

Multivariate Gaussian Analysis

Steffen Lauritzen, University of Oxford

BS2 Statistical Inference, Lecture 12, Hilary Term 2008

February 25, 2008

For a positive definite covariance matrix Σ , the multivariate Gaussian distribution has density on \mathcal{R}^d

$$f(x | \xi, \Sigma) = (2\pi)^{-d/2} (\det K)^{1/2} e^{-(x-\xi)^\top K(x-\xi)/2}, \quad (1)$$

where $K = \Sigma^{-1}$ is the *concentration matrix* of the distribution.

If $X_1 \sim \mathcal{N}_d(\xi_1, \Sigma_1)$ and $X_2 \sim \mathcal{N}_d(\xi_2, \Sigma_2)$ and $X_1 \perp\!\!\!\perp X_2$

$$X_1 + X_2 \sim \mathcal{N}_d(\xi_1 + \xi_2, \Sigma_1 + \Sigma_2).$$

If A is an $r \times d$ matrix, $b \in \mathcal{R}^r$ and $X \sim \mathcal{N}_d(\xi, \Sigma)$, then

$$Y = AX + b \sim \mathcal{N}_r(A\xi + b, A\Sigma A^\top).$$

Partition X into X_1 and X_2 , where $X_1 \in \mathcal{R}^r$ and $X_2 \in \mathcal{R}^s$ with $r + s = d$ and partition mean vector, concentration and covariance matrix accordingly.

Then, if $X \sim \mathcal{N}_d(\xi, \Sigma)$

$$X_2 \sim \mathcal{N}_s(\xi_2, \Sigma_{22}).$$

If Σ_{22} is regular, it further holds that

$$X_1 | X_2 = x_2 \sim \mathcal{N}_r(\xi_{1|2}, \Sigma_{1|2}),$$

where

$$\xi_{1|2} = \xi_1 + \Sigma_{12}\Sigma_{22}^{-1}(x_2 - \xi_2) \quad \text{and} \quad \Sigma_{1|2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}.$$

In particular, *if $\Sigma_{12} = 0$ if and only if X_1 and X_2 are independent.*

From the matrix identities

$$K_{11}^{-1} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} = \Sigma_{1|2} \quad (2)$$

and

$$K_{11}^{-1}K_{12} = -\Sigma_{12}\Sigma_{22}^{-1}, \quad (3)$$

it follows that then the conditional expectation and concentrations also can be calculated as

$$\xi_{1|2} = \xi_1 - K_{11}^{-1}K_{12}(x_2 - \xi_2) \quad \text{and} \quad K_{1|2} = K_{11}.$$

Note that the *marginal covariance is simply expressed in terms of Σ* where as the *conditional concentration is simply expressed in terms of K* .

A square matrix A has *trace*

$$\text{tr}(A) = \sum_i a_{ii}.$$

The trace has a number of properties:

1. $\text{tr}(\gamma A + \mu B) = \gamma \text{tr}(A) + \mu \text{tr}(B)$ for γ, μ being scalars;
2. $\text{tr}(A) = \text{tr}(A^\top)$;
3. $\text{tr}(AB) = \text{tr}(BA)$
4. $\text{tr}(A) = \sum_i \lambda_i$ where λ_i are the *eigenvalues* of A .

For symmetric matrices the last statement follows from taking an orthogonal matrix O so that $OA O^\top = \text{diag}(\lambda_1, \dots, \lambda_d)$ and using

$$\text{tr}(OA O^\top) = \text{tr}(A O^\top O) = \text{tr}(A).$$

The trace is thus *orthogonally invariant*, as is the determinant:

$$\det(OA O^\top) = \det(O) \det(A) \det(O^\top) = 1 \det(A) 1 = \det(A).$$

There is an important trick that we shall use again and again: For $\lambda \in \mathcal{R}^d$

$$\lambda^\top A \lambda = \text{tr}(\lambda^\top A \lambda) = \text{tr}(A \lambda \lambda^\top)$$

since $\lambda^\top A \lambda$ is a scalar.

Consider first the case where $\xi = 0$ and a sample $X_1 = x_1, \dots, X_n = x_n$ from a multivariate Gaussian distribution $\mathcal{N}_d(0, \Sigma)$ with Σ regular. Using (1), we get the likelihood function

$$\begin{aligned} L(K) &= (2\pi)^{-nd/2} (\det K)^{n/2} e^{-\sum_{\nu=1}^n x_{\nu}^{\top} K x_{\nu} / 2} \\ &\propto (\det K)^{n/2} e^{-\sum_{\nu=1}^n \text{tr}\{K x_{\nu} x_{\nu}^{\top}\} / 2} \\ &= (\det K)^{n/2} e^{-\text{tr}\{K \sum_{\nu=1}^n x_{\nu} x_{\nu}^{\top}\} / 2} \\ &= (\det K)^{n/2} e^{-\text{tr}(KW) / 2}. \end{aligned} \tag{4}$$

where

$$W = \sum_{\nu=1}^n x_{\nu} x_{\nu}^{\top}$$

is the matrix of *sums of squares and products*.

Writing the trace out

$$\text{tr}(KW) = \sum_i \sum_j k_{ij} W_{ji}$$

emphasizes that it is linear in both K and W and we can recognize this as a linear and canonical exponential family with K as the canonical parameter and $-W/2$ as the canonical sufficient statistic. Thus, the likelihood equation becomes

$$\mathbf{E}(-W/2) = -n\Sigma/2 = -W/2$$

since $\mathbf{E}(W) = n\Sigma$. Solving, we get

$$\hat{K}^{-1} = \hat{\Sigma} = W/n$$

in analogy with the univariate case.

Rewriting the likelihood function as

$$\log L(K) = \frac{n}{2} \log(\det K) - \text{tr}(KW)/2$$

we can of course also differentiate to find the maximum, leading to

$$\frac{\partial}{\partial k_{ij}} \log(\det K) = w_{ij}/n,$$

which in combination with the previous result yields

$$\frac{\partial}{\partial K} \log(\det K) = K^{-1}.$$

This can also be derived directly by writing out the determinant, and it holds for any non-singular square matrix!

The Wishart distribution is the sampling distribution of the matrix of sums of squares and products. More precisely:

A random $d \times d$ matrix W has a *d -dimensional Wishart distribution* with parameter Σ and n *degrees of freedom* if

$$W \stackrel{D}{=} \sum_{i=1}^n X_i X_i^\top$$

where $X_i \sim \mathcal{N}_d(0, \Sigma)$. We then write

$$W \sim \mathcal{W}_d(n, \Sigma).$$

The Wishart is the multivariate analogue to the χ^2 :

$$\mathcal{W}_1(n, \sigma^2) = \sigma^2 \chi^2(n).$$

If $W \sim \mathcal{W}_d(n, \Sigma)$ its mean is $\mathbf{E}(W) = n\Sigma$.

If W_1 and W_2 are independent with $W_i \sim \mathcal{W}_d(n_i, \Sigma)$, then

$$W_1 + W_2 \sim \mathcal{W}_d(n_1 + n_2, \Sigma).$$

If A is an $r \times d$ matrix and $W \sim \mathcal{W}_d(n, \Sigma)$, then

$$AWA^\top \sim \mathcal{W}_r(n, A\Sigma A^\top).$$

For $r = 1$ we get that when $W \sim \mathcal{W}_d(n, \Sigma)$ and $\lambda \in R^d$,

$$\lambda^\top W \lambda \sim \sigma_\lambda^2 \chi^2(n),$$

where $\sigma_\lambda^2 = \lambda^\top \Sigma \lambda$.

If $W \sim \mathcal{W}_d(n, \Sigma)$, where Σ is regular, then W is regular with probability one if and only if $n \geq d$.

When $n \geq d$ the Wishart distribution has density

$$\begin{aligned} f_d(w | n, \Sigma) \\ = c(d, n)^{-1} (\det \Sigma)^{-n/2} (\det w)^{(n-d-1)/2} e^{-\text{tr}(\Sigma^{-1}w)/2} \end{aligned}$$

for w positive definite, and 0 otherwise.

The Wishart constant $c(d, n)$ is

$$c(d, n) = 2^{nd/2} (2\pi)^{d(d-1)/4} \prod_{i=1}^d \Gamma\{(n+1-i)/2\}.$$