

# MULTIVARIATE VARIANCE COMPONENTS IN TWIN STUDIES

JAN DE LEEUW AND AGATHA LEE

ABSTRACT. We outline the computation of both FIML and REML estimates in homogeneous sibgroups of a multivariate quantitative characteristic, in the case in which there are no covariates. A program in R is provided.

## 1. SINGLE POPULATION

Suppose  $\underline{x}$  and  $\underline{y}$  are  $m$ -dimensional normal random vectors<sup>1</sup> with

$$(1a) \quad \mathbf{E}(\underline{x}) = \mathbf{E}(\underline{y}) = \underline{\mu},$$

$$(1b) \quad \mathbf{E}(\underline{x} - \underline{\mu})(\underline{x} - \underline{\mu})' = \mathbf{E}(\underline{y} - \underline{\mu})(\underline{y} - \underline{\mu})' = \Sigma,$$

$$(1c) \quad \mathbf{E}(\underline{x} - \underline{\mu})(\underline{y} - \underline{\mu})' = \mathbf{E}(\underline{y} - \underline{\mu})(\underline{x} - \underline{\mu})' = \Omega.$$

This implies  $\Sigma \succeq \Omega$  in the Loewner sense, i.e.  $\Sigma - \Omega$  must be positive semi-definite. The intraclass correlation matrix [Rao, 1945] is the matrix  $\Gamma = \Sigma^{-\frac{1}{2}}\Omega\Sigma^{-\frac{1}{2}}$ .

---

*Date:* Saturday 10<sup>th</sup> October, 2009 — 14h 24min — Typeset in KEPLER.

*Key words and phrases.* Multivariate Variance Components, Maximum Likelihood, Twin Studies.

<sup>1</sup>When calculating within a model we underline [Hemelrijk, 1966] the hypothetical random variables whose realizations we observe. The realizations themselves use the same symbol as that for the random variable, but without the underlining. In addition we use overlining for means of vectors of observed quantities, and tildes for the elements of the vector in deviations from the mean.

Define the new random vectors  $\underline{u} = \frac{1}{2}(\underline{x} - \underline{y})$  and  $\underline{v} = \frac{1}{2}(\underline{x} + \underline{y})$ . Then

$$(2a) \quad \mathbf{E}(\underline{u}) = 0,$$

$$(2b) \quad \mathbf{E}(\underline{v}) = \mu,$$

$$(2c) \quad \mathbf{E}(\underline{u}\underline{u}') = \frac{1}{2}(\Sigma - \Omega),$$

$$(2d) \quad \mathbf{E}(\underline{u}\underline{v}') = 0,$$

$$(2e) \quad \mathbf{E}(\underline{v} - \mu)(\underline{v} - \mu)' = \frac{1}{2}(\Sigma + \Omega).$$

Rao et al. [1987]; Srivastava et al. [1988] If we have  $n$  independent realizations  $(u_i, v_i)$  of  $(\underline{u}, \underline{v})$ , then the deviance <sup>2</sup> is

$$(3) \quad \mathcal{D}(\Delta, \Xi, \mu) = n \log \mathbf{det}(\Delta) + n \log \mathbf{det}(\Xi) + \\ + \sum_{i=1}^n u_i' \Delta^{-1} u_i + \sum_{i=1}^n (v_i - \mu)' \Xi^{-1} (v_i - \mu).$$

where  $\Delta = \frac{1}{2}(\Sigma - \Omega)$  and  $\Xi = \frac{1}{2}(\Sigma + \Omega)$ .

Thus the maximum likelihood estimate of  $\mu$  is

$$\hat{\mu} = \bar{v} = \frac{1}{2}(\bar{x} + \bar{y}),$$

which is the mean vector of all  $2n$  vectors.

Define  $G = \frac{1}{n} \sum_{i=1}^n u_i u_i'$  and  $H = \frac{1}{n} \sum_{i=1}^n \tilde{v}_i \tilde{v}_i'$ . The concentrated negative log-likelihood is

$$(4) \quad \mathcal{D}_\star(\Delta, \Xi) \stackrel{\Delta}{=} \min_{\mu} \mathcal{D}(\Delta, \Xi, \mu) = n[\log \mathbf{det}(\Delta) + \log \mathbf{det}(\Xi) + \\ + \mathbf{tr} \Delta^{-1} G + \mathbf{tr} \Xi^{-1} H].$$

Maximum likelihood estimates of the variance components can be computed as

$$\hat{\Sigma} = \hat{\Xi} + \hat{\Lambda} = H + G,$$

$$\hat{\Omega} = \hat{\Xi} - \hat{\Lambda} = H - G.$$

---

<sup>2</sup>The deviance is twice the negative log-likelihood, except for irrelevant constants.

Note that  $\hat{\Sigma} - \hat{\Omega} = 2G \succeq 0$ , but it is not guaranteed that  $\hat{\Omega} \succeq 0$ . Maximum likelihood estimates that guarantee positive semi-definiteness are discussed in the Appendix.

Since  $\mathbf{E}(\underline{G}) = \frac{1}{2}(\Sigma - \Omega)$  and  $\mathbf{E}(\underline{H}) = \frac{1}{2} \frac{n-1}{n}(\Sigma + \Omega)$  we can compute unbiased estimates using

$$\begin{aligned}\hat{\Sigma} &= \frac{n}{n-1}H + G, \\ \hat{\Omega} &= \frac{n}{n-1}H - G.\end{aligned}$$

These unbiased estimates are also the restricted maximum likelihood or REML estimates.

## 2. HERITABILITY

The ACE model for monozygotic and dizygotic twins decomposes variation into additive genetic, common environmental, and unique environmental variation. In the simplest case we have for the two types of twins

$$\begin{aligned}\Sigma_M &= \Sigma_D = \Theta_A + \Theta_C + \Theta_E, \\ \Omega_D &= \frac{1}{2}\Theta_A + \Theta_C, \\ \Omega_M &= \Theta_A + \Theta_C.\end{aligned}$$

Solving gives the unbiased estimates

$$\begin{aligned}\hat{\Theta}_A &= 2(\Omega_M - \Omega_D) = 2\left(\frac{n}{n-1}(H_M - H_D) - (G_M - G_D)\right), \\ \hat{\Theta}_C &= 2\Omega_D - \Omega_M = \frac{n}{n-1}(2H_D - H_M) - (2G_D - G_M), \\ \hat{\Theta}_E &= \Sigma_M - \Omega_M = 2G_M.\end{aligned}$$

## APPENDIX A. A MATRIX PROBLEM

Consider the problem of minimizing

$$\log \mathbf{det}(\Delta) + \log \mathbf{det}(\Xi) + \mathbf{tr} \Delta^{-1}G + \mathbf{tr} \Xi^{-1}H$$

over the two matrices  $\Delta$  and  $\Xi$ , with  $\Xi \succeq \Delta > 0$ . Clearly if  $H \succeq G > 0$  the solution is  $\hat{\Xi} = H$  and  $\hat{\Delta} = G$ .

Assume  $H > 0$ . Then there is a nonsingular  $S$  such that  $G = S\Phi S'$  and  $H = SS'$  with  $\Phi \succeq 0$  diagonal. Define  $\bar{\Delta} = S'\Delta^{-1}S$  and  $\bar{\Xi} = S'\Xi^{-1}S$ . Suppose  $\tilde{\xi}_i$  and  $\tilde{\delta}_i$  are the diagonal elements and  $\lambda_i$  and  $\theta_i$  are the eigenvalues of  $\tilde{\Xi}$  and  $\tilde{\Delta}$ . Then

$$\sum_{i=1}^n -\log \lambda_i - \log \theta_i + \xi_i \phi_i + \lambda_i$$

## REFERENCES

- J. Hemelrijk. Underlining Random Variables. *Statistica Neerlandica*, 20:1–7, 1966.
- C.R. Rao. Familial Correlations or the Multivariate Generalization of the Intraclass Correlation. *Current Science*, 14:66–67, 1945.
- D.C. Rao, G.P. Vogler, M. McGue, and J.M. Russell. Maximum Likelihood Estimation of Familial Correlations from Multivariate Quantitative Data on Pedigrees: A General Method and Examples. *American Journal of Human Genetics*, 41:1104–1116, 1987.
- M.S. Srivastava, K.J. Keen, and R.S. Katapa. Estimation of Interclass and Intraclass Correlations in Multivariate Familial Data. *Biometrics*, 44:141–150, 1988.

(Jan de Leeuw) DEPARTMENT OF STATISTICS, UNIVERSITY OF CALIFORNIA LOS ANGELES

*E-mail address*, Jan de Leeuw: [deleeuw@stat.ucla.edu](mailto:deleeuw@stat.ucla.edu)

(Agatha Lee) LABORATORY OF NEURO IMAGING, UNIVERSITY OF CALIFORNIA, LOS ANGELES

*E-mail address*, Agatha Lee: [alee.loni@gmail.com](mailto:alee.loni@gmail.com)