Sigma}^{-1} (\underline{x}_k - \hat{\mu}_s).$$

Demostración. Es obvio que las ecuaciones (3.10) y (3.11) siguen siendo válidas en este caso. Demostraremos ahora la ecuación (3.12). La ecuación de verosimilitud en este caso es

$$\frac{\partial L}{\partial \hat{V}_{ij}^s} = \sum_{s=1}^g \frac{\partial L}{\partial \hat{V}_{ij}^s} \frac{\partial \hat{V}_{ij}^s}{\partial \hat{V}_{ij}^s} = 0$$

Ya tenemos en el caso general que

$$\frac{\partial L}{\partial \hat{V}_{ij}^s} = \sum_{k=1}^N \frac{\lambda_s a_s (\underline{x}_k, \mu_s, \Sigma_s)}{f(\underline{x}_k)} (1 - \delta_{ij}^s) \left[\sqrt{V_s}^{-1} \delta_{ij}^s - (x_{ik} - \mu_{si})(x_{jk} - \mu_{sj}) \right],$$

Y además en este caso

$$\frac{\partial \hat{V}_{ij}^s}{\partial \hat{V}_{ij}^s} = \frac{\partial}{\partial \hat{V}_{ij}^s} (a^s \hat{V}_{ij}^s) = a^s.$$

Por lo tanto

$$\frac{\partial L}{\partial \sigma_{ij}} = \sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s ds}{f(z_k)} (1 - \delta_{ij}) \left[\hat{v}_s^{ij} - (x_{ik} - \hat{\mu}_{si})(x_{jk} - \hat{\mu}_{sj}) \right] \hat{a}^s = 0$$

y como $\nabla_s^i \hat{a}^s = a_s \nabla^s \hat{a}^s$, tenemos

$$(1 - \delta_{ij}) \sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s ds}{f(z_k)} \left[\hat{a}_s \hat{v}_s^{ij} - (x_{ik} - \hat{\mu}_{si})(x_{jk} - \hat{\mu}_{sj}) \right] \hat{a}^s = 0$$

i.e.,

$$\sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s ds}{f(z_k)} \left[\hat{v}_s^{ij} - \hat{a}_s (x_{ik} - \hat{\mu}_{si})(x_{jk} - \hat{\mu}_{sj}) \right] = 0.$$

Entonces

$$\begin{aligned} \hat{v}_s^{ij} \sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s ds}{f(z_k)} &= \sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s ds}{f(z_k)} \hat{a}^s (x_{ik} - \hat{\mu}_{si})(x_{jk} - \hat{\mu}_{sj}) \\ &= \sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s ds}{f(z_k)} \hat{a}^s x_{ik} x_{jk} - \sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s ds}{f(z_k)} \hat{a}^s (x_{ik} \hat{\mu}_{sj} + x_{jk} \hat{\mu}_{si}) \\ &\quad + \sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s ds}{f(z_k)} \hat{a}^s \hat{\mu}_{si} \hat{\mu}_{sj} \\ &= \sum_{s=1}^g \sum_{k=1}^N \hat{a}^s \hat{p}(s|z_k) x_{ik} x_{jk} - \sum_{s=1}^g \hat{a}^s \left[\sum_{k=1}^N \frac{\hat{\lambda}_s ds}{f(z_k)} (x_{ik} \hat{\mu}_{sj} + x_{jk} \hat{\mu}_{si}) \right] \\ &\quad + \sum_{s=1}^g \hat{a}^s \hat{\mu}_{si} \hat{\mu}_{sj} \sum_{k=1}^N \frac{\hat{\lambda}_s ds (z_k, \mu_s, \Sigma_s)}{f(z_k)} \\ &= \sum_{s=1}^g \sum_{k=1}^N \hat{a}^s \hat{p}(s|z_k) (x_{ik} x_{jk}) - \sum_{s=1}^g \hat{a}^s \left[\hat{\mu}_{sj} \sum_{k=1}^N \hat{p}(s|z_k) x_{ik} + \hat{\mu}_{si} \sum_{k=1}^N \hat{p}(s|z_k) x_{jk} \right] \\ &\quad + \sum_{s=1}^g \hat{a}^s \hat{\mu}_{si} \hat{\mu}_{sj} N \hat{\lambda}_s = \end{aligned}$$

$$\begin{aligned}
&= \sum_{s=1}^g \sum_{k=1}^N \hat{a}^s \hat{p}(s|x_k) x_{ik} x_{jk} - \sum_{s=1}^g \hat{a}^s [\hat{\mu}_{sj} N \hat{\lambda}_s \hat{\mu}_{si} + \hat{\mu}_{si} N \hat{\lambda}_s \hat{\mu}_{sj}] \\
&\quad + N \sum_{s=1}^g \hat{\lambda}_s \hat{a}^s \hat{\mu}_{si} \hat{\mu}_{sj} \\
&= \sum_{s=1}^g \sum_{k=1}^N \hat{a}^s \hat{p}(s|x_k) x_{ik} x_{jk} - 2N \sum_{s=1}^g \hat{a}^s \hat{\lambda}_s \hat{\mu}_{si} \hat{\mu}_{sj} + N \sum_{s=1}^g \hat{a}^s \hat{\lambda}_s \hat{\mu}_{si} \hat{\mu}_{sj} \\
&= \sum_{s=1}^g \sum_{k=1}^N \hat{a}^s \hat{p}(s|x_k) x_{ik} x_{jk} - N \sum_{s=1}^g \hat{a}^s \hat{\lambda}_s \hat{\mu}_{si} \hat{\mu}_{sj}
\end{aligned}$$

De donde obtenemos finalmente que

$$\hat{\gamma}^i \cdot \hat{\sigma} \cdot N = \sum_{s=1}^g \sum_{k=1}^N \hat{a}^s \hat{p}(s|x_k) x_{ik} x_{jk} - N \sum_{s=1}^g \hat{a}^s \hat{\lambda}_s \hat{\mu}_{si} \hat{\mu}_{sj}$$

Al despejar de esta ecuación obtenemos (3.12).

Finalmente, demostramos la ecuación (3.13). Para

ello:

$$\begin{aligned}
\frac{\partial L}{\partial a^s} &= \frac{\partial}{\partial a^s} \left[\sum_{k=1}^N \log \sum_{s=1}^g \lambda_s a^s (x_k, \mu_s, \Sigma_s) \right] \\
&= \sum_{k=1}^N \frac{\partial}{\partial a^s} \log \sum_{s=1}^g \lambda_s a^s (x_k, \mu_s, \Sigma_s) \\
&= \sum_{k=1}^N \lambda_s \frac{\partial}{\partial a^s} \frac{N s (x_k, \mu_s, \Sigma_s)}{f(x_k)}
\end{aligned}$$

Ahora bien,

$$\begin{aligned}
\frac{\partial}{\partial a^s} \kappa_s &= \frac{\partial}{\partial a^s} \exp(\log \kappa_s) = \kappa_s \frac{\partial}{\partial a^s} \log \kappa_s = \\
&= \kappa_s \frac{\partial}{\partial a^s} \left[\frac{1}{2} \log |a^s \Sigma^{-1}| - \frac{1}{2} (\mathbf{x}_k - \mu_s)' a^s \Sigma^{-1} (\mathbf{x}_k - \mu_s) \right] \\
&= \kappa_s \left[\frac{1}{2} \frac{\partial}{\partial a^s} \log \left((a^s)^p |\Sigma^{-1}| \right) - \frac{1}{2} \frac{\partial}{\partial a^s} \left(a^s (\mathbf{x}_k - \mu_s)' \Sigma^{-1} (\mathbf{x}_k - \mu_s) \right) \right] \\
&= \kappa_s \left[\frac{1}{2} \frac{p (a^s)^{p-1} |\Sigma^{-1}|}{(a^s)^p |\Sigma^{-1}|} - \frac{1}{2} (\mathbf{x}_k - \mu_s)' \Sigma^{-1} (\mathbf{x}_k - \mu_s) \right] \\
&= \frac{1}{2} \kappa_s (\mathbf{x}_k, \mu_s, \Sigma_s) \left[\frac{p}{a^s} - (\mathbf{x}_k - \mu_s)' \Sigma^{-1} (\mathbf{x}_k - \mu_s) \right]
\end{aligned}$$

Por lo tanto

$$\frac{\partial \mathcal{L}}{\partial a^s} = \frac{1}{2} \sum_{k=1}^N \frac{\hat{\lambda}_s ds}{f(\mathbf{x}_k)} \left[p \hat{a}_s - (\mathbf{x}_k - \hat{\mu}_s)' \hat{\Sigma}^{-1} (\mathbf{x}_k - \hat{\mu}_s) \right] = 0$$

De donde

$$\sum_{k=1}^N \hat{p}(s|\mathbf{x}_k) p \hat{a}_s = \sum_{k=1}^N \hat{p}(s|\mathbf{x}_k) (\mathbf{x}_k - \hat{\mu}_s)' \hat{\Sigma}^{-1} (\mathbf{x}_k - \hat{\mu}_s)$$

Despejando \hat{a}_s obtenemos finalmente

$$\hat{a}_s = \frac{1}{p N \hat{\lambda}_s} \sum_{k=1}^N \hat{p}(s|\mathbf{x}_k) (\mathbf{x}_k - \hat{\mu}_s)' \hat{\Sigma}^{-1} (\mathbf{x}_k - \hat{\mu}_s)$$

que es la ecuación (3.13). \square

Notemos que si $Q_s = I$, $s=1, 2, \dots, g$, la ecuación (3.12) es la ecuación (3.9), que es el caso en que se tienen matrices de dispersión iguales.

Observemos también que la estimación máximo verosímil de los Q'_s está en términos de la distancia de cada objeto a la media del grupo en cuestión con la inversa de la matriz estimada de dispersión como métrica, esto es, la distancia de Mahalanobis, y de las probabilidades de pertenencia al grupo en cuestión para cada uno de los objetos.

Finalmente notemos que si los grupos estuvieran muy separados tendríamos que

$$\hat{Q}'_s = \frac{1}{PN \hat{\lambda}_s} \sum_{j \in S} (x_j - \hat{\mu}_s) \hat{\Sigma}^{-1} (x_j - \hat{\mu}_s)$$

donde

$$S = \{i: x_i \text{ está en el grupo } s\}$$

4. MATRICES DE CORRELACIÓN IGUALES.

Existen situaciones en donde la suposición de iguales matrices de dispersión en los grupos no es válida, porque las variables en cada grupo tienen distintas varianzas, pero podemos suponer que tienen la misma matriz de correlación; por ejemplo, en Taxonomía nos podemos encontrar con una situación donde al clasificar diferentes especies animales de acuerdo a un conjunto de medidas sobre varios órganos del

cuerpo, estas medidas estén muy correlacionadas en cada grupo pero difieren en medias y varianzas para distintos grupos; o en Medicina, tratando de agrupar pacientes de acuerdo a un conjunto de variables como temperatura, presión, etc., las variables en distintos grupos pueden tener las mismas correlaciones pero diferentes medias y varianzas. Tal y como ya indicamos, en estas situaciones tenemos $(g-1)p(p-1)/2$ parámetros menos por estimar que en el caso general, lo cual es una gran ventaja.

Desde luego que los estimadores máximo verosímiles de los λ'_i y de los μ'_s siguen siendo los mismos que en el caso general. Los estimadores de las varianzas se obtienen de la ecuación (3.6) y están dados por

$$\hat{V}_i'' = \frac{1}{N \hat{\lambda}_i} \sum_{k=1}^N \hat{P}(s | X_k) (X_{ik} - \hat{\mu}_{si})^2.$$

Ahora bien, como estamos suponiendo que los distintos grupos tienen una matriz de correlación común $C = \{c^{ij}\}$ entonces para la población s , tenemos que la matriz de dispersión se puede escribir de la manera siguiente:

$$\Sigma_s = D_s C D_s$$

donde $C = \{c^{ij}\}$ y D_s es la matriz diagonal dada por

$$D_s = \begin{bmatrix} d_s^1 & & & \\ & d_s^2 & & \\ & & \dots & \\ 0 & & & d_s^p \end{bmatrix} = \begin{bmatrix} \sqrt{V_s^{11}} & & & 0 \\ & \sqrt{V_s^{22}} & & \\ & & \dots & \\ 0 & & & \sqrt{V_s^{pp}} \end{bmatrix}$$

Si $\{C_{ij}\} = \{C_{ij}^s\}^{-1}$ y $\{d_i^s\} = \{d_i^s\}^{-1}$, podemos expresar a $f(\underline{x})$ de la manera siguiente

$$f(\underline{x}) = \sum_{s=1}^g \lambda_s (2\pi)^{-p/2} |D_s^{-1} C^{-1} D_s^{-1}|^{\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\underline{x} - \mu_s)' D_s^{-1} C D_s^{-1} (\underline{x} - \mu_s) \right\}.$$

En forma análoga a las dos secciones anteriores las ecuaciones máximo verosímiles son

$$\frac{\partial L}{\partial C_{ij}} = \sum_{s=1}^g \frac{\partial L}{\partial \hat{v}_{ij}^s} \frac{\partial \hat{v}_{ij}^s}{\partial C_{ij}} = 0$$

y dado que

$$C_{ij} = \frac{\hat{v}_{ij}^s}{\sqrt{\hat{v}_{ii}^s \hat{v}_{jj}^s}} = \frac{\hat{v}_{ij}^s}{d_i^s d_j^s}$$

entonces

$$\frac{\partial \hat{v}_{ij}^s}{\partial C_{ij}} = d_i^s d_j^s$$

Luego, las ecuaciones de máxima verosimilitud son

$$\frac{\partial L}{\partial C_{ij}} = \sum_{s=1}^g \frac{\partial L}{\partial \hat{v}_{ij}^s} \hat{d}_i^s \hat{d}_j^s = 0$$

Procediendo de manera idéntica que en la sección 2 de este capítulo, obtenemos

$$\sum_{s=1}^g \sum_{k=1}^N \frac{\partial}{\partial \hat{v}_{ij}^s} \frac{f(\underline{x}_k)}{f(\underline{x}_k)} \hat{d}_i^s \hat{d}_j^s = 0$$

i. e.,

$$\sum_{s=1}^g \sum_{k=1}^N \left(1 - \frac{\hat{d}_{ij}^s}{2}\right) \frac{\lambda_s \alpha_s}{f(\underline{x}_k)} \left[\hat{v}_{ij}^s - (x_{ik} - \mu_{si})(x_{jk} - \mu_{sj}) \right] \hat{d}_i^s \hat{d}_j^s =$$

$$\begin{aligned}
 &= (1 - \frac{d_{ij}}{2}) \sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s \alpha_s}{f(x_k)} \left[\hat{c}^{ij} \hat{d}_s^i \hat{d}_s^j - (\alpha_{ik} - \hat{\mu}_{si})(\alpha_{jk} - \hat{\mu}_{sj}) \right] \hat{d}_i^s \hat{d}_j^s \\
 &= (1 - \frac{d_{ij}}{2}) \sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s \alpha_s}{f(x_k)} \left[\hat{c}^{ij} - (\alpha_{ik} - \hat{\mu}_{si})(\alpha_{jk} - \hat{\mu}_{sj}) \right] \hat{d}_i^s \hat{d}_j^s = 0.
 \end{aligned}$$

De donde

$$\begin{aligned}
 N \hat{c}^{ij} &= \sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s \alpha_s}{f(x_k)} (\alpha_{ik} - \hat{\mu}_{si})(\alpha_{jk} - \hat{\mu}_{sj}) \hat{d}_i^s \hat{d}_j^s = \\
 &= \sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s \alpha_s}{f(x_k)} \alpha_{ik} \alpha_{jk} \hat{d}_i^s \hat{d}_j^s + \sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s \alpha_s}{f(x_k)} \hat{\mu}_{si} \hat{\mu}_{sj} \hat{d}_i^s \hat{d}_j^s \\
 &\quad - \sum_{s=1}^g \sum_{k=1}^N \frac{\hat{\lambda}_s \alpha_s}{f(x_k)} (\alpha_{ik} \hat{\mu}_{sj} + \alpha_{jk} \hat{\mu}_{si}) \hat{d}_i^s \hat{d}_j^s \\
 &= \sum_{s=1}^g \sum_{k=1}^N \hat{P}(s|x_k) \alpha_{ik} \alpha_{jk} \hat{d}_i^s \hat{d}_j^s + N \sum_{s=1}^g \hat{\lambda}_s \hat{\mu}_{si} \hat{\mu}_{sj} \hat{d}_i^s \hat{d}_j^s \\
 &\quad - 2N \sum_{s=1}^g \hat{\lambda}_s \hat{\mu}_{si} \hat{\mu}_{sj} \hat{d}_i^s \hat{d}_j^s \\
 &= \sum_{s=1}^g \sum_{k=1}^N \hat{P}(s|x_k) \alpha_{ik} \alpha_{jk} \hat{d}_i^s \hat{d}_j^s - N \sum_{s=1}^g \hat{\lambda}_s \hat{\mu}_{si} \hat{\mu}_{sj} \hat{d}_i^s \hat{d}_j^s
 \end{aligned}$$

Por lo tanto

$$\hat{c}^{ij} = \frac{1}{N} \sum_{s=1}^g \sum_{k=1}^N \hat{P}(s|z_k) x_{ik} x_{jk} \hat{d}_i^s \hat{d}_j^s - \sum_{s=1}^g \hat{\lambda}_s \hat{\mu}_{si} \hat{\mu}_{sj} \hat{d}_i^s \hat{d}_j^s$$

que es la estimación máxima verosímil de las correlaciones.

Estos resultados los podemos condensar en el siguiente teorema.

Teorema 4. Los estimadores máximo verosímiles de los parámetros de una mezcla de distribuciones normales multivariadas con matriz de correlación común a todos los grupos son

$$(3.14) \quad \hat{\lambda}_s = \frac{1}{N} \sum_{k=1}^N \hat{P}(s|z_k)$$

$$(3.15) \quad \hat{\mu}_{si} = \frac{1}{N \hat{\lambda}_s} \sum_{k=1}^N \hat{P}(s|z_k) x_{ik}$$

$$(3.16) \quad \hat{\tau}_s^{ii} = \frac{1}{N \hat{\lambda}_s} \sum_{k=1}^N \hat{P}(s|z_k) (x_{ik} - \mu_{si})^2$$

$$(3.17) \quad \hat{c}^{ij} = \frac{1}{N} \sum_{s=1}^g \sum_{k=1}^N \hat{P}(s|z_k) x_{ik} x_{jk} \hat{d}_i^s \hat{d}_j^s - \sum_{s=1}^g \hat{\lambda}_s \hat{\mu}_{si} \hat{\mu}_{sj} \hat{d}_i^s \hat{d}_j^s$$

donde $\hat{d}_i^s = (d_s^i)^{-1} = (\sqrt{\hat{\tau}_s^{ii}})^{-1} = \sqrt{\hat{\tau}_s^{ii}^{-1}}$; c^{ij} son las correlaciones.

Es interesante hacer notar que si los grupos en la mezcla tienen varianzas comunes τ^{ii} , $i=1, 2, \dots, p$, entonces

$$\hat{c}^{ij} = \left\{ \frac{1}{N} \sum_{k=1}^N x_{ik} x_{jk} - \sum_{s=1}^g \hat{\mu}_s \hat{\mu}_{si} \hat{\mu}_{sj} \right\} / \sqrt{\hat{v}_s^{ii} \hat{v}_s^{jj}} = \frac{\hat{v}_s^{ij}}{\sqrt{\hat{v}_s^{ii} \hat{v}_s^{jj}}}$$

donde \hat{v}_s^{ij} son las estimaciones máximo verosímiles de v^{ij} suponiendo que las poblaciones en la mezcla tienen una matriz de dispersión común $\Sigma = \{v^{ij}\}$. Los estimadores \hat{c}_i^j no son estrictamente los estimadores máximo verosímiles para este modelo; sin embargo, son buenos estimadores.

Es importante señalar que el tercer caso ($\Sigma_s = a_s \Sigma$) es un caso particular del cuarto caso (matrices de correlación iguales). Para ver esto, notemos que si las matrices de dispersión son de la forma $\Sigma_s = a_s \Sigma$, entonces

$$\Sigma_s = \{v_s^{ij}\} = a_s \Sigma = a_s \{v^{ij}\} = \{a_s v^{ij}\}$$

y las correlaciones son

$$c_s^{ij} = \frac{v_s^{ij}}{\sqrt{v_s^{ii} v_s^{jj}}} = \frac{a_s v^{ij}}{\sqrt{a_s v^{ii} a_s v^{jj}}} = \frac{v^{ij}}{\sqrt{v^{ii} v^{jj}}}$$

para todo $s=1, 2, \dots, g$; es decir, que las correlaciones entre las variables en los distintos grupos son iguales.

En el tercer caso ($\Sigma_s = a_s \Sigma$) se estiman menos parámetros que en el caso de matrices de correlación iguales siempre y cuando $g(p-1) \geq p$. Para ver la justificación de esto recordemos que el número de parámetros a estimar cuando se tienen matrices de correlación iguales es $g-1 + gp + \frac{1}{2} p(p-1) + gp$, mientras que para el caso $\Sigma_s = a_s \Sigma$, se tienen que estimar

$g-1 + gp + \frac{1}{p}(p+1) + g$ y restando ambas expresiones obtenemos $g(p-1) - p$; para que esto sea positivo es necesario que $g(p-1) \geq p$, como quedamos verificar.

CAPITULO IV

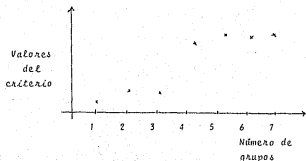
PRUEBAS SOBRE EL NUMERO DE GRUPOS EN UNA MEZCLA

Un problema común a todas las técnicas de análisis de conglomerados consiste en decidir el número de grupos a formarse. Este problema será tratado en este capítulo.

1. REVISION DE ALGUNOS METODOS.

Para decidir cuántos grupos formar al hacer un agrupamiento de objetos o individuos, varios autores han sugerido diversos métodos que dependen de la técnica utilizada al hacer dicho agrupamiento.

Para aquellas técnicas que buscan optimizar un criterio, Friedman y Rubin [1967] [8], sugieren que una gráfica del valor del criterio contra el número de grupos indicará el número más adecuado de éstos en donde haya un cambio muy pronunciado de dicho valor. Por ejemplo, la siguiente gráfica sugeriría que el número más adecuado de grupos a formar es 4:



Sin embargo, este procedimiento no ha sido satisfactorio por completo. Otros intentos se han hecho para resolver este problema. Por ejemplo, Beale (1969) [2] sugiere una estadística F definida de la manera siguiente:

$$F(C_2, C_1) = \frac{R_{C_2} - R_{C_1}}{R_{C_2}} \left/ \left[\frac{N - C_1}{N - C_2} \left(\frac{C_2}{C_1} \right)^{2/p} - 1 \right] \right.$$

basada en $p(C_2 - C_1)$ y $p(N - C_2)$ grados de libertad. En esta fórmula $R_C = (N - C)S_C^2$ y S_C^2 es la desviación cuadrática media respecto a los centros de los grupos en la muestra. Un resultado significativo sería indicador de que una división en C_2 grupos es significativamente mejor que una división en C_1 grupos ($C_2 > C_1$). Calinski y Harabasz (1971), sugirieron usar el cociente R dado por

$$R = \frac{\text{traza}(B) / (g - 1)}{\text{traza}(W) / (N - g)}$$

donde B y W son las matrices de dispersión entre grupos y dentro de grupos respectivamente. Si en esta expresión el valor de R crece monótonamente con el de g , entonces se sugiere que no existe estructura de grupos en los datos; si por el contrario, R decrece monótonamente con g , se sugiere la presencia de una estructura jerárquica en los grupos con g^* grupos, donde g^* es aquel valor de g donde se produce el cambio más pronunciado conforme variamos g . Tanto el método de Beale, como el de Calinski y Harabasz son útiles, pero sólo en situaciones donde los grupos son aproximadamente esféricos y de igual tamaño.

En 1971, Harriot investigó en detalle las propiedades del 'Criterio $|W|$ '. El propone el uso de $g^2/|W|$ donde g es el número correcto de grupos, como un indicador de la estructura de datos, tomando como el número correcto de grupos aquel valor de g para el que $g^2/|W|$ es mínimo. Asimismo, Harriot investigó las propiedades muestrales de la estadística asociada, es decir de $g^2/|W|/T$, donde la población tiene una distribución uniforme, y concluye que a partir de ésta se puede saber si hay alguna estructura de grupos en los datos, o bien, si se está tratando con un solo grupo. Este método puede aplicarse a conglomerados más generales que en los casos anteriores, puesto que no hay suposición alguna sobre la forma de los grupos, aunque sí se requiere que éstos sean aproximadamente del mismo tamaño.

Para las técnicas jerárquicas de análisis de conglomerados no existen indicadores claros del número más adecuado de grupos, y puede ser, que en los tipos de datos para los que estas técnicas sean las más apropiadas, el problema del número de grupos no sea de gran interés. Pero en aquellos casos en donde sí necesitamos tener algún indicador del número correcto de grupos, el simple análisis de los dendogramas no es muy satisfactorio.

En el caso de ajustar mezclas de distribuciones multivariadas a un conjunto de datos, existe un criterio, mucho más razonable que los anteriores, para investigar el número correcto de grupos en la muestra; este criterio es la Prueba del Cociente de Verosimilitud, el cual estudiaremos a continuación.

2. PRUEBA DEL COCIENTE DE VEROSIMILITUD.

Para probar la hipótesis de r grupos contra s grupos ($r < s$), podemos usar la estadística

$$\Lambda = - 2 \log \lambda$$

donde λ es el cociente de las verosimilitudes al tomar

s y r grupos, es decir

$$\lambda = L_s / L_r$$

donde L_s y L_r son las funciones de verosimilitud de la muestra bajo las hipótesis de s y r grupos, respectivamente.

Tanto Wilks (1938) [16] como Silvey (1959) [17], demostraron que bajo ciertas condiciones de regularidad, la estadística Λ se distribuye como una χ^2 con grados de libertad igual a la diferencia de los números de parámetros bajo las dos hipótesis. Bartlett (1959) [7], al investigar el problema de probar la igualdad de r medias en análisis de varianza multivariado, mejoró el resultado de Wilks, para muestras pequeñas usando la estadística modificada

$$\Lambda' = -2 C \log \lambda$$

con $p(s - r)$ grados de libertad y $C = (N - 1 - (p+1)/2)/N$, donde N es el tamaño de la muestra y p es el número de variables.

En 1970, Wolfe [19] notó que la fórmula de Wilks no funciona correctamente. Tanto para los datos de "Fisher-Iris" como para unos conglomerados generados artificialmente, la prueba de Wilks rechazaba la hipótesis nula cuando era verdadera. Un año después, el mismo Wolfe

señaló que la causa de que hubiera ocurrido esto, es que algunas de las condiciones de regularidad pedidas en el Teorema de Wilks no son satisfechas por el cociente de verosimilitud en el problema de mezclas.

Para probar la hipótesis de r grupos contra s grupos, Wolfe recomienda una modificación al cociente de verosimilitud. Esta modificación consiste en usar la estadística

$$\Lambda^* = -\frac{2}{N} \left(N - 1 - \rho - \frac{s}{2} \right) \log \lambda$$

como una distribución χ^2 con $2\rho(s-r)$ grados de libertad si se hace la suposición de matrices de dispersión iguales, $2 \left[\rho \left(\frac{\rho+1}{2} + \rho \right) \right] (s-r)$ grados de libertad en el caso general y $2(2\rho)(s-r)$ grados de libertad para el caso de matrices de correlación iguales. Esta recomendación es sólo heurística, no tiene justificación teórica o numérica.

CAPITULO V

USO DE NORMIX

Tal y como lo señalamos en los capítulos anteriores, en este capítulo presentamos la manera como debe hacerse un agrupamiento de un conjunto de datos usando la técnica de mezclas. Para ello usamos un programa cuya versión original se debe a J.H. Wolfe (1967) (Ver 18) titulado NORMIX. La versión utilizada en el TIMAS incorpora algunas modificaciones hechas por A. Hernández.

Este último capítulo consta de dos partes: en la primera presentamos las características generales de este programa así como un manual de usuario para el mismo; en la segunda parte presentamos un ejemplo de un agrupamiento hecho en los datos de Iris de Fisher.

1. CARACTERISTICAS GENERALES Y MANUAL DE USUARIO.

El objetivo de este programa consiste en encontrar los estimadores máximo verosímiles de una mezcla de distribuciones normales multivariadas; una vez teniendo éstos se

lleva a cabo el agrupamiento de los datos de acuerdo a las probabilidades de pertenencia de cada uno de los objetos o individuos a cada uno de los grupos. Existe una opción que permite al usuario especificar si las matrices de dispersión son iguales o diferentes (aún no han sido incluidas las opciones para los casos 3 y 4 del capítulo anterior).

Tal y como vimos en el Capítulo II, para hacer un agrupamiento en un conjunto de datos por medio de la técnica de mezclas, es necesario contar con estimaciones iniciales para todos los parámetros involucrados en la mezcla. En la versión de Wolfe se obtenían estas estimaciones usando técnicas jerárquicas de agrupamiento. La versión existente en el THAS, incluye la opción de obtener estas estimaciones minimizando la traza de W o bien minimizando el determinante de W , donde W es la matriz de dispersión dentro de grupos. Esta modificación provoca que la velocidad de convergencia de los métodos iterativos al resolver las ecuaciones máximo verosímiles sea mucho más rápida.

Si el número de variables es muy grande, los cálculos dentro de la computadora son muy pesados. Sin embargo, esta no es una dificultad tan seria, puesto que dado que los resultados son independientes de transformaciones lineales, las variables observadas pueden ser reemplazadas por un número menor de componentes principales, y el programa da los mismos resultados, modificados solamente por la pér-

didada de información en los componentes principales descartados.

Tal y como ya lo indicamos, existe una opción que permite especificar si las matrices de dispersión son iguales o distintas. Las ventajas de los dos enfoques son similares a aquellas de usar funciones discriminantes lineales o cuadráticas; el análisis de conglomerados trata con diferencias en la localización y en la práctica la versión con matrices de dispersión distintas puede producir estimaciones muy diferentes a las obtenidas con matrices de dispersión iguales, aunque las reglas de colocación o asignación al formar los grupos sean muy similares para ambos casos.

La salida del programa nos da para cada valor de g (el número de grupos), la verosimilitud, la estimación de las proporciones y de las probabilidades de pertenencia para cada uno de los grupos. Si los grupos están muy separados la convergencia es muy rápida y la mayoría de los puntos son asignados a un grupo con probabilidad cercana a uno. Si no lo están, la convergencia es muy lenta, además de que pueden resultar varios máximos locales; las estimaciones dependen de los valores iniciales escogidos, por lo que es muy recomendable hacer varias corridas para un mismo conjunto de datos con diferentes estimaciones iniciales. El hecho de que haya varios máximos locales, no es una dificultad muy seria, desde el punto de vista de Análisis de Conglomerados,

pues la mayoría de las veces se trata de dividir distribuciones multimodales más que de dividir distribuciones platidénticas.

Para correr NORMIX en el IHAS, la instrucción para hacerlo es:

```
R$BINNORMIX4;FILE FILE5(DISK, TITLE=nombre del programa,
FILETYPE=7), FILE1(DISK, TITLE=nombre del archivo de datos,
FILETYPE=7).
```

Las instrucciones para cada uno de los renglones del programa son:

RENGLÓN	ESPECIFICACIONES
100	Los primeros cuatro renglones se usan para algún título del programa, comentarios u observaciones que se desee aparezcan al principio de cada corrida.
200	
300	Si no desean utilizarse tienen que dejarse en blanco.
400	

500	Columnas 34-35: Número de variables medidas a cada individuo.
	Columnas 50-53: Número de individuos u observaciones; tiene que ser un número del 1 al 9999.
	Columnas 69-70: Número del archivo de datos. Generalmente se elige el 01.

- 600 Columnas 32 en adelante: número de grupos a formar; se escriben en forma ascendiente, separados por una coma; deben ser números del 01 al 99.
-
- 700 Columna 44: se escribe un 0 si se desea la opción de matrices de dispersión iguales; se escribe un 1 si se desea la opción de matrices de dispersión diferentes.
- Columnas 68-70: mínimo número de puntos que debe tener un grupo para que pueda ser estimada su matriz de dispersión. Si el tamaño de algún grupo resulta menor que este número se toma a la matriz de dispersión entre grupos como matriz de dispersión para dicho grupo.
-
- 800 Columnas 70-72: el programa formará un número mayor de grupos si la prueba resulta significativa. Esta prueba consiste en ir probando un grupo más que el último especificado. En estas columnas se pone el valor de α , la significancia de la prueba.
-
- 900 Columnas 36-38: Se usa solamente en el caso en que

se desea un agrupamiento inicial de tipo jerárquico. Se especifica el tamaño de la submuestra que se usa para generar estimaciones iniciales. Debe ser un número del 040 al 250.

Columna 58: si se escribe un 1 entonces en la salida del programa se imprimen los resultados de cada iteración. Desde luego esto origina que se tenga una salida muy grande.

Columna 68: si se escribe un 1 entonces el usuario puede identificar cada observación. Para esto las etiquetas deben estar en las primeras n columnas ($1 \leq n \leq 8$) del archivo de datos. Por lo tanto, el formato en los renglones 13 y/o 14 debe ser de la forma $(A_n, \text{formato de las observaciones})$ y la etiqueta debe ocupar n caracteres en el archivo de datos. si se escribe un 0 no hay alteraciones y las observaciones son numeradas e identificadas de acuerdo al orden en que se leen. Desde luego que es posible modificar el orden en que se leen las etiquetas en el archivo de datos alterando la instrucción de lectura del archivo de datos.

1000 Columnas 25-16: En estas columnas se especifica un límite de tiempo (en minutos) para que el programa tenga tiempo de imprimir los resultados. Este tiempo debe ser alrededor de un minuto menos que el límite de tiempo del JOB correspondiente, con el objeto de que se puedan imprimir los resultados antes de que el sistema operativo concluya el JOB. Si se dejan en blanco, es probable que el listado de resultados no esté completo.

Columnas 57-59: Si se dejan en blanco el programa iterará hasta que se obtenga convergencia. Esto toma de 1 a 100 iteraciones. Si no se dispone de tiempo máquina suficiente se pueden indicar en estas columnas un número menor de iteraciones. Generalmente con 010 iteraciones es más que suficiente. Si se escribe aún un número menor de iteraciones es posible que ya no se tengan los estimadores máximo verosímiles, aunque los obtenidos serán siempre mejores que los que se obtienen del agrupamiento inicial.

- 1100 Columna 34: Si se deja en blanco se hace un agrupamiento inicial jerárquico. Si se escribe un 1 se hará un agrupamiento inicial no jerárquico: si se escribe un 1 en la columna 46 y un 0 en la 58 el método es minimizando la traza de W ; si se escribe un 2 en la columna 46 y un 1 en la 58 el método es minimizando el determinante de W .
-
- 1200 Se deja en blanco
-
- 1300 Columna 1 en adelante: Formato de entrada; o sea el formato con el que se leerá el archivo de datos. Va escrito entre paréntesis.
-
- 1400 Sólo se usa en caso de que el usuario dé estimaciones iniciales para los parámetros de la mezcla. En este caso, en las columnas 26-27 se escribe el número de grupos del conjunto de estimaciones iniciales. Además en las columnas 44, 55 y 70 se escribe un 1 si medias iniciales, desviaciones estándar iniciales o correlaciones iniciales van a ser leídas respectivamente. En este caso el formato de

entrada es 8610.4 y ya no se tiene que especificar. Las estimaciones iniciales se escriben los renglones sucesivos con este formato, escribiendo primero todas las medias, después desviaciones estándar y finalmente correlaciones. En caso de que no se den estimaciones iniciales, este renglón se deja en blanco.

1500 Columna 1-2: Se escribe lo siguiente: /* (diagonal asterisco) que indica el fin de los comandos. Se debe tener un renglón en blanco antes de éste.

Después de esto ya se puede correr el programa con la instrucción dada para hacerlo.

En la siguiente sección de este capítulo presentamos un ejemplo de como hacer un agrupamiento en un conjunto de datos usando este programa. Se presentan también los resultados obtenidos.

2. ANALISIS DE LOS DATOS DE IRIS DE FISHER.

En esta sección se analiza un conjunto de datos publicados por Fisher en 1936, correspondientes a tres varie-

dades de Iris, a saber, Iris Setosa, Iris Versicolor e Iris Virginica. Se tienen 50 observaciones para cada una de las variedades y para facilitar el análisis las hemos numerado de la manera siguiente: con los números 1 al 50 las observaciones correspondientes a Iris setosa; con los números 51 al 100 las correspondientes a Iris versicolor y con los números 101 al 150 las correspondientes a Iris virginica. Cada observación consta de 4 mediciones hechas a cada uno de los ejemplares de Iris: longitud del sépalo (LS), ancho del sépalo (AS), longitud del pétalo (LP) y ancho del pétalo (AP). Dos de las especies, Iris setosa e Iris Versicolor fueron encontradas en la misma colonia por el Botánico E. Anderson. La especie Iris virginica fue encontrada en una colonia distinta. El mismo Fisher (1936) aplicó su procedimiento de análisis discriminante a los tres grupos anteriores y encontró que la Iris setosa podía separarse muy bien de las otras especies, sin embargo, las especies Iris Versicolor e Iris Virginica tenían un traslape considerable y podrían confundirse algunos elementos de una especie con los de la otra.

En la siguiente tabla se muestran este conjunto de datos.

DATOS DE LAS TRES VARIEDADES DE TRIS

TRIS SETOSA					TRIS VERSICOLOR					TRIS VIRGINICA				
	LS	AS	LP	AP		LS	AS	LP	AP		LS	AS	LP	AP
1	5.1	3.5	1.4	0.2	51	7.0	3.2	4.7	1.4	101	6.3	3.3	6.0	2.5
2	4.9	3.0	1.4	0.2	52	6.4	3.2	4.5	1.5	102	5.8	2.7	5.1	1.9
3	4.7	3.2	1.3	0.2	53	6.9	3.1	4.9	1.5	103	7.1	3.0	5.9	2.1
4	4.6	3.1	1.5	0.2	54	5.5	2.3	4.0	1.3	104	6.3	2.9	5.6	1.8
5	5.0	3.6	1.4	0.2	55	6.5	2.4	4.6	1.5	105	6.5	3.0	5.8	2.2
6	5.4	3.9	1.7	0.4	56	5.7	2.8	4.5	1.3	106	7.6	3.0	6.6	2.1
7	4.6	3.4	1.4	0.3	57	6.3	3.3	4.7	1.6	107	4.9	2.5	4.5	1.7
	5.0	3.4	1.5	0.2	58	4.9	2.4	3.3	1.0	108	7.3	2.9	6.3	1.8
	4.4	2.9	1.4	0.2	59	6.6	2.9	4.6	1.3	109	6.7	2.5	5.8	1.8
10	4.9	3.1	1.5	0.1	60	5.2	2.7	3.9	1.4	110	7.2	3.6	6.1	2.5
11	5.4	3.7	1.5	0.2	61	5.0	2.0	3.5	1.0	111	6.5	3.2	5.1	2.0
	4.8	3.4	1.6	0.2	62	5.9	3.0	4.2	1.5	112	6.4	2.7	5.3	1.9
	4.8	3.0	1.4	0.1	63	6.0	2.2	4.0	1.0	113	6.8	3.0	5.5	2.1
	4.3	3.0	1.1	0.1	64	6.1	2.9	4.7	1.4	114	5.7	2.5	5.0	2.0
	5.8	4.0	1.2	0.2	65	5.6	2.9	3.6	1.3	115	5.8	2.8	5.1	2.4
16	5.7	4.4	1.5	0.4	66	6.7	3.1	4.4	1.4	116	6.4	3.2	5.3	2.3
17	5.4	3.9	1.3	0.4	67	5.6	3.0	4.5	1.5	117	6.5	3.0	5.5	1.8
18	5.1	3.5	1.4	0.3	68	5.8	2.7	4.1	1.0	118	7.7	3.8	6.7	2.2
19	5.7	3.8	1.7	0.3	69	6.2	2.2	4.5	1.5	119	7.7	2.6	6.9	2.3
20	5.1	3.8	1.5	0.3	70	5.6	2.5	3.9	1.1	120	6.0	2.2	5.0	1.5
21	5.4	3.4	1.7	0.2	71	5.9	3.2	4.8	1.8	121	6.9	3.2	5.7	2.3
22	5.1	3.7	1.5	0.4	72	6.1	2.8	4.0	1.3	122	5.6	2.8	4.9	2.0
23	4.6	3.6	1.0	0.2	73	6.3	2.5	4.9	1.5	123	7.7	2.8	6.7	2.0
	5.1	3.3	1.7	0.5	74	6.1	2.8	4.7	1.2	124	6.3	2.7	4.9	1.8
25	4.8	3.4	1.9	0.2	75	6.4	2.9	4.3	1.3	125	6.7	3.3	5.7	2.1
26	5.0	3.0	1.6	0.2	76	6.6	3.0	4.4	1.4	126	7.2	3.2	6.0	1.8
27	5.0	3.4	1.6	0.4	77	6.8	2.8	4.6	1.4	127	6.2	2.8	4.8	1.8
28	5.2	3.5	1.5	0.2	78	6.7	3.0	5.0	1.7	128	6.1	3.0	4.9	1.8
29	5.2	3.4	1.4	0.2	79	6.0	2.9	4.5	1.5	129	6.4	2.8	5.6	2.1
30	4.7	3.2	1.6	0.2	80	5.7	2.6	3.5	1.0	130	7.2	3.0	5.8	1.6
31	4.8	3.1	1.6	0.2	81	5.5	2.4	3.8	1.1	131	7.4	2.8	6.1	1.9
32	5.4	3.4	1.5	0.4	82	5.5	2.4	3.7	1.0	132	7.9	3.8	6.4	2.0
33	5.2	4.1	1.5	0.1	83	5.8	2.7	3.9	1.2	133	6.4	2.8	5.6	2.2
34	5.5	4.2	1.4	0.2	84	6.0	2.7	5.1	1.6	134	6.3	3.8	5.1	1.5
35	4.9	3.1	1.5	0.2	85	5.4	3.0	4.5	1.5	135	6.1	2.6	5.6	1.4
36	5.0	3.2	1.2	0.2	86	6.0	3.4	4.5	1.6	136	7.7	3.0	6.1	2.3
37	5.5	3.5	1.3	0.7	87	6.7	3.1	4.7	1.5	137	6.5	3.4	5.4	2.4
38	4.9	3.6	1.4	0.1	88	6.3	2.3	4.4	1.3	138	6.4	3.1	5.5	1.8
39	4.4	3.0	1.3	0.2	89	5.6	3.0	4.1	1.3	139	6.0	3.0	4.8	1.8
40	5.1	3.4	1.5	0.2	90	5.5	2.5	4.0	1.3	140	6.9	3.1	5.4	2.1
41	5.0	3.5	1.3	0.3	91	5.5	2.6	4.4	1.2	141	6.7	3.1	5.6	2.4
42	4.5	2.3	1.3	0.3	92	6.1	3.0	4.6	1.4	142	6.9	3.1	5.1	2.3
43	4.4	3.2	1.3	0.2	93	5.9	2.6	4.0	1.2	143	5.8	2.7	5.1	1.9
44	5.0	3.5	1.6	0.6	94	5.0	2.3	3.3	1.0	144	6.8	3.2	5.9	2.3
45	5.1	3.6	1.9	0.4	95	5.6	2.7	4.2	1.3	145	6.7	3.3	5.7	2.5
46	4.8	3.0	1.4	0.3	96	5.7	3.0	4.2	1.2	146	6.7	3.0	5.2	2.3
47	5.1	3.8	1.6	0.2	97	5.7	2.9	4.2	1.3	147	6.5	2.5	5.0	1.9
48	4.6	3.2	1.4	0.2	98	6.1	2.8	4.3	1.3	148	6.5	3.0	5.2	2.0
49	5.3	3.7	1.5	0.2	99	5.1	2.5	3.0	1.1	149	6.7	3.4	5.4	2.3
50	5.0	3.3	1.4	0.2	100	5.7	2.8	4.1	1.3	150	5.9	3.0	5.1	1.8

Para este conjunto de datos hicimos varias corridas de NORMIX. Estas corridas fueron hechas con las opciones de matrices de dispersión iguales y diferentes; también variamos los métodos de agrupamiento inicial: jerárquico, minimización del determinante de w y minimización de la traza de W . A continuación presentamos los resultados de estas corridas. En la tabla 1 presentamos las pruebas sobre el número de grupos, usando hipótesis de un grupo, dos grupos, tres grupos y cuatro grupos. Cada hipótesis se prueba contra la anterior usando la modificación propuesta por Wolfe al cociente de verosimilitud, a saber

$$\chi^2 = -\frac{2}{N} (N-1-p-\frac{s}{2}) \log \frac{L_s}{L_r}$$

con $2p(s-1)$ grados de libertad si se hace la suposición de matrices de dispersión iguales y $2(\frac{p(p+1)}{2} + p)(s-1)$ grados de libertad si se hace la suposición de matrices de dispersión diferentes. Desde luego que en este caso tenemos $p=4$ y las pruebas de hipótesis fueron hechas con $s-1=1$, por lo que tenemos 8 grados de libertad en el caso de matrices de dispersión iguales y 28 grados de libertad en el caso de matrices de dispersión diferentes. En la tabla 2 presentamos las estimaciones para algunos parámetros de la mezcla (de las proporciones y las medias) y los comparamos con los valores reales que el mismo Fisher publicó en su citado artículo. Esta tabla incluye también el número de individuos mal clasificados con cada uno

de los métodos y cuales fueron estos individuos mal clasificados en cada caso para poder hacer una comparación posterior de los métodos.

TABLA 1

Pruebas de significancia sobre el número de grupos

MATRICES DE DISPERSION IGUALES

Número de grupos: H_0/H_a	JERARQUICO		MIN DET W		MIN TR W	
		P		P		P
1/2	-	-	151.61	10^{-8}	151.61	10^{-8}
2/3	40.71	.000002	59.75	10^{-8}	59.75	10^{-8}
3/4	95.93	10^{-8}	12.04	.034	12.04	.034

MATRICES DE DISPERSION DIFERENTES

Número de grupos: H_0/H_a	JERARQUICO		MIN DET W		MIN TR W	
		P		P		P
1/2	-	-	-	-	-	-
2/3	38.18	0.09505	49.98	.00640	49.25	.00783
3/4	58.64	.000606	20.89	.82992	16.92	.95010

Notas: 1) La P es la probabilidad de la hipótesis nula ante la alternativa.

2) Los valores faltantes no los dio la máquina; son muy altos los valores de χ^2 .

TABLA 2
Estimación de parámetros de Iris

I. IRIS SETOSA

PARAMETRO	FISHER	MATRICES DE DISPERSION IGUALES			MATRICES DE DISPERSION DIFERENTES		
		Jordanquico	Mín det W	Mín et W	Jordanquico	Mín det W	Mín et W
Proporción	0.3333	0.3333	0.3333	0.3333	0.3333	0.3333	0.3333
Medía [LS]	5.0060	5.0060	5.0060	5.0060	5.0060	5.0060	5.0060
Medía [AS]	3.4280	3.4280	3.4280	3.4280	3.4280	3.4280	3.4280
Medía [LP]	1.4620	1.4620	1.4620	1.4620	1.4620	1.4620	1.4620
Medía [AP]	0.2460	0.2460	0.2460	0.2460	0.2460	0.2460	0.2460
# mat clasificados	0	0	0	0	0	0	0

II. IRIS VESICOLOR

PARAMETRO	FISHER	MATRICES DE DISPERSION IGUALES			MATRICES DE DISPERSION DIFERENTES		
		Jordanquico	Mín det W	Mín et W	Jordanquico	Mín det W	Mín et W
Proporción	0.3333	0.5200	0.5153	0.2973	0.7470	0.2850	0.2870
Medía [LS]	6.5880	7.0045	6.5982	6.6454	6.9766	6.6733	6.6875
Medía [AS]	2.9740	2.9318	3.0081	3.0120	2.9216	2.9980	3.0005
Medía [LP]	5.5520	5.5591	5.5587	5.5926	5.3929	5.6413	5.6500
Medía [AP]	0.2060	1.6591	2.0567	2.0757	1.6607	2.0679	2.0697
# mat clasificados	0	3	1	0	12	0	0

III. IRIS VIRGINICA

PARAMETRO	FISHER	MATRICES DE DISPERSION IGUALES			MATRICES DE DISPERSION DIFERENTES		
		Jordanquico	Mín det W	Mín et W	Jordanquico	Mín det W	Mín et W
Proporción	0.3333	0.1467	0.3514	0.3694	0.5260	0.5810	0.3850
Medía [LS]	5.9360	6.0526	5.9607	5.9535	6.0721	5.9544	5.9550
Medía [AS]	2.7700	2.8531	2.7555	2.7594	2.8587	2.7788	2.7781
Medía [LP]	4.2600	4.7782	4.3244	4.3545	4.7755	4.3561	4.3625
Medía [AP]	1.3260	1.6802	1.3349	1.3544	1.6801	1.3829	1.3864
# mat clasificados	0	56	4	4	40	6	6

Total de mat clasificados	0	39	5	4	52	6	6
---------------------------	---	----	---	---	----	---	---

Las pruebas de significancia sobre el número de grupos indican que la hipótesis de un grupo debe ser rechazada ante la alternativa de dos grupos; también indican que la hipótesis de dos grupos debe ser rechazada ante la alternativa de tres grupos; si acaso el único método donde se presenta una duda es el de agrupamiento jerárquico con matrices de dispersión diferentes donde la P tiene un valor de 0.09505 por lo que no podíamos rechazar al 95 % de confianza. Asimismo, estas pruebas indican (con excepción de dos métodos) que no debe rechazarse la hipótesis de 3 grupos ante la alternativa de 4 grupos. Los dos métodos en donde esto no es tan claro son los de agrupamiento por \min det W y \min tr W con matrices de dispersión iguales, en donde la P tiene un valor de .034.

La tabla 2 presenta los estimadores máximo verosímiles de algunos parámetros de la mezcla (proporciones y medias) que se obtienen con NORMIX y se comparan con los valores reales de dichos parámetros que el mismo Fisher proporcionó. El renglón de individuos mal clasificados da el número de ejemplares de Iris cuyas probabilidades de pertenencia fueron más altas para otros grupos que para el que pertenecían en realidad por lo que quedaron mal clasificados. Excepto para los métodos en donde se realizó un agrupamiento inicial jerárquico, las estimaciones que se obtienen con NORMIX fueron muy parecidas a las reales; con la opción de matrices iguales se obtienen

resultados ligeramente mejores que con la opción de matrices de dispersión diferentes. El método que hizo el mejor agrupamiento fue el de un agrupamiento inicial por $\min tr W$ con matrices de dispersión iguales en donde solamente se malclasificaron 4 ejemplares de Iris, a saber 120, 130, 134 y 135 que siendo de la especie Iris Virgínica los incluyó el programa en la especie Iris Versicolor. El método de $\min det W$ inicial con matrices de dispersión iguales, clasificó mal a los mismos individuos además del individuo 71. Los métodos $\min det W$ y $\min tr W$ inicial con matrices de dispersión diferentes clasificaron mal a seis ejemplares, a saber, los marcados con los números 124, 127, 128, 134, 139 y 150 que siendo de la especie Iris Virgínica fueron incluidos en la especie Iris Versicolor. Como un comentario final, debe indicarse que los métodos con agrupamiento jerárquico inicial dieron tan malos agrupamientos debido a que los datos no presentan para nada un esquema jerárquico de agrupamiento.

Con el objeto de averiguar que sucedía con estos mismos datos si primero se obtenían componentes principales, se corrió también NORRIX para los primeros dos componentes principales que resultan de estos datos y que explican el 96 % de la variabilidad total existente en estos datos, resultando que el método que dio mejores resultados fue $\min tr W$ inicial y matrices de dispersión diferentes. Con este método se mal clasi-

ficaron 13 individuos, los marcados con los números: 78, 102, 107, 114, 122, 124, 127, 128, 134, 135, 139, 143, 150. Es decir, se volvieron a traslapar algunos ejemplares de Iris Versicolor e Iris Virginica, mientras que los de Iris Setosa fueron agrupados perfectamente.

Finalmente y con propósitos de ilustración, presentamos una corrida de NORNIX así como sus resultados. Elegimos la opción de matrices de dispersión diferentes y un agrupamiento inicial por minimización de la traza de W. Limitamos también el número de iteraciones a 5 para que no se tuviera una salida muy grande.

Las instrucciones de entrada quedan de la manera siguiente:

```

100 ANALISIS DE LOS DATOS DE IRIS DE FISHER
200 AGRUPAMIENTO NO JERARQUICO INICIAL; MATRICES DE DISPERSION DIFERENTES
300 5 ITERACIONES COMO MAXIMO
400 SEGUNDA CORRIDA
500                                04                                0150                                01
600                                02,03,04
700                                1                                005
800                                1                                .90
900
1000                               01
1100                               1                                1                                005
1200                               1                                1                                0
1300 (2X,4F4.1)
1400
1500 /*

```

La salida del programa queda así: (con esto concluimos el capítulo)

PROGRAM NAME: FSC CLUSTER ANALYSIS
 MAXIMUM LEVEL: 1000 (INITIAL) OF 4 (1000) OF MULTIVARIATE NORMAL DISTRIBUTIONS

ANALYSIS OF LOS DATOS DE IRIS DE FISHER
 PROGRAM: FSC CLUSTER ANALYSIS
 INITIAL: 1000 (INITIAL) OF 4 (1000) OF MULTIVARIATE NORMAL DISTRIBUTIONS
 SECOND: 1000 (INITIAL) OF 4 (1000) OF MULTIVARIATE NORMAL DISTRIBUTIONS

STORAGE ALLOCATION: 141204 (INITIAL) OF 141204
 EXPANSION AT 141204

1	7	31	91	141	183	231	215	231	247	263
413	563	713	863	1464	1612	1765	1845	1924	1991	2067

15727

A NON-HIERARCHICAL INITIAL GROUPING WAS PROVIDED

(THE INITIAL DISTANCE WILL BE USED
 THE DISTANCE IS IN ARBITRARY UNITS)

THE CRITERION VALUE FOR THE INITIAL GROUPING FOR 2 TYPES IS .122242145

SAMPLE SIZE = 15
 NUMBER OF TYPES = 2

ITERATION NUMBER = 1

LOG-LIKELIHOOD OF 2 TYPES IN THIS SAMPLE = .0

CHARACTERISTICS OF THE WHOLE SAMPLE

MEANS			
1.8423	2.1573	3.7381	4.1993
STANDARD DEVIATIONS			
1.8241	1.4350	1.7655	1.7427
CORRELATIONS			
1.0000	-.1176	-.0343	-.0179
-.1176	1.0000	-.0343	-.0179
-.0343	-.0343	1.0000	-.0179
-.0179	-.0179	-.0179	1.0000

CHARACTERISTICS OF TYPE 1

THE PROPORTION OF THE POPULATION FROM THIS TYPE = .667

MEANS			
1.5230	2.0866	4.9588	4.6950
STANDARD DEVIATIONS			
1.8227	0.3266	1.7318	1.6156
CORRELATIONS			
1.0000	0.2301	-.0179	0.0113
0.2301	1.0000	-.0179	0.0113
-.0179	-.0179	1.0000	-.0179
0.0113	0.0113	-.0179	1.0000

CHARACTERISTICS OF TYPE 2

THE PROPORTION OF THE POPULATION FROM THIS TYPE = .333

MEANS			
2.0037	2.2698	1.5814	1.2936
STANDARD DEVIATIONS			
1.7427	1.4350	1.4414	1.2116
CORRELATIONS			
1.0000	0.4230	-.0343	-.0179
0.4230	1.0000	-.0343	-.0179
-.0343	-.0343	1.0000	-.0179
-.0179	-.0179	-.0179	1.0000

ITERATION 1
 LOG-LIKELIHOOD OF 2 TYPES IN THIS SAMPLE = .0
 LOG-LIKELIHOOD OF 2 TYPES IN THIS SAMPLE = .0
 LOG-LIKELIHOOD OF 2 TYPES IN THIS SAMPLE = .0

ITERATION TIME = 2.02 SECONDS
 TOTAL TIME USED = 2.02 SECONDS



SAMPLE SIZE N = 101
 NUMBER OF VARIABLES = 3
 NUMBER OF TYPES = 3

ESTIMATION METHOD =

LOGLIKELIHOOD OF 3 TYPES IN THIS SAMPLE = .3371962112

CHARACTERISTICS OF THE WHOLE SAMPLE

MEANS			
1	2	3	4
5.6633	5.532	5.754	5.1991
STANDARD DEVIATIONS			
0.7265	0.4259	0.7655	0.7922
CORRELATIONS			
1	2	3	4
1	-.1170	-.0713	-.0177
-.1170	1	-.4245	-.0288
-.0713	-.4245	1	-.0420
-.0177	-.0288	-.0420	1

CHARACTERISTICS OF TYPE 1

THE PROPORTION OF THE POPULATION FROM THIS TYPE = .333

MEANS			
1	2	3	4
5.6416	5.4296	5.451	5.2466
STANDARD DEVIATIONS			
0.3525	0.3791	0.1727	0.1154
CORRELATIONS			
1	2	3	4
1	0.7435	0.2077	0.2248
0.7435	1	0.1771	0.2316
0.2077	0.1771	1	0.1171
0.2248	0.2316	0.1171	1

CHARACTERISTICS OF TYPE 2

THE PROPORTION OF THE POPULATION FROM THIS TYPE = .412

MEANS			
1	2	3	4
5.9316	5.7484	6.2925	5.4159
STANDARD DEVIATIONS			
0.4664	0.2963	0.1199	0.2975
CORRELATIONS			
1	2	3	4
1	0.4619	0.1022	0.2530
0.4619	1	0.1199	0.1171
0.1022	0.1199	1	0.1171
0.2530	0.1171	0.1171	1

CHARACTERISTICS OF TYPE 3

THE PROPORTION OF THE POPULATION FROM THIS TYPE = .255

MEANS			
1	2	3	4
6.8527	5.737	6.7671	6.711
STANDARD DEVIATIONS			
0.4942	0.2911	0.4064	0.2799
CORRELATIONS			
1	2	3	4
1	0.1727	0.7365	0.2511
0.1727	1	0.462	0.1441
0.7365	0.462	1	0.1441
0.2511	0.1441	0.1441	1

ESTIMATION	1	LOG LIKELIHOOD OF 3 TYPES IN THIS SAMPLE =	.3371962112
ESTIMATION	2	LOG LIKELIHOOD OF 3 TYPES IN THIS SAMPLE =	.3271962112
ESTIMATION	3	LOG LIKELIHOOD OF 3 TYPES IN THIS SAMPLE =	.3171962112
ESTIMATION	4	LOG LIKELIHOOD OF 3 TYPES IN THIS SAMPLE =	.3071962112

ESTIMATION TIME = 0.57 SECONDS
 TOTAL TIME USED = 0.45 SECONDS

SAMPLE SIZE N = 101
 NUMBER OF VARIABLES = 3
 NUMBER OF TYPES = 3

ESTIMATION METHOD =

LOGLIKELIHOOD OF 3 TYPES IN THIS SAMPLE = .3621962112

CHARACTERISTICS OF THE WHOLE SAMPLE

1	2	3	4
0.0433	0.0573	0.7514	0.1971
STANDARD DEVIATIONS			
0.0081	0.0359	1.7655	0.7632
CORRELATIONS			
-0.1977	-0.1976	-0.0272	-0.0127
-0.0194	-0.1267	0.0000	0.0000
-0.0194	-0.1267	0.0000	0.0000

CHARACTERISTICS OF TYPE 1

THE PROPORTION OF THE POPULATION FROM THIS TYPE = 0.153

1	2	3	4
0.0760	0.0250	0.4621	0.2443
STANDARD DEVIATIONS			
0.0148	0.0353	0.7179	0.3483
CORRELATIONS			
0.1977	0.1976	0.0272	0.0127
0.0194	0.1267	0.0000	0.0000
0.0194	0.1267	0.0000	0.0000

CHARACTERISTICS OF TYPE 2

THE PROPORTION OF THE POPULATION FROM THIS TYPE = 0.207

1	2	3	4
0.0955	0.0711	0.3025	0.2394
STANDARD DEVIATIONS			
0.0353	0.0532	0.5107	0.2432
CORRELATIONS			
0.1977	0.1976	0.0272	0.0127
0.0194	0.1267	0.0000	0.0000
0.0194	0.1267	0.0000	0.0000

CHARACTERISTICS OF TYPE 3

THE PROPORTION OF THE POPULATION FROM THIS TYPE = 0.207

1	2	3	4
0.0810	0.0405	0.6519	0.2067
STANDARD DEVIATIONS			
0.0135	0.0244	0.5751	0.2609
CORRELATIONS			
0.1977	0.1976	0.0272	0.0127
0.0194	0.1267	0.0000	0.0000
0.0194	0.1267	0.0000	0.0000

PROBABILITIES OF TYPE MEMBERSHIP

1	2	3
---	---	---

PROBABILITIES OF TYPE MEMBERSHIP

1	2	3
---	---	---

PROBABILITIES OF TYPE MEMBERSHIP

1	2	3
---	---	---

PROBABILITIES OF TYPE MEMBERSHIP

1	2	3
---	---	---

INDIVIDUALS IN GROUP NUMBER 1

1	2	3	4	5	6	7	8	9	10	11	12	13	14
0.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0

INDIVIDUALS IN GROUP NUMBER 2

1	2	3	4	5	6	7	8	9	10	11	12	13	14
0.0	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1

101	102	103	104	105	106	107	108	109	110	111	112	113	114	115	116	117	118	119	120
-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

REGRESSION OF LEVEL 1 ON RATED OF 3 TO 2 TYPES = .2565 732-1-2
 STANDARD DEVIATION OF ESTIMATE OF REGRESSION = .40125
 CORRELATION OF NULL HYPOTHESIS = .1733342

THE CRITERION VALUE FOR THE CRITERION GROUPING FOR 4 TYPES IS .572650-1-1

SAMPLE SIZE = 151
 NUMBER OF VARIABLES = 2

ITERATION NUMBER =

REGRESSION OF 4 TYPES IN THIS SAMPLE = .56259152-1-2

CHARACTERISTICS OF THE WHOLE SAMPLE

1	2	3	4
5.3622	5.2573	5.7571	5.1943
STANDARD DEVIATIONS			
0.8265	0.4359	1.7651	1.7622
CORRELATIONS			
-0.1179	-0.1179	-0.8210	-0.8272
-0.1179	-0.1179	-0.8210	-0.8272
-0.1179	-0.1179	-0.8210	-0.8272

CHARACTERISTICS OF TYPE 1

THE PROPORTION OF THE POPULATION FROM THIS TYPE = 0.11

1	2	3	4
5.1148	5.2583	5.9165	5.1526
STANDARD DEVIATIONS			
1.4829	0.2323	1.4551	0.2362
CORRELATIONS			
-0.1179	0.1179	-0.8210	-0.8272
-0.1179	0.1179	-0.8210	-0.8272
-0.1179	0.1179	-0.8210	-0.8272

CHARACTERISTICS OF TYPE 2

THE PROPORTION OF THE POPULATION FROM THIS TYPE = 0.11

1	2	3	4
5.0116	5.4215	5.4825	5.2637
STANDARD DEVIATIONS			
0.3329	0.3791	0.1727	0.1154
CORRELATIONS			
0.1179	0.2425	0.2071	-0.2711
0.1179	0.2425	0.2071	-0.2711
0.1179	0.2425	0.2071	-0.2711

CHARACTERISTICS OF TYPE 3

THE PROPORTION OF THE POPULATION FROM THIS TYPE = 0.12

1	2	3	4
5.5321	5.2757	5.4627	5.2794
STANDARD DEVIATIONS			
0.3178	0.2657	0.2015	0.1563
CORRELATIONS			
0.1179	0.1179	-0.8210	-0.8272
0.1179	0.1179	-0.8210	-0.8272
0.1179	0.1179	-0.8210	-0.8272

101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120

THE PROPORTION OF THE POPULATION FROM THIS TYPE * 0.21

MEANS			
1	2	3	4
6.2664	2.2146	6.5567	1.6657
STANDARD DEVIATIONS			
1	2	3	4
0.3206	0.2178	0.2660	0.2492
CORRELATIONS			
1	2	3	4
0.1011	0.1977	-0.1531	-0.2652
-0.1818	-0.1732	0.1572	-0.2422
-0.1818	0.1646	0.0252	0.1111

PROPORTION	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50
PROPORTION	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50
PROPORTION	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50
PROPORTION	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50

TOTAL TYPE 4617 * 1.00000000

SAMPLE SIZE * 35
NUMBER OF VARIABLES * 2
NUMBER OF TYPES * 2

PROPORTION FROM TYPE 1

LIKELIHOOD OF 4 TYPES IN THIS SAMPLE * .572734612
CHARACTERISTICS OF THE WHOLE SAMPLE

MEANS			
1	2	3	4
5.6423	3.0572	2.7373	1.9493
STANDARD DEVIATIONS			
1	2	3	4
0.4621	0.4359	0.7655	0.7622
CORRELATIONS			
1	2	3	4
0.1011	-0.1176	0.8113	-0.1170
-0.1818	-0.1572	-0.4584	-0.2161
-0.1818	-0.2041	0.7629	0.2876

CHARACTERISTICS OF TYPE 1

THE PROPORTION OF THE POPULATION FROM THIS TYPE * 0.21

MEANS			
1	2	3	4
6.7822	3.0666	2.8128	2.3523
STANDARD DEVIATIONS			
1	2	3	4
0.5542	0.3112	0.4459	0.2200
CORRELATIONS			
1	2	3	4
0.1011	0.3645	0.8162	-0.2163
-0.1818	0.1542	0.4584	-0.2161
-0.1818	0.2452	-0.4116	0.1111

CHARACTERISTICS OF TYPE 2

THE PROPORTION OF THE POPULATION FROM THIS TYPE * 0.23

MEANS			
1	2	3	4
5.0141	2.4261	3.6821	6.2661
STANDARD DEVIATIONS			
1	2	3	4
1.2649	1.2753	1.3750	1.143
CORRELATIONS			
1	2	3	4
0.1011	0.2475	0.2977	0.2788
-0.1818	0.1572	0.1777	0.2161
0.1818	0.2041	0.2516	0.2161

CHARACTERISTICS OF TYPE 3

THE PROPORTION OF THE POPULATION FROM THIS TYPE * 0.21

MEANS			
1	2	3	4
5.4253	2.4976	6.1164	1.2501
STANDARD DEVIATIONS			
1	2	3	4
0.2636	0.2944	0.4216	0.1090
CORRELATIONS			
1	2	3	4
0.1011	0.1371	0.1517	0.2161
-0.1818	0.1572	0.1777	0.2161
0.1818	0.2041	0.2516	0.2161

CHARACTERISTICS OF TYPE 4

THE PROPORTION OF THE POPULATION FROM THIS TYPE * 0.25

MEANS			
1	2	3	4
6.2862	2.8369	6.8756	1.6273
STANDARD DEVIATIONS			
1	2	3	4
0.3644	0.4669	0.1111	0.1111

BIBLIOGRAFIA

1. Bartlett H. S. (1939), "A note on tests of significance in multivariate Analysis", Proc. Camb. Vol 8, pp 376-86.
2. Beale E.M.L. (1969), Cluster Analysis. London: Scientific Control Systems.
3. Cohen A. (1967), "Estimation in mixtures of two normal distributions", Technometrics, Vol. 9, No. 1, pp 15-28.
4. Day E. (1969) "Estimating the components of a mixture of normal distributions". Biometrika, Vol 56, pp 465-72.
5. Dempster et Al. (1976), "Maximum likelihood from incomplete data via the EM algorithm", Journal of the Royal Statistical Society, pp 3 -40.
6. Everitt B., (1974), Cluster Analysis, SSRC, Heinemann, London.
7. Fisher R. A. (1936), "Multiple measurements in taxonomic problems", Annals of Eugenics, Vol VII, pp 179-88.
8. Friedman y Rubin (1967), "On some invariant criteria for grouping data", JASA, Vol. 62, pp 1159-78.
9. Gnanadesikan R. (1977), Methods for Statistical data analysis of multivariate observations" John Wiley & Sons New York, pp 82-120.
10. Hernández A. y Harriot F.H.C. (1981), "The theory of distributions mixtures", Comunicaciones técnicas del IIMAS, No. 228, Serie Naranja.

11. Hernández A. (1979), "Problems in Cluster Analysis" Ph. D. Thesis, Oxford University.
12. Hasselblad, V. (1966), "Estimation of parameters for a mixture of normal distributions, *Technometrics*, Vol. 8, No. 3, pp 431-46.
13. Kale B.K. (1962), "On the solution of likelihood equations by iteration processes. The multiparametric case", *Biometrika*, Vol. 49, pp 479-86.
14. Marriot F.H.C. (1971), "Practical problems in a method of cluster analysis". *Biometrics*, Vol. 27, pp 501-14.
15. Marriot F.H.C. (1975), "Separating mixtures of normal distributions", *Biometrics*, 31, pp 767-69.
16. Silvey S.D. (1959), "Lagrangian Multiple tests", *Annals of Mathematical Statistics*, Vol. 30, pp 389-407.
17. Wilks S.S. (1938), "A large sample distribution of the likelihood ratio for testing composite hypothesis", *Annals of Mathematical Statistics*, Vol. 9, pp 60-62.
18. Wolfe J.H. (1967). "Normix; computational methods for estimating the parameters of multivariate normal mixtures of distributions". Research memorandum, SRM 68-2, U.S. Naval Personnel Research Activity, San Diego.
19. Wolfe J.H. (1970), "Pattern clustering by multivariate mixture analysis", U.S. Naval Personnel Research Laboratory, pp 329-50.