

**UTILIZACIÓN DE MÉTRICAS RIEMANNIANAS  
EN ANÁLISIS DE DATOS MULTIDIMENSIONALES  
Y SU APLICACIÓN A LA BIOLOGÍA**

**JOSE M<sup>a</sup> OLLER SALA**

**BARCELONA, 25 de NOVIEMBRE de 1982.**

donde  $c_1$  es una constante de integración. Nótese que  $c_1 > 0$  si y sólo si  $\mu_B > \mu_A$ ,  $c_1 < 0$  si y sólo si  $\mu_B < \mu_A$ . Para el caso  $\frac{d\mu}{ds} = 0$ ,  $\mu = \text{cte.}$  es una solución de la primera ecuación de (28) y puede considerarse incluida en (30), haciendo  $c_1 = 0$ . Entonces  $\mu = \mu_A = \mu_B$ .

Sustituyendo (30) en la segunda de (28), obtenemos:

$$\frac{d^2\sigma}{ds^2} - \frac{1}{\sigma} \left(\frac{d\sigma}{ds}\right)^2 + \frac{c_1^2}{2} \sigma^3 = 0 \quad (31)$$

Si efectuamos el cambio:

$$\frac{d\sigma}{ds} = \pm \sqrt{y} \quad (32)$$

donde el signo de  $\sqrt{y}$  depende del signo de  $\frac{d\sigma}{ds}$ , resulta:

$$\frac{d^2\sigma}{ds^2} = \frac{1}{\pm 2\sqrt{y}} \frac{dy}{d\sigma} \frac{d\sigma}{ds} = \frac{1}{2} \frac{dy}{d\sigma} \quad (33)$$

por tanto obtenemos:

$$\frac{dy}{d\sigma} - \frac{2}{\sigma} y + c_1^2 \sigma^3 = 0 \quad (34)$$

ecuación lineal cuya integral general viene dada por:

$$y = e^{\int 2/\sigma d\sigma} \left( -c_1^2 \int \sigma^3 e^{-\int 2/\sigma d\sigma} d\sigma + c_2 \right) = \sigma^2 \left( c_2 - \frac{c_1^2}{2} \sigma^2 \right) \quad (35)$$

Recordemos ahora que el sistema debe resolverse con la condición adicional que el vector tangente a la curva  $\left( \frac{d\mu}{ds}, \frac{d\sigma}{ds} \right)$  sea uni-

tario en un punto, (28), y por tanto, en todos los puntos de la curva (por transporte paralelo). De (28), (30) y (32) se sigue que:

$$c_2 = \frac{1}{2} \quad (36)$$

por tanto (35) se transforma, teniendo en cuenta (32), en:

$$\frac{d\sigma}{ds} = \pm \sigma \sqrt{\frac{1}{2} - \frac{c_1}{2} \sigma^2} \quad (37)$$

Vamos a resolver (37) distinguiendo varios casos.

Caso a  $c_1=0$ , ocurre como hemos visto cuando  $\mu_A=\mu_B$ , resulta:

$$\frac{d\sigma}{ds} = \pm \frac{\sigma}{\sqrt{2}} \quad (38)$$

puesto que  $\sigma > 0$  eligiremos la solución positiva si  $\sigma_B > \sigma_A$  y la negativa en caso contrario:

$$\int_{\sigma_A}^{\sigma_B} \pm \frac{d\sigma}{\sigma} = \int_0^{s_{AB}} \frac{ds}{\sqrt{2}} \quad (39)$$

por tanto:

$$s_{AB} = \sqrt{2} \left| \ln \frac{\sigma_B}{\sigma_A} \right| \quad (40)$$

y las ecuaciones de las geodésicas que pasan por el punto  $(\mu_A, \sigma_A)$  son:

$$\begin{aligned} \sigma &= \sigma_A e^{\pm s/\sqrt{2}} \\ \mu &= \mu_A \end{aligned} \quad (41)$$

eligiendo el signo en la primera, positivo si  $\sigma_B > \sigma_A$  y negativo si  $\sigma_B < \sigma_A$ . Es claro que siempre existe una geodésica que une  $(\mu_A, \sigma_A)$  con otro punto de la forma  $(\mu_B, \sigma_B)$ , con  $\mu_B = \mu_A$ .

Caso b Supongamos que  $\mu_A \neq \mu_B$  entonces  $c_1 \neq 0$  y supongamos además que  $\frac{d\sigma}{ds} \neq 0$  en todos los puntos de la curva comprendidos entre  $(\mu_A, \sigma_A)$  y  $(\mu_B, \sigma_B)$ . Entonces  $\frac{d\sigma}{ds} > 0$  si  $\sigma_B > \sigma_A$  y  $\frac{d\sigma}{ds} < 0$  si  $\sigma_B < \sigma_A$ , habiendo de elegir el signo en (37) de acuerdo con estas consideraciones. Resulta:

$$\left| \int_{\sigma_A}^{\sigma_B} \frac{d\sigma}{\sigma \sqrt{1 - c_1^2 \sigma^2}} \right| = \int_0^{s_{AB}} \frac{ds}{\sqrt{2}} \quad (42)$$

y por tanto:

$$s_{AB} = \sqrt{2} \left| \operatorname{argsech}(|c_1| \sigma_A) - \operatorname{argsech}(|c_1| \sigma_B) \right| \quad (43)$$

Podemos expresar  $\sigma$  en función de  $s$ , de forma que para  $s=0$  la curva pase por  $\sigma_A$ :

$$\sigma = \frac{1}{|c_1|} \operatorname{sech} \left( \pm \frac{s}{\sqrt{2}} + \operatorname{argsech}(|c_1| \sigma_A) \right) \quad (44)$$

eligiendo el signo negativo si queremos que  $\sigma > \sigma_A$  y positivo en caso contrario.

En cuanto a  $\mu$ , teniendo en cuenta (30), podemos escribir:

$$\frac{d\mu}{ds} = \frac{1}{c_1} \operatorname{sech}^2 \left( \pm \frac{s}{\sqrt{2}} + \operatorname{argsech}(|c_1| \sigma_A) \right) \quad (45)$$

y por tanto:

$$\int_{\mu_A}^{\mu_B} d\mu = \frac{1}{c_1} \int_0^{s_{AB}} \operatorname{sech}^2\left(\pm \frac{s}{\sqrt{2}} + \operatorname{argsech}(|c_1| \sigma_A)\right) ds \quad (46)$$

Teniendo en cuenta que:

$$\operatorname{tgh}(\operatorname{arg sech}(|c_1| \sigma)) = \sqrt{1 - c_1^2 \sigma^2} \quad (47)$$

y las consideraciones efectuadas acerca del signo en la expresión (44), al integrar (46) obtenemos:

$$\mu_B - \mu_A = \frac{\sqrt{2}}{c_1} \left| \sqrt{1 - c_1^2 \sigma_B^2} - \sqrt{1 - c_1^2 \sigma_A^2} \right| \quad (48)$$

Podemos obtener también  $\mu$  en función de  $s$ , de forma que para  $s=0$  la curva pase por  $\mu_A$ :

$$\mu = \mu_A + \frac{\sqrt{2}}{c_1} \left| \operatorname{tgh}\left(\mp \frac{s}{\sqrt{2}} + \operatorname{arg sech}(|c_1| \sigma_A)\right) - \sqrt{1 - c_1^2 \sigma_A^2} \right| \quad (49)$$

Nótese que el signo de  $c_1$  es positivo si  $\mu > \mu_A$  y negativo en caso contrario. Las ecuaciones (44) y (49) son las ecuaciones de las geodésicas, en el caso b, y la distancia entre dos puntos medida sobre la geodésica viene dada por (43).

A continuación intentemos encontrar  $c_1$  y a la vez caracterizar el caso b, es decir, dados dos puntos de la variedad paramétrica  $(\mu_A, \sigma_A)$ ,  $(\mu_B, \sigma_B)$  que condiciones deben cumplir para que exista una geodésica que los una, cumpliendo las hipótesis del caso b. Elevando al cuadrado (48) obtenemos:



$$c_1^2 \left( \frac{(\mu_B - \mu_A)^2}{2} + \sigma_A^2 + \sigma_B^2 \right) - 2 = -2 \sqrt{1 - c_1^2 \sigma_B^2} \sqrt{1 - c_1^2 \sigma_A^2} \quad (50)$$

para que exista solución, bajo las hipótesis consideradas, debe cumplirse:

$$c_1^2 \left( \frac{(\mu_B - \mu_A)^2}{2} + \sigma_A^2 + \sigma_B^2 \right) \leq 2 \quad (51)$$

Hallemos ahora  $c_1$ . Elevando al cuadrado (50) y reordenando términos, obtenemos:

$$c_1^2 \left( \frac{(\mu_B - \mu_A)^4}{4} + (\sigma_A^2 - \sigma_B^2)^2 + (\sigma_A^2 + \sigma_B^2)(\mu_B - \mu_A)^2 \right) = 2(\mu_B - \mu_A)^2 \quad (52)$$

y finalmente, como  $c_1 > 0$  si  $\mu_B - \mu_A > 0$ , resulta:

$$c_1 = \frac{\sqrt{2} (\mu_B - \mu_A)}{\sqrt{\frac{(\mu_B - \mu_A)^4}{4} + (\sigma_A^2 - \sigma_B^2)^2 + (\sigma_A^2 + \sigma_B^2)(\mu_B - \mu_A)^2}} \quad (53)$$

y por tanto la condición (51) puede reescribirse, al sustituir  $c_1$  por (53), como:

$$(\mu_B - \mu_A)^2 \left( \frac{(\mu_B - \mu_A)^4}{2} + \sigma_A^2 + \sigma_B^2 \right) \leq \frac{(\mu_B - \mu_A)^4}{4} + (\sigma_A^2 - \sigma_B^2)^2 + (\sigma_A^2 + \sigma_B^2)(\mu_B - \mu_A)^2 \quad (54)$$

equivalente a:

$$(\mu_B - \mu_A)^2 \leq 2 |\sigma_A^2 - \sigma_B^2| \quad (55)$$

La expresión (55) nos da la condición que deben verificar dos puntos de la variedad para que pueda existir una geodésica que los una cum

pliendo las hipótesis del caso b.

Caso c Consideremos ahora el caso  $\mu_A \neq \mu_B$  y  $\frac{d\sigma}{ds} > 0$  al inicio de la curva y  $\frac{d\sigma}{ds} < 0$  al final. Por razones de continuidad existirá un punto intermedio de la curva con  $\frac{d\sigma}{ds} = 0$ , en este punto  $\sigma$  toma el valor máximo. Nótese que no es posible que  $\frac{d\sigma}{ds} < 0$  al inicio de la curva y  $\frac{d\sigma}{ds} > 0$  al final, ya que en el punto donde se anularía la derivada habría un mínimo, para  $\sigma$ , hecho que estaría en contradicción con (31) que implica  $\frac{d^2\sigma}{ds^2} < 0$ .

Dividiremos la geodésica en dos trozos. Es claro que  $\sigma_{\max} = \frac{1}{|c_1|}$ , por tanto:

$$\int_{\sigma_A}^{1/|c_1|} \frac{d\sigma}{\sigma\sqrt{1-c_1^2\sigma^2}} - \int_{1/|c_1|}^{\sigma_B} \frac{d\sigma}{\sigma\sqrt{1-c_1^2\sigma^2}} = \int_0^{s_M} \frac{ds}{\sqrt{2}} + \int_{s_M}^{s_{AB}} \frac{ds}{\sqrt{2}} \quad (56)$$

integrando resulta:

$$s_{AB} = \sqrt{2} (\arg \operatorname{sech}(|c_1| \sigma_A) + \arg \operatorname{sech}(|c_1| \sigma_B)) \quad (57)$$

Podemos expresar  $\sigma$  en función de  $s$ , de forma que para  $s=0$  la curva pase por  $\sigma_A$ :

$$\sigma = \frac{1}{|c_1|} \operatorname{sech} \left( \left| \frac{s}{\sqrt{2}} - \arg \operatorname{sech}(|c_1| \sigma_A) \right| \right) \quad (58)$$

En cuanto a  $\mu$ , podemos escribir:

$$\int_{\mu_A}^{\mu_M} d\mu + \int_{\mu_M}^{\mu_B} d\mu = \frac{1}{c_1} \int_0^{s_M} \operatorname{sech}^2 \left( -\frac{s}{\sqrt{2}} + \operatorname{argsech}(|c_1| \sigma_A) \right) ds +$$

$$+ \frac{1}{c_1} \int_{s_M}^{s_{AB}} \operatorname{sech}^2 \left( \frac{s}{\sqrt{2}} - \operatorname{argsech} (|c_1| \sigma_A) \right) ds \quad (59)$$

teniendo en cuenta (47) y que para  $s_M$ ,  $\sigma$  alcanza  $\frac{1}{|c_1|}$ , integrando ob-  
tenemos:

$$\mu_B - \mu_A = \frac{\sqrt{2}}{c_1} \left( \sqrt{1 - c_1^2 \sigma_B^2} + \sqrt{1 - c_1^2 \sigma_A^2} \right) \quad (60)$$

Podemos también obtener  $\mu$  en función de  $s$ , de forma que para  $s=0$ ,  
la curva pase por  $\mu_A$ :

$$\mu = \mu_A + \frac{\sqrt{2}}{c_1} \left| \operatorname{tgh} \left( -\frac{s}{\sqrt{2}} + \operatorname{argsech} (|c_1| \sigma_A) \right) - \sqrt{1 - c_1^2 \sigma_A^2} \right| \quad (61)$$

La expresión (58) junto con (61) constituyen las ecuaciones de las  
geodésicas, en el caso c, y la distancia entre dos puntos medida sobre  
la geodésica viene dada por (57).

Intentemos determinar, a continuación, la condición que deben ve-  
rificar los puntos  $(\mu_A, \sigma_A)$  y  $(\mu_B, \sigma_B)$  para que se de el caso c. Elevando  
al cuadrado (60) obtenemos:

$$c_1^2 \left( \frac{(\mu_B - \mu_A)^2}{2} + \sigma_B^2 + \sigma_A^2 \right) - 2 = 2 \sqrt{1 - c_1^2 \sigma_B^2} \sqrt{1 - c_1^2 \sigma_A^2} \quad (62)$$

por tanto debe verificarse:

$$c_1^2 \left[ \frac{(\mu_B - \mu_A)^2}{2} + \sigma_B^2 + \sigma_A^2 \right] > 2 \quad (63)$$

Podemos hallar  $c_1$  a partir de (62), elevando al cuadrado y reordenan-  
do términos obtenemos finalmente el mismo valor que en el caso anterior,



en otras palabras podemos calcular  $c_1$  a partir de (53). Al sustituir en (63) y operar obtenemos finalmente:

$$(\mu_B - \mu_A)^2 > 2 |\sigma_A^2 - \sigma_B^2| \quad (64)$$

La expresión (64) nos da la condición que deben verificar dos puntos de la variedad para que pueda existir una geodésica que los una cumpliendo las hipótesis del caso c.

Observemos que dados dos puntos  $(\mu_A, \sigma_A)$  y  $(\mu_B, \sigma_B)$  siempre existirá una geodésica que los una. Bastará comprobar que:

$$|c_1| \sigma_A \leq 1 \quad |c_1| \sigma_B \leq 1 \quad (65)$$

será suficiente probarlo para el máximo de  $\sigma_A, \sigma_B$ . Sea, por ejemplo,  $\sigma_A$ , la primera expresión de (65) es equivalente a:

$$2(\mu_B - \mu_A)^2 \sigma_A^2 \leq \frac{(\mu_B - \mu_A)^4}{4} + (\sigma_A^2 - \sigma_B^2) + (\sigma_A^2 + \sigma_B^2) \cdot (\mu_B - \mu_A)^2 \quad (66)$$

sea  $\sigma_A^2 = \sigma_B^2 + \epsilon$ ,  $\epsilon \geq 0$ , entonces, si llamamos  $\Delta = \mu_B - \mu_A$ , resulta:

$$2\Delta^2 \sigma_A^2 \leq \frac{\Delta^4}{4} + \epsilon^2 + 2\sigma_A^2 \Delta^2 - \epsilon \Delta^2 \quad (67)$$

equivalente a:

$$0 \leq \frac{\Delta^4}{4} + \epsilon^2 - \epsilon \Delta^2 \quad (68)$$

La igualdad en la expresión anterior sólo se verifica cuando  $\epsilon = \frac{\Delta^2}{2}$  lo que implica que la expresión (68) no puede tomar valores negativos,

quedando por tanto probada (65).

Notemos que las geodésicas verifican, en el caso a, la propiedad de ser rectas perpendiculares al eje de las  $\mu$ . En los demás casos, de (47), (59) y (61), se deduce:

$$(\mu - \mu_M)^2 = \frac{2}{c_1} (1 - c_1^2 \sigma^2) \quad (69)$$

y como  $\sigma_M = \frac{1}{c_1}$  resulta:

$$\frac{(\mu - \mu_M)^2}{2\sigma_M^2} + \frac{\sigma^2}{\sigma_M^2} = 1 \quad (70)$$

es decir, las geodésicas son arcos de elipse con centro en  $(\mu_M, 0)$ . Es fácil comprobar que:

$$\mu_M = \frac{\mu_A + \mu_B}{2} - \frac{\sigma_B^2 - \sigma_A^2}{\mu_B - \mu_A} \quad (71)$$

Notemos también una transformación admisible de los parámetros que deja constante al tensor métrico excepto en una componente. La transformación viene definida por:

$$\begin{aligned} (\mu, \sigma) &\longrightarrow (\mu, v) \\ 2\sigma^2 &= e^{-\sqrt{2}v} \end{aligned} \quad (72)$$

Bajo las nuevas coordenadas el tensor métrico viene definido por:

$$\bar{g}_{11} = 2e^{\sqrt{2}v} \quad \bar{g}_{12} = \bar{g}_{21} = 0 \quad \bar{g}_{22} = 1 \quad (73)$$

Notemos además que como las geodésicas están definidas para todo valor de su parámetro, estamos en condiciones de afirmar que la variedad Riemanniana E definida en (17) y (20) es completa, por lo que podemos afirmar que la distancia entre dos puntos de la variedad vendrá dada por integración del elemento de línea a lo largo de una geodésica. Resumiendo los resultados anteriores, la distancia entre dos puntos  $(\mu_A, \sigma_A)$  y  $(\mu_B, \sigma_B)$ ,  $d_{AB}$ , viene dada por:

$$\begin{aligned} \mu_A = \mu_B & \Rightarrow d_{AB} = \sqrt{2} \left| \ln \left( \frac{\sigma_B}{\sigma_A} \right) \right| \\ \left. \begin{aligned} \mu_A \neq \mu_B \\ (\mu_B - \mu_A)^2 \leq 2|\sigma_A^2 - \sigma_B^2| \end{aligned} \right\} & \Rightarrow d_{AB} = \sqrt{2} |\operatorname{argsech}(|c_1|\sigma_A) - \operatorname{argsech}(|c_1|\sigma_B)| \\ \left. \begin{aligned} \mu_A \neq \mu_B \\ (\mu_B - \mu_A)^2 > 2|\sigma_A^2 - \sigma_B^2| \end{aligned} \right\} & \Rightarrow d_{AB} = \sqrt{2} (\operatorname{argsech}(|c_1|\sigma_A) + \operatorname{argsech}(|c_1|\sigma_B)) \end{aligned} \quad (74)$$

donde  $|c_1|$  viene dado por (53).

Una vez resuelto el caso normal univariante, por el resultado 3.1.1 queda resuelto el caso de la distribución log-normal, y en general de el caso de todas aquellas variables aleatorias que puedan transformarse mediante una transformación admisible, en una variable aleatoria con distribución normal.

### 5.3. DISTRIBUCION NORMAL MULTIVARIANTE

Sea  $X = (X_1, \dots, X_n)$  un vector aleatorio cuya función de densidad conjunta viene definida por:

$$f(X, \mu, \Sigma) = (2\pi)^{-\frac{1}{2}n} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} (x-\mu)^t \Sigma^{-1} (x-\mu)\right) \quad (75)$$

donde  $\mu$  es el vector formado con las esperanzas de las  $X_i$  y  $\Sigma$  es la matriz de varianzas y covarianzas entre las  $X_i$ . El espacio paramétrico vendrá definido por:

$$E = \left\{ (\mu_1, \dots, \mu_n, \sigma_{11}, \dots, \sigma_{nn}) \in \mathbb{R}^{n + \frac{n(n+1)}{2}} \mid \Sigma = (\sigma_{ij}) \text{ es simétrica y definida positiva} \right\} \quad (76)$$

Se comprueba fácilmente que  $E$  es una variedad  $n + \frac{n(n+1)}{2}$  dimensional arco-conexa.

Como el logaritmo de la función de densidad viene dado por:

$$\ln f(X, \mu, \Sigma) = -\frac{n}{2} \ln(2\pi) - \frac{1}{2} \ln|\Sigma| - \frac{1}{2} (x-\mu)^t \Sigma^{-1} (x-\mu) \quad (77)$$

al derivar respecto algún parámetro  $\alpha$ , ( $\mu_i$  ó  $\sigma_{ij}$ ) resulta:

$$\frac{\partial \ln f}{\partial \alpha} = -\frac{1}{2} \left[ \frac{1}{|\Sigma|} \frac{\partial |\Sigma|}{\partial \alpha} + \sum_{i=1}^n \sum_{j=1}^n \left( \frac{\partial \sigma^{ij}}{\partial \alpha} (x_i - \mu_i) (x_j - \mu_j) - 2\sigma^{ij} \frac{\partial \mu_i}{\partial \alpha} (x_j - \mu_j) \right) \right] \quad (78)$$

siendo los  $\sigma^{ij}$  los elementos de  $\Sigma^{-1}$ , matriz también simétrica. Volviendo a derivar, ahora respecto a  $\beta$ , y cambiando el signo, resultará:

$$\begin{aligned} -\frac{\partial^2 \ln f}{\partial \beta \partial \alpha} = \frac{1}{2} \left[ -\frac{1}{|\Sigma|^2} \frac{\partial |\Sigma|}{\partial \beta} \frac{\partial |\Sigma|}{\partial \alpha} + \frac{1}{|\Sigma|} \frac{\partial^2 |\Sigma|}{\partial \beta \partial \alpha} + \sum_{i=1}^n \sum_{j=1}^n \left( \frac{\partial^2 \sigma^{ij}}{\partial \beta \partial \alpha} (x_i - \mu_i) (x_j - \mu_j) - \right. \right. \\ \left. \left. - 2 \frac{\partial \sigma^{ij}}{\partial \alpha} \frac{\partial \mu_i}{\partial \beta} (x_j - \mu_j) - 2 \frac{\partial \sigma^{ij}}{\partial \beta} \frac{\partial \mu_i}{\partial \alpha} (x_j - \mu_j) + 2\sigma^{ij} \frac{\partial \mu_i}{\partial \alpha} \frac{\partial \mu_j}{\partial \beta} \right) \right] \quad (79) \end{aligned}$$

Hallando la esperanza y teniendo en cuenta que  $E(x_i - \mu_i) = 0$ , resulta:

$$-E\left(\frac{\partial^2 \ln f}{\partial \beta \partial \alpha}\right) = \frac{1}{2} \left[ -\frac{1}{|\Sigma|^2} \frac{\partial |\Sigma|}{\partial \beta} \frac{\partial |\Sigma|}{\partial \alpha} + \frac{1}{|\Sigma|} \frac{\partial^2 |\Sigma|}{\partial \beta \partial \alpha} + \sum_{i=1}^n \sum_{j=1}^n \left( \frac{\partial^2 \sigma^{ij}}{\partial \beta \partial \alpha} \sigma_{ij} + 2\sigma^{ij} \frac{\partial \mu_i}{\partial \alpha} \frac{\partial \mu_j}{\partial \beta} \right) \right] \quad (80)$$

Consideremos los siguientes casos:

Caso a     $\alpha = \mu_k$              $\beta = \mu_h$              $k, h = 1, \dots, n$

$$-E\left(\frac{\partial^2 \ln f}{\partial \mu_k \partial \mu_h}\right) = \sigma^{kh} \quad (81)$$

Caso b     $\alpha = \mu_k$              $\beta = \sigma_{hm}$     ( $\delta \beta = \mu_k$      $\alpha = \sigma_{hm}$ )     $k, h, m = 1, \dots, n$

$$-E\left(\frac{\partial^2 \ln f}{\partial \mu_k \partial \sigma_{hm}}\right) = 0 \quad (82)$$

Caso c     $\alpha = \sigma_{kh}$              $\beta = \sigma_{pq}$      $k, h, p, q = 1, \dots, n$

$$-E\left(\frac{\partial^2 \ln f}{\partial \sigma_{kh} \partial \sigma_{pq}}\right) = \frac{1}{2} \left[ -\frac{1}{|\Sigma|^2} \frac{\partial |\Sigma|}{\partial \sigma_{kh}} \frac{\partial |\Sigma|}{\partial \sigma_{pq}} + \frac{1}{|\Sigma|} \frac{\partial^2 |\Sigma|}{\partial \sigma_{kh} \partial \sigma_{pq}} + \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2 \sigma^{ij}}{\partial \sigma_{kh} \partial \sigma_{pq}} \sigma_{ij} \right] \quad (83)$$

Observemos que:

$$\sum_{i=1}^n \sum_{j=1}^n \sigma^{ij} \sigma_{ij} = \sum_{i=1}^n \delta_{ii} = n \quad (84)$$

luego su derivada segunda, respecto de  $\sigma_{kh}$  ó  $\sigma_{pq}$ , se anulará:

$$0 = \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2}{\partial \sigma_{kh} \partial \sigma_{pq}} (\sigma^{ij} \sigma_{ij}) = \sum_{i=1}^n \sum_{j=1}^n \left( \frac{\partial^2 \sigma^{ij}}{\partial \sigma_{kh} \partial \sigma_{pq}} \sigma_{ij} + \frac{\partial \sigma^{ij}}{\partial \sigma_{pq}} \frac{\partial \sigma_{ij}}{\partial \sigma_{kh}} + \frac{\partial \sigma_{ij}}{\partial \sigma_{kh}} \cdot \frac{\partial \sigma^{ij}}{\partial \sigma_{pq}} \right) \quad (85)$$

y por tanto resulta:

$$\sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2 \sigma^{ij}}{\partial \sigma_{kh} \partial \sigma_{pq}} \sigma_{ij} = - \frac{\partial \sigma^{kh}}{\partial \sigma_{pq}} (2 - \delta_{kh}) - \frac{\partial \sigma^{pq}}{\partial \sigma_{kh}} (2 - \delta_{pq}) \quad (86)$$

y teniendo en cuenta que:

$$\sigma^{ij} = \frac{\Sigma^{ji}}{|\Sigma|} \quad (87)$$

siendo  $\Sigma^{ji}$  el adjunto del elemento  $\sigma_{ji}$ . Por tanto podemos escribir:

$$\frac{\partial \sigma^{ij}}{\partial \sigma_{pq}} = - \frac{1}{|\Sigma|^2} \frac{\partial |\Sigma|}{\partial \sigma_{pq}} \Sigma^{ji} + \frac{1}{|\Sigma|} \frac{\partial \Sigma^{ji}}{\partial \sigma_{pq}} \quad (88)$$

y debido a que:

$$\frac{\partial |\Sigma|}{\partial \sigma_{pq}} = \Sigma^{qp} (2 - \delta_{pq}) \quad \frac{\partial^2 |\Sigma|}{\partial \sigma_{kh} \partial \sigma_{pq}} = \frac{\partial \Sigma^{qp}}{\partial \sigma_{kh}} (2 - \delta_{pq}) \quad (89)$$

resulta:

$$-E \left( \frac{\partial^2 \ln f}{\partial \sigma_{kh} \partial \sigma_{pq}} \right) = \frac{1}{2 |\Sigma|^2} \frac{\partial |\Sigma|}{\partial \sigma_{kh}} \frac{\partial |\Sigma|}{\partial \sigma_{pq}} - \frac{1}{2 |\Sigma|} \frac{\partial^2 |\Sigma|}{\partial \sigma_{kh} \partial \sigma_{pq}} = - \frac{1}{2} \frac{\partial^2 \ln |\Sigma|}{\partial \sigma_{kh} \partial \sigma_{pq}} \quad (90)$$

por tanto el tensor métrico asociado a la variedad E, definida en (76)

matricialmente podrá expresarse como:

$$G = \begin{pmatrix} \Sigma^{-1} & 0 \\ 0 & A \end{pmatrix} \quad (91)$$

siendo  $\Sigma^{-1}$  la matriz inversa de la matriz de varianzas y covarianzas,

y A una matriz cuadrada de orden  $n(n+1)/2$  y cuyos elementos son de la

forma:

$$a_{\mu\nu} = -\frac{1}{2} \frac{\partial^2 \ln|\Sigma|}{\partial\sigma_{kh} \partial\sigma_{pq}} \quad \begin{array}{l} k \leq h = 1, \dots, n \\ p \leq q = 1, \dots, n \end{array} \quad (92)$$

los índices  $\mu$  y  $\nu$  dependen de  $k, h, p$  y  $q$ , según ordenamos la parte triangular superior (o inferior) de la matriz de varianzas y covarianzas, al considerar a sus elementos como las coordenadas  $n+1, \dots, n + \frac{n(n+1)}{2}$  de la variedad  $E$ .

Hasta el momento no se ha logrado invertir explícitamente la matriz  $A$ , por lo que no se ha procedido al cálculo de los símbolos de Christoffel de segunda especie, necesarios para obtener las ecuaciones de las geodésicas de la forma habitual. Para obtener una expresión de las ecuaciones diferenciales de las geodésicas puede plantearse el problema variacional asociado al cálculo de las mismas. El cuadrado del elemento de línea será:

$$ds^2 = \sum_{p,q} \sigma^{pq} d\mu_p d\mu_q - \frac{1}{2} \sum_{\alpha < \beta} \sum_{r < s} \frac{\partial^2 \ln|\Sigma|}{\partial\sigma_{\alpha\beta} \partial\sigma_{rs}} d\sigma_{\alpha\beta} d\sigma_{rs} \quad (93)$$

por tanto habrá que resolver:

$$\delta \int_0^s \sqrt{\sum_{p,q} \sigma^{pq} \dot{\mu}_p \dot{\mu}_q - \frac{1}{2} \sum_{\alpha < \beta} \sum_{r < s} \frac{\partial^2 \ln|\Sigma|}{\partial\sigma_{\alpha\beta} \partial\sigma_{rs}} \dot{\sigma}_{\alpha\beta} \dot{\sigma}_{rs}} ds = 0 \quad (94)$$

donde  $\dot{\mu}_p$  y  $\dot{\sigma}_{rs}$  son las derivadas de  $\mu_p$  y  $\sigma_{rs}$  respecto  $s$ .

Si llamamos  $G$  a la expresión interior a la raíz, planteando las ecuaciones de Euler, respecto  $\mu_\gamma$ , resulta:

$$-\frac{d}{ds} \left( \frac{1}{2\sqrt{G}} \frac{\partial G}{\partial \dot{\mu}_\gamma} \right) = 0 \quad (95)$$



y si  $s$  es la longitud medida sobre el arco, como  $\sqrt{G} = 1$ , tendremos, debido a que  $\sigma^{p\gamma} = \sigma^{\gamma p}$ ,

$$\sum_{p=1}^n \sigma^{p\gamma} \dot{\mu}_\gamma = A_\gamma \quad (96)$$

siendo los  $A_\gamma$  constantes de integración. Equivalentemente podremos escribir:

$$\dot{\mu}_q = \sum_{\gamma=1}^n \sigma_{q\gamma} A_\gamma \quad (97)$$

Al plantear las ecuaciones de Euler respecto las coordenadas  $\sigma_{\tau\nu}$ , con  $\tau \leq \nu$ , resulta:

$$\frac{\partial \sqrt{G}}{\partial \sigma_{\tau\nu}} - \frac{d}{ds} \left( \frac{\partial \sqrt{G}}{\partial \dot{\sigma}_{\tau\nu}} \right) = \frac{1}{2\sqrt{G}} \frac{\partial G}{\partial \sigma_{\tau\nu}} - \frac{d}{ds} \left( \frac{1}{2\sqrt{G}} \frac{\partial G}{\partial \dot{\sigma}_{\tau\nu}} \right) = 0 \quad (98)$$

pero como  $\sqrt{G} = 1$ , podemos escribir:

$$\frac{\partial G}{\partial \sigma_{\tau\nu}} - \frac{d}{ds} \left( \frac{\partial G}{\partial \dot{\sigma}_{\tau\nu}} \right) = 0 \quad (99)$$

Como resulta que:

$$\frac{\partial G}{\partial \sigma_{\tau\nu}} = \sum_{p,q} \frac{\partial \sigma^{pq}}{\partial \sigma_{\tau\nu}} \dot{\mu}_p \dot{\mu}_q - \frac{1}{2} \sum_{\alpha < \beta} \sum_{r < s} \frac{\partial^3 \ln |\Sigma|}{\partial \sigma_{\tau\nu} \partial \sigma_{\alpha\beta} \partial \sigma_{rs}} \dot{\sigma}_{\alpha\beta} \dot{\sigma}_{rs} \quad (100)$$

y además:

$$\frac{\partial G}{\partial \dot{\sigma}_{\tau\nu}} = - \sum_{\alpha < \beta} \frac{\partial^2 \ln |\Sigma|}{\partial \sigma_{\alpha\beta} \partial \sigma_{\tau\nu}} \dot{\sigma}_{\alpha\beta} \quad (101)$$

y

$$- \frac{d}{ds} \left( \frac{\partial G}{\partial \dot{\sigma}_{\tau\nu}} \right) = \sum_{\alpha < \beta} \left( \sum_{r < s} \frac{\partial^3 \ln |\Sigma|}{\partial \sigma_{rs} \partial \sigma_{\alpha\beta} \partial \sigma_{\tau\nu}} \dot{\sigma}_{\alpha\beta} \dot{\sigma}_{rs} \right) + \frac{\partial^2 \ln |\Sigma|}{\partial \sigma_{\alpha\beta} \partial \sigma_{\tau\nu}} \ddot{\sigma}_{\alpha\beta} \quad (102)$$

podemos escribir:

$$\sum_{p,q} \frac{\partial \sigma^{pq}}{\partial \sigma_{\tau\nu}} \dot{\mu}_p \dot{\mu}_q + \frac{1}{2} \sum_{\alpha < \beta} \sum_{r < s} \frac{\partial^3 \ln |\Sigma|}{\partial \sigma_{rs} \partial \sigma_{\alpha\beta} \partial \sigma_{\tau\nu}} \dot{\sigma}_{\alpha\beta} \dot{\sigma}_{rs} + \sum_{\alpha < \beta} \frac{\partial^2 \ln |\Sigma|}{\partial \sigma_{\alpha\beta} \partial \sigma_{\tau\nu}} \ddot{\sigma}_{\alpha\beta} = 0$$

$$\tau \leq \nu = 1, \dots, n \quad (103)$$

Teniendo en cuenta que al cumplirse:

$$\sum_{q=1}^n \sigma^{pq} \sigma_{qr} = \delta_r^p \quad (104)$$

se verifica que:

$$\sum_{q=1}^n \frac{\partial \sigma^{pq}}{\partial \sigma_{\tau\nu}} \sigma_{qr} = - \sum_{q=1}^n \sigma^{pq} \frac{\partial \sigma_{qr}}{\partial \sigma_{\tau\nu}} \quad (105)$$

si sustituimos (97) en (103) y simplificamos, resulta finalmente:

$$-(2-\delta_{\tau\nu}) A_{\tau} A_{\nu} + \frac{1}{2} \sum_{r < s} \sum_{\alpha < \beta} \frac{\partial^3 \ln |\Sigma|}{\partial \sigma_{rs} \partial \sigma_{\alpha\beta} \partial \sigma_{\tau\nu}} \dot{\sigma}_{\alpha\beta} \dot{\sigma}_{rs} + \sum_{\alpha < \beta} \frac{\partial^2 \ln |\Sigma|}{\partial \sigma_{\alpha\beta} \partial \sigma_{\tau\nu}} \ddot{\sigma}_{\alpha\beta} = 0$$

$$\tau \leq \nu = 1, \dots, n \quad (106)$$

El sistema de ecuaciones diferenciales de las geodésicas viene dado pues por (97) y (106), y hay que resolverlo con la condición suplementaria:

$$\sum_{p,q} \sigma^{pq} \dot{\mu}_p \dot{\mu}_q - \frac{1}{2} \sum_{\alpha < \beta} \sum_{r < s} \frac{\partial^2 \ln |\Sigma|}{\partial \sigma_{\alpha\beta} \partial \sigma_{rs}} \dot{\sigma}_{\alpha\beta} \dot{\sigma}_{rs} = 1 \quad (107)$$

para que la norma del vector tangente, respecto al parámetro natural, sea unitaria. La expresión (107), teniendo en cuenta (97), puede ser escrita como:

$$\sum_{pq} \sigma_{pq} A_p A_q - \frac{1}{2} \sum_{\alpha < \beta} \sum_{r < s} \frac{\partial^2 \ln |\Sigma|}{\partial \sigma_{\alpha\beta} \partial \sigma_{rs}} \ddot{\sigma}_{\alpha\beta} \dot{\sigma}_{rs} = 1 \quad (108)$$

Es posible expresar de otra forma a (106), teniendo en cuenta que:

$$\frac{d^2}{ds^2} \left( \frac{\partial \ln |\Sigma|}{\partial \sigma_{\tau\nu}} \right) = \sum_{\alpha < \beta} \left( \sum_{r < s} \frac{\partial^3 \ln |\Sigma|}{\partial \sigma_{rs} \partial \sigma_{\alpha\beta} \partial \sigma_{\tau\nu}} \dot{\sigma}_{\alpha\beta} \dot{\sigma}_{rs} \right) + \frac{\partial^2 \ln |\Sigma|}{\partial \sigma_{\alpha\beta} \partial \sigma_{\tau\nu}} \ddot{\sigma}_{\alpha\beta} \quad (109)$$

podemos escribir:

$$-(2-\delta_{\tau\nu}) A_{\tau} A_{\nu} + \frac{1}{2} \frac{d^2}{ds^2} \left( \frac{\partial \ln |\Sigma|}{\partial \sigma_{\tau\nu}} \right) + \frac{1}{2} \sum_{\alpha < \beta} \frac{\partial^2 \ln |\Sigma|}{\partial \sigma_{\alpha\beta} \partial \sigma_{\tau\nu}} \ddot{\sigma}_{\alpha\beta} = 0 \quad (110)$$

$\tau < \nu = 1, \dots, n$

y como resulta:

$$\frac{\partial \ln |\Sigma|}{\partial \sigma_{\tau\nu}} = (2-\delta_{\tau\nu}) \sigma^{\tau\nu} \frac{\partial^2 \ln |\Sigma|}{\partial \sigma_{\alpha\beta} \partial \sigma_{\tau\nu}} = (2-\delta_{\tau\nu}) \frac{\partial \sigma^{\tau\nu}}{\partial \sigma_{\alpha\beta}} \quad (111)$$

obtenemos:

$$-2A_{\tau} A_{\nu} + \ddot{\sigma}^{\tau\nu} + \sum_{\alpha < \beta} \frac{\partial \sigma^{\tau\nu}}{\partial \sigma_{\alpha\beta}} \ddot{\sigma}_{\alpha\beta} = 0 \quad (112)$$

$\tau < \nu = 1, \dots, n$

multiplicando por  $\sigma_{\nu\lambda}$ ,  $\lambda \geq \tau$ , y sumando respecto a  $\nu$ , resulta:

$$-2 \sum_{\nu=1}^n A_{\tau} A_{\nu} \sigma_{\nu\lambda} + \sum_{\nu=1}^n \ddot{\sigma}^{\tau\nu} \sigma_{\nu\lambda} + \sum_{\alpha < \beta} \left( \sum_{\nu=1}^n \frac{\partial \sigma^{\tau\nu}}{\partial \sigma_{\alpha\beta}} \sigma_{\nu\lambda} \right) \ddot{\sigma}_{\alpha\beta} = 0 \quad (113)$$

$\tau < \lambda = 1, \dots, n$

pero como se cumple (105) se obtiene:

$$-2 \sum_{\nu=1}^n A_{\tau} A_{\nu} \sigma_{\nu\lambda} + \sum_{\nu=1}^n \ddot{\sigma}^{\tau\nu} \sigma_{\nu\lambda} - \sum_{\nu=1}^n \sigma^{\tau\nu} \ddot{\sigma}_{\nu\lambda} = 0 \quad (114)$$

$\tau < \lambda = 1, \dots, n$

Además, teniendo en cuenta (104), se cumple:

$$\sum_{\nu=1}^n A_{\tau}^{\tau} A_{\nu}^{\nu} \sigma_{\nu\lambda} + \sum_{\nu=1}^n \dot{\sigma}^{\tau\nu} \dot{\sigma}_{\nu\lambda} + \sum_{\nu=1}^n \sigma^{\tau\nu} \ddot{\sigma}_{\nu\lambda} = 0 \quad (115)$$

$$\tau \leq \lambda = 1, \dots, n$$

o equivalentemente:

$$\sum_{\nu=1}^n A_{\tau}^{\tau} A_{\nu}^{\nu} \sigma_{\nu\lambda} + \frac{d}{ds} \left( \sum_{\nu=1}^n \sigma^{\tau\nu} \dot{\sigma}_{\nu\lambda} \right) = 0 \quad (116)$$

$$\tau \leq \lambda = 1, \dots, n$$

que puede a su vez ser expresada matricialmente, junto con (97):

$$\dot{\mu} = \Sigma A$$

$$AA^t \Sigma + \frac{d}{ds} (\Sigma^{-1} \dot{\Sigma}) = 0 \quad (117)$$

sistema de ecuaciones diferenciales que definen las geodésicas para el caso de la distribución normal multivariante. Hasta el momento, dicho sistema no ha sido integrado por lo que no disponemos de una expresión que permita encontrar la distancia entre dos puntos medida sobre una geodésica. Se ha obtenido sin embargo, una acotación de la misma.

Sean  $(\mu_A, \Sigma_A)$  y  $(\mu_B, \Sigma_B)$  dos puntos del espacio paramétrico. Es bien sabido que siempre existe una transformación lineal tal que reduce simultáneamente a  $\Sigma_A$  y  $\Sigma_B$  a la forma diagonal, Mardia (1979):

$$\begin{aligned} \bar{\mu}_A &= T^t \mu_A & \bar{\Sigma}_A &= T^t \Sigma_A T = I \\ \bar{\mu}_B &= T^t \mu_B & \bar{\Sigma}_B &= T^t \Sigma_B T = D \end{aligned} \quad (118)$$

donde la matriz  $T$  es una matriz cuyas columnas son los vectores propios de  $\Sigma_B$  respecto  $\Sigma_A$  y:

$$D = \text{diag} (\lambda_1, \dots, \lambda_n) \quad (119)$$

siendo los  $\lambda_i$  los valores propios.

Intentemos ahora resolver el problema con la restricción suplementaria de que las variables aleatorias sean independientes en todo el trazo de la curva que une ambos puntos. En otras palabras, que la matriz de varianzas y covarianzas sea diagonal.

El caso de  $n$  variables aleatorias normales independientes ya está resuelto, por haber solucionado el caso de una distribución normal univariante y por el resultado 3.3.5. Por tanto una cota superior de la distancia entre  $(\mu_A, \Sigma_A)$  y  $(\mu_B, \Sigma_B)$  calculando las distancias entre cada una de las funciones de densidad marginales (normales univariantes), y sacando la raíz cuadrada de la suma de sus cuadrados. La distancia  $d_j$  entre la  $j$ -ésima marginal de  $A$  y de  $B$  es la distancia calculable a partir de la expresión (74) entre los puntos  $(\bar{\mu}_j, 1)_A$  y  $(\bar{\mu}_j, \lambda_j)_B$ , y por tanto una cota superior de la distancia vendrá dada por:

$$d^* = \sqrt{d_1^2 + \dots + d_n^2} \quad (120)$$

Nótese que en el caso  $\mu_A = \mu_B$  (y por tanto  $\bar{\mu}_A = \bar{\mu}_B$ ), si existe una solución del sistema (117) con  $\hat{\sigma}_{\alpha\beta} = 0$  si  $\alpha \neq \beta$ . En efecto, en este caso basta coger  $A=0$ , y resulta:

$$\frac{d^2 \ln \sigma_{ii}}{ds^2} = 0 \quad (121)$$

por tanto:

$$\ln \sigma_{ii} = A_i s + B_i \quad (122)$$

y al resolverlo con la condición que para  $s=0$   $\sigma_{ii}=1$  y para un cierto  $d$ ,  $\sigma_{ii} = \lambda_i$ , resulta

$$\sum_{i=1}^n \ln^2 \sigma_{ii} = \left( \sum_{i=1}^n A_i^2 \right) s^2 \quad (123)$$

y al ser por (107):

$$\sum_{i=1}^n A_i^2 = 2 \quad (124)$$

resulta:

$$d = \sqrt{\frac{1}{2} \sum_{i=1}^n \ln^2 \lambda_i} \quad (125)$$

por tanto, bajo la condición  $\mu_A = \mu_B$ , la distancia entre ambos puntos no puede ser mayor que  $d$  definida en (125).

Algunas aproximaciones a la distancia buscada pueden ser encontradas en el siguiente capítulo.

## 6. REPRESENTACION EN DIMENSION REDUCIDA.

Resumen:

En este capítulo se examina el problema de representar  $k$  puntos de una variedad riemanniana, en un espacio euclideo de dimensión  $q$ , generalmente  $q=2$ , de forma que la distancia euclidea entre ellos se asemeje a la distancia riemanniana original.

Vamos a efectuar tal estudio considerando varias situaciones posibles, atendiendo en primer lugar a la naturaleza de la geometría riemanniana definida en el espacio paramétrico y en segundo lugar al hecho de disponer, o no, de una expresión analítica de la distancia entre dos puntos de la variedad paramétrica, así como a criterios de tipo práctico.

Sumario:

6.1. ESPACIO PARAMETRICO EUCLIDEO.

6.2. ESPACIO PARAMETRICO NO EUCLIDEO.

6.2.1. Caso de disponer de una expresión analítica de la distancia.

6.2.2. Caso de no disponer de una expresión analítica de la distancia.

6.2.2.a. Utilización de técnicas de cálculo numérico.

6.2.2.b. Euclidización local.

6.2.2.c. Utilización del Test de la Razón de verosimilitud.



### 6.1. ESPACIO PARAMETRICO EUCLIDEO

Corresponde al caso más sencillo. En estas condiciones siempre existe un sistema de coordenadas bajo el cual el tensor métrico es un tensor constante, en particular la identidad. Esto se logra resolviendo el sistema (61) del capítulo 2, eligiendo las constantes de integración de forma adecuada para que el tensor métrico, en el nuevo sistema de coordenadas, sea la identidad. Ejemplo de ello puede verse en el capítulo 4, (91), (137) y (265) y en el capítulo 5, (8).

Si tenemos  $k$  puntos de coordenadas  $P_i: (x_{i1}, \dots, x_{in})$   $i=1, \dots, k$ , el primer paso será efectuar la transformación de coordenadas antes indicada, obteniendo las nuevas coordenadas de los puntos,  $P_i: (y_{i1}, \dots, y_{in})$   $i=1, \dots, k$ . El problema que se plantea ahora es el siguiente: dados  $k$  puntos de un espacio afín euclídeo, cuyas coordenadas están referidas a una base ortonormal, hallar una variedad lineal, de dimensión  $q$ , tal que la suma de los cuadrados de las distancias de los puntos (vectores) a dicha variedad sea mínima. Dicho en otras palabras, hallar el hiperplano,  $q$  dimensional, que "mejor se ajuste" a la nube de puntos  $P_1, \dots, P_k$ . Este problema ha sido ampliamente estudiado dentro del Análisis Multivariante, Mardia (1979), y se resuelve con la técnica conocida como Análisis de Componentes Principales. Si con las nuevas coordenadas de los  $P_i$  formamos una matriz  $Y$  de  $k$  filas y  $n$  columnas:

$$Y = (y_{ij})_{k \times n} \quad (1)$$

y definimos el vector columna, de  $n$  filas  $Z$ :

$$z_j = \frac{1}{k} \sum_{h=1}^k y_{hj} \quad (2)$$

entonces es bien sabido que el problema se resuelve diagonalizando la matriz:

$$A = Y^t Y - Z Z^t \quad (3)$$

respecto la matriz identidad I:

$$A v_i = \lambda_i I v_i = \lambda_i v_i \quad (4)$$

La variedad lineal buscada será de la forma:

$$V_q = \{x/ x = \alpha_1 v_1 + \dots + \alpha_q v_q + z \quad \alpha_i \in \mathbb{R}\} \quad (5)$$

siendo  $z$  el vector definido en (2), y  $v_1, \dots, v_q$   $q$  vectores propios, ortonormales, asociados a los valores propios  $\lambda_1, \dots, \lambda_q$  con tal de que éstos verifiquen:

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_q \quad (6)$$

Nótese que por ser  $A$  simétrica es siempre diagonalizable, por tanto la expresión (5) tendrá sentido para cualquier  $q$ ,  $1 \leq q \leq n$ . En cuanto a las coordenadas de los puntos proyectados en la variedad lineal, pueden obtenerse a través de la ecuación matricial:

$$U = (Y - E Z^t) V \quad (7)$$

donde  $E$  es el vector  $(k \times 1)$  cuyas componentes son todas iguales a 1, y  $V$  una matriz  $(n \times q)$  cuyas columnas son las componentes de los vectores  $v_1, \dots, v_q$ . La matriz  $U$  será por tanto una matriz  $(k \times q)$ , cada una de sus filas define las coordenadas de un punto proyectado en  $V_q$ . Así la proyección del punto  $P_i$  en la variedad lineal  $V_q$ , puede expresarse como:

$$\hat{P}_i = u_{i1}v_1 + \dots + u_{iq}v_q + z \quad (8)$$

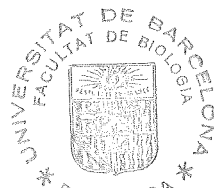
$$i = 1, \dots, k$$

y para las representaciones gráficas, podemos prescindir del vector constante  $z$ , ya que no afecta para nada las distancias entre los puntos proyectados en  $V_q$  (equivale esta supresión a una traslación del sistema de coordenadas), y asociar a cada punto  $P_i$  las coordenadas  $(u_{i1}, \dots, u_{iq})$  de un sistema de referencia cartesiano.

También se verifica que la suma de los cuadrados de las distancias entre las proyecciones de los  $P_i, P_j$ , extendida a todo par  $(i, j)$ , es máxima dentro de la clase formada por todas las variedades lineales de dimensión  $q$ , Cuadras (1981).

## 6.2. ESPACIO PARAMETRICO NO EUCLIDEO

En este caso habrá que distinguir varias situaciones, atendiendo en primer lugar si disponemos o no de una expresión analítica para la distancia.



### 6.2.1. Caso de disponer de una expresión analítica de la distancia

En estas condiciones, dados  $k$  puntos de la variedad riemanniana,  $P_1, \dots, P_k$ , pueden darse dos situaciones: pueden existir  $Q_1, \dots, Q_k$  puntos de un espacio euclídeo, de dimensión  $m \leq k-1$ , tales que la distancia euclídea entre ellos coincida con la distancia riemanniana original, o bien puede darse el caso contrario.

Si definimos:

$$D = (d_{ij})_{k \times k} \quad d_{ij} = d(P_i, P_j) \quad (9)$$

$$A = (a_{ij})_{k \times k} \quad a_{ij} = -\frac{1}{2} d_{ij}^2 \quad (10)$$

$$H = (h_{ij})_{k \times k} \quad H = I - \frac{1}{k} EE^t \quad (11)$$

$$B = H A H \quad (12)$$

donde  $E = (1, \dots, 1)^t$ .

Entonces es bien sabido, Cuadras (1981), que si la matriz  $B$  es semidefinida positiva, de rango  $m$ ,  $m \leq k-1$ , es posible encontrar  $k$  puntos de un espacio euclídeo,  $Q_1, \dots, Q_k$ , tales que la distancia euclídea entre ellos, verifique:

$$d_E(Q_i, Q_j) = d_{ij} \quad (13)$$

Además, es posible obtener explícitamente una configuración de  $k$  puntos de  $\mathbb{R}^m$  que verifiquen (13), mediante la diagonalización de  $B$ .

En efecto, por ser  $B$  simétrica, siempre diagonaliza ortogonalmente, por tanto:

$$B = T D_{\lambda} T^t = (T D_{\lambda}^{1/2}) (D_{\lambda}^{1/2} T^t) = (T D_{\lambda}^{1/2}) (T D_{\lambda}^{1/2})^t \quad (14)$$

con

$$D_{\lambda} = \text{diag}(\lambda_1, \dots, \lambda_m) \quad \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m > 0 \quad (15)$$

siendo  $T$  una matriz  $k \times m$  cuyas columnas son los vectores propios asociados a  $\lambda_1, \dots, \lambda_m$  respectivamente y ortonormales. Si definimos:

$$X = T D_{\lambda}^{1/2} \quad (16)$$

entonces resulta que las filas de  $X$  definen  $k$  puntos  $Q_1, \dots, Q_k$  de  $\mathbb{R}^m$ , al considerar cada fila como las coordenadas de un punto respecto la base canónica de  $\mathbb{R}^m$ , y estos  $k$  puntos de  $\mathbb{R}^m$  verifican (13).

Además, con el cálculo de  $X$  resolvemos a la vez el problema de la representación de los puntos  $P_1, \dots, P_k$  en un subespacio de  $\mathbb{R}^m$  de dimensión  $q$ , con  $q$  generalmente igual a 2, de forma que la suma de los cuadrados de las distancias de los  $k$  puntos al subespacio sea mínima ó bien que la suma de los cuadrados de las interdistancias entre los puntos proyectados en dicho subespacio sea máxima. Bastará escoger las  $q$  primeras columnas de la matriz  $X$ .

En el caso que  $B$  no sea semidefinida positiva, existirán algunos valores propios negativos:

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0 > \lambda_{r+1} \dots \geq \lambda_m \quad (17)$$

en estas condiciones no será posible encontrar una configuración de puntos de un espacio euclídeo que verifique (13). Si procedieramos de forma análoga al caso anterior, al hallar  $D_\lambda^{1/2}$  aparecerían números imaginarios puros y las filas de  $X$  serían de la forma:

$$(x_{j_1}, \dots, x_{j_r}, i x_{j_{r+1}}, \dots, i x_{j_m}) \quad (18)$$

siendo  $i = \sqrt{-1}$ . Cada punto estaría representado en un espacio con ejes reales y con ejes imaginarios puros. Representaciones en dimensión reducida pueden obtenerse escogiendo, por separado, los ejes reales asociados a los mayores valores propios y los ejes imaginarios asociados a los menores valores propios. Sin embargo, la representación respecto los ejes imaginarios presenta problemas de interpretación: cuanto más separados dos puntos, más negativa, y por tanto menor, es su distancia.

Para solventar esta dificultad, se utilizan otros métodos que se basan en transformar la matriz de interdistancias  $D$ , definida en (9), en otra matriz  $\hat{D}$  tal que al repetir el proceso anterior con ésta última, obtengamos que exista una configuración de puntos en un espacio euclídeo tales que sus interdistancias coincidan con los correspondientes elementos de  $\hat{D}$ , en otras palabras los valores propios de la matriz  $B$ , calculada a partir de  $\hat{D}$ , sean no negativos. Para que todo el proceso tenga utilidad la matriz  $\hat{D}$  debe ser "parecida" a la matriz  $D$ , en un sentido que vamos a precisar, introduciendo el concepto de preordenación asociada a una distancia.