

## Bayesian discrimination with longitudinal data

P. J. BROWN

*University of Kent at Canterbury, Institute of Mathematics and Statistics, Cornwallis Building,  
Canterbury, CT2 7NF, UK  
Philip.J.Brown@ubc.ac.uk*

M. G. KENWARD

*London School of Hygiene and Tropical Medicine, Department of Epidemiology and Population  
Health, Kennel Street, London, WC1E 7HT, UK*

E. E. BASSETT

*University of Kent at Canterbury, Institute of Mathematics and Statistics, Cornwallis Building,  
Canterbury, CT2 7NF, UK*

### SUMMARY

The motivation for the methodological development is a double-blind clinical trial designed to estimate the effect of regular injection of growth hormone, with the purpose of identifying growth hormone abusers in sport. The data formed part of a multicentre investigation jointly sponsored by the European Union and the International Olympic Committee. The data are such that for each individual there is a matrix of marker variables by time point (nominally 8 markers at each of 7 time points). Data arise out of a double-blind trial in which individuals are given growth hormone at one of two dose levels or placebo daily for 28 days. Monitoring by means of blood samples is at 0, 21, 28, 30, 33, 42 and 84 days.

We give a new method of Bayesian discrimination for multivariate longitudinal data. This involves a Kronecker product covariance structure for the time by measurements (markers) data on each individual. This structure is estimated by an empirical Bayes approach, using an ECM algorithm, within a Bayesian Gaussian discrimination model. In future one may have markers for an individual at one or more time points. The method gives probabilities that an individual is on placebo or on one of the two dose regimes.

*Keywords:* Bayesian methods; Discrimination; Double blind clinical trial; ECM algorithm; Empirical Bayes; Kronecker covariance structures; Repeated measures; Smoothing; Type II likelihood.

### 1. INTRODUCTION

Multivariate data on individuals are often collected repeatedly so that an individual's observation consists of a data matrix whose covariability is across time and variables. Each individual may belong to one of several different groups. In future an individual may present some or all of the variables at a subset of the times and one wishes to calculate the probability that the individual belongs to each group, and perhaps even to assess where in the time trajectory the individual fits.

Earlier comprehensive references to the general area of multivariate discrimination are provided by McLachlan (1992); Krzanowski and Marriott (1995). There is relatively little work on discrimination involving repeated measures. In the two-population case, a components of variance model is developed

by Logan and Gupta (1993). This extends Geisser (1964), as does Brown (1993, chapter 8) in the direction of informative prior structures, albeit natural conjugate ones.

Two-way covariance structures have been of interest in other areas: those of repeated measures and spatial statistics. For the latter field see Brown *et al.* (1994). Many aspects of our modelling joint covariance structures have this more general applicability. Our emphasis on discrimination arises out of our work on the detection of growth hormone abuse in sport (GH-2000 EU/IOC 1999 BMH14 CT950678, Co-ordinator P. H. Sönksen), a multi-centre project analysing the effects of growth hormone in a range of trials, as described in more detail in Section 2. In this paper we are concerned with the methodology for dealing with multivariate discrimination when repeated measures are involved, using the growth hormone data as a motivating example.

## 2. GH2000 DOUBLE-BLIND STUDY

This section describes the application of our methodology to one of the datasets generated through the GH2000 project mentioned in the introduction: a double-blind multi-centre study. Pituitary GH is secreted in pulses mainly in response to exercise, sleep and stress. GH acts on the liver and on bone, stimulating the formation of a number of substances ('markers') which are then secreted into the blood stream producing a characteristic response. Exogenous recombinant human growth hormone (rhGH) administration, or something stimulating endogenous GH release, but nothing else, produces a characteristic pattern in these markers. Four of our chosen eight markers, IGF-1, IGFBP<sub>2</sub>, IGFBP<sub>3</sub> and ALS, enter the blood stream via the liver, whereas the other four, Osteocalcin, P-III-P, PICP, ICTP, are bone related. Whereas GH disappears very rapidly, these markers persist much longer, extinguishing themselves at different rates and producing a characteristic fingerprint. Thus we use the 'downstream' effects of growth hormone as the basis of testing for GH abuse. The trajectory of the markers for individuals over time is the starting point for developing a methodology for Bayesian discrimination on future presentation of a blood sample by an athlete.

### 2.1 *The data*

These consisted of around 100 individual volunteers who received by injection a single dose, a double dose or placebo daily for 28 days. The individuals and the researchers taking blood samples were both blind as to which of the three dose regimes the individual had been randomized. Blood samples were taken at baseline, time 0, just before the first dose, at 21 and 28 days and subsequent to growth hormone intake at 30, 33, 42 and 84 days. These blood samples were split, refrigerated, and sent to two laboratories to measure the markers, each laboratory concentrating on four of the markers.

### 2.2 *Data exploration*

Several analyses have been carried out on these and related data. Our example here is for illustrative purposes concentrating on males. The best known of the markers, the insulin-like growth factor I (IGF-I) shows a characteristically raised response to administration of recombinant human growth hormone. This can be seen from the plot of IGF-I against time for the males of the study in Figure 1. Those individuals on placebo have a flat trajectory, whereas the responses of those on a single or double dose of GH rise sharply over the first 28 days.

A plot of IGF-I against P-III-P at 21 days in Figure 2 shows there is complete separation at this time for this pair of markers in combination although not for either individually. A linear discriminant function optimized for day 21 is applied to the data and given as a dotplot in Figure 3. The 'clean' individual's data, the baseline measurements for all individuals and all time points for placebo measurements have

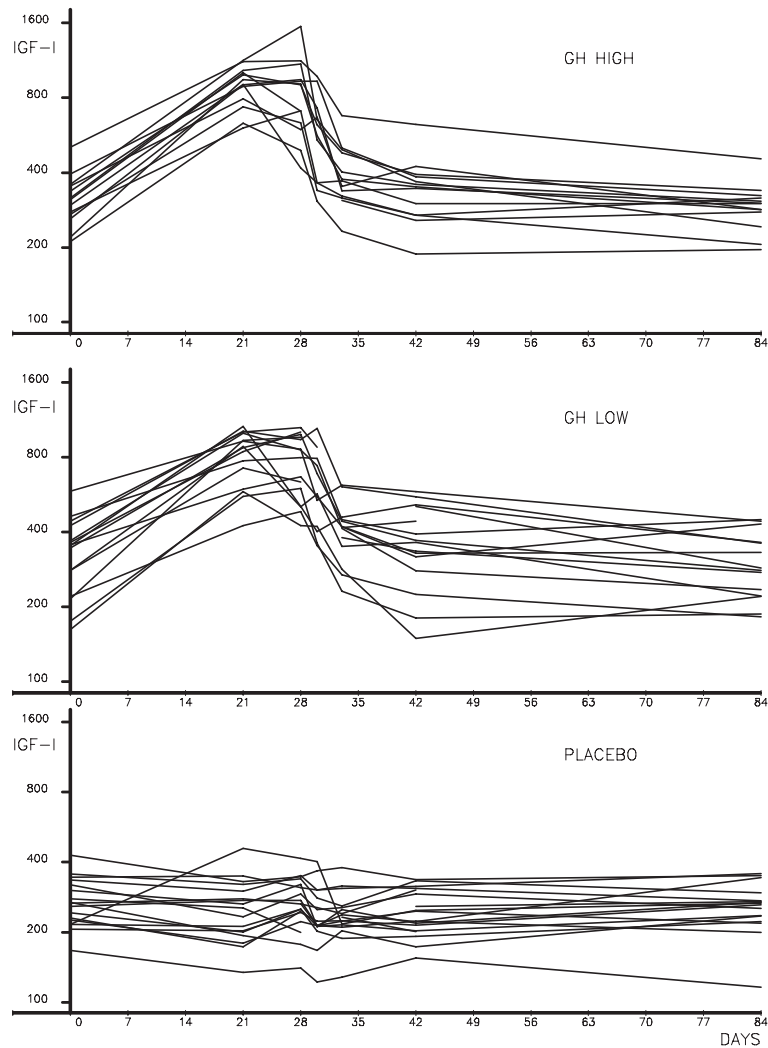


Fig. 1. IGF-I against time in days for placebo, GH low, GH high.

been combined, centred on zero, and scaled so as to have unit variance. They are all categorized as 'base'. Day 21 and 28 dosed individuals have discriminant function values on average four or five standard deviations above the placebo mean, and although there is some overlap of upper tail of placebos and the lower tail of dosed individuals it is easy to see the potential for setting an allocation rule with very high sensitivity, say 1 in 10 000, which would correspond to 3.5 standard deviations above the placebo mean. Rather than discriminate at a single time point using only data from that time, our approach uses the whole trajectory of time points to obtain an optimal discriminator at any particular time.

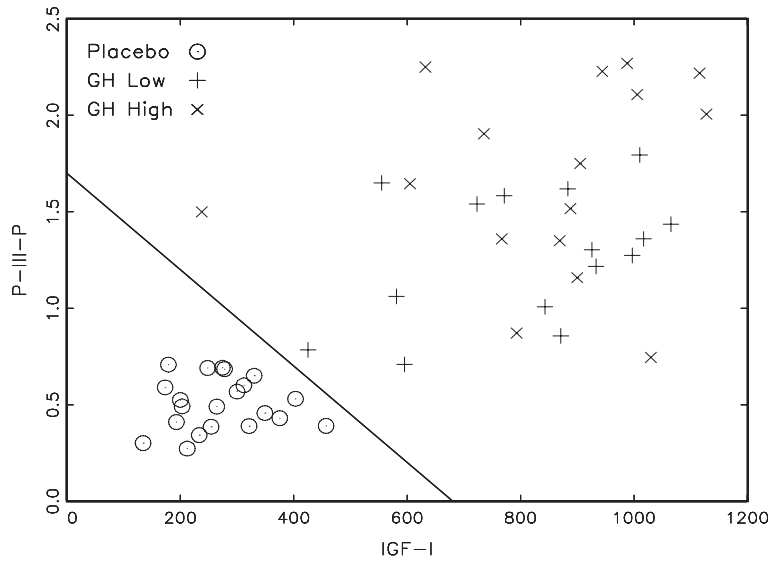


Fig. 2. IGF-I versus PIIP, day 21. Groups: 0 (placebo), 1 (GH low), 2 (GH high).

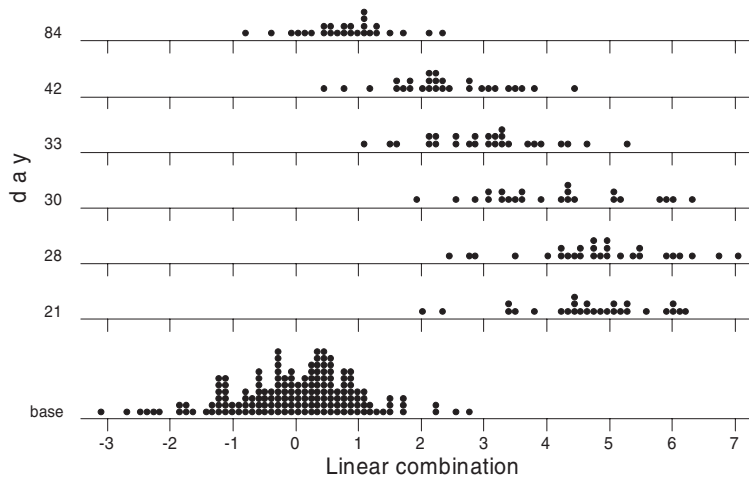


Fig. 3. Discriminant scores evaluated at 21 days plotted by time for all males, low and high dose at days 21–84, base is ‘clean’ data (baseline and placebo).

### 2.3 Training and validation

The data from the double-blind study consisted of 43 male individuals with complete sets of blood samples taken at all seven time points. The other nine male individuals had samples at fewer than seven points. We used the complete data as training data and the individuals with incomplete records for validation. The pattern of missing measurements was not generally one of dropout, rather some intervening time point was missed, so we felt somewhat reassured that perhaps those with missing data were not particularly different from those without. We did check back as far as was feasible to find out whether there were

particular reasons for being missing, but none emerged.

It may be noted that even though these validation records are incomplete, since we are using them at a single future time point, they provide several time points for validation.

Females in the study behaved quite differently in their response to growth hormones, and their results will be reported elsewhere. A more thorough analysis will need to adjust also for age as a covariate. Adjustment for age would seem to be desirable because of strong age dependence of some of the biochemical markers (for example, IGF-I). This may not be important for validation of the model for the double-blind data alone but it will be needed to form a common ground for adjustment in the light of the cross-sectional study of elite athletes, another component of the GH2000 project. All eight markers were characteristically positively skewed but were symmetric and normal when logged, and this form was used throughout.

### 3. THE MODEL

First let us establish some notation. The number of groups is  $g + 1$ , which will be 3 in the example considered, corresponding to doses  $d = 0$  (placebo),  $d = 1$  and  $d = 2$ . The training data, for which we know both variables and population group, comprise  $n_d$  independent observations forming a  $\sum_0^g n_d \times qp$  matrix where the  $qp$  observations on an individual are ordered by time and marker, all the  $p$  biochemical markers being measured at  $q$  time points. We may also need to include individual-level covariates such as age, sex or weight. For the  $d$ th population we assume multivariate normality, perhaps after a suitable transformation. Conditional on parameters  $\alpha_d, B, \Omega$  the standard multivariate normal regression model assumed is (in the matrix variate notation discussed in Appendix B)

$$Y_d - 1\alpha'_d - X_d B \sim \mathcal{N}(I_{n_d}, \Omega), \tag{1}$$

with  $n_d \times qp$  random matrix  $Y_d$ ,  $1$  an  $n_d \times 1$  vector of ones,  $n_d \times c$  model matrix  $X_d$  regarded as fixed and corresponding to  $c$  covariates,  $\alpha_d$  a  $qp \times 1$  vector of intercepts, and  $B$  the  $c \times qp$  matrix of regression coefficients for these covariates (for example, age, body mass index, etc). The matrix  $Y_d$  could have been structured as a three-way array of dimensions  $n_d \times q \times p$ , and although we will avoid such notational considerations at this point our model and analysis will separately structure the times and variables.

Here, by virtue of the identity matrix first argument of  $\mathcal{N}(\cdot, \cdot)$ , rows of  $Y_d$  are independently distributed. We have chosen a relatively simple covariance structure where the covariance matrix  $\Omega$  does not depend on dose level. Straightforward model elaboration would have  $\Omega_d$  and possibly also  $B_d$ . The former implicitly allows quadratic discrimination, the latter allows differing dependence on the covariates by groups. Without loss of generality we assume that columns of the concatenated matrix  $X$ , ( $\sum_0^g n_d \times c$ ) where  $X' = (X'_0, X'_1, \dots, X'_g)$ , have been centred by subtracting their column means. Similarly, the matrix  $Y$  ( $\sum_0^g n_d \times qp$ ) is post hoc centred by subtracting means of columns.

The special forms of prior distributions for parameters  $\alpha_d, B, \Omega$  are given as follows. First, given  $\Omega$ , for the intercepts  $\alpha_d$  ( $qp \times 1$ ) with prior mean  $\alpha_{d0}$ ,

$$\alpha'_d - \alpha'_{d0} \sim \mathcal{N}(h, \Omega); \tag{2}$$

secondly and independently, given  $\Omega$ , for the regression coefficients  $B$  ( $c \times qp$ ) with prior mean  $B_0$ ,

$$B - B_0 \sim \mathcal{N}(H, \Omega). \tag{3}$$

Notice that from our matrix variate characterization, both priors (2) and (3) have covariances dependent on  $\Omega$  in a way that directly extends the univariate regression natural conjugate prior distributions. We could have combined  $\alpha$  and  $B$  but we felt it helpful to emphasize the different status of the covariates to which

they relate. Often a different prior distribution would be appropriate and scalar  $h$  would be constructed differently from  $H$  ( $c \times c$ ), although here we assume uninformative prior distributions in both cases. Thus the scalar hyperparameter  $h$  of the prior distributions for  $\alpha_d$  given in (2) will be assigned a large value, tending to infinity, when the value ascribed to the prior mean,  $\alpha_{d0}$ , becomes irrelevant. We will also not be specially interested in imposing anything other than vague prior structure on the regression coefficient matrix  $B$ , achieved by setting  $H$  to be very large. This has the added simplifying advantage of decoupling the estimation of these parameters from that of the covariance structure  $\Omega$ .

Now the prior marginal distribution of  $\Omega$  ( $qp \times qp$ ) is taken to be inverse Wishart,

$$\Omega \sim \mathcal{IW}(\delta; Q), \quad (4)$$

in the notation of Appendix B. Here  $Q$  is a  $qp \times qp$  scale matrix of the inverse Wishart, and  $\Omega$  has prior expectation  $Q/(\delta - 2)$  for  $\delta > 2$ . The shape parameter  $\delta$  is such that  $\delta > 0$  for a proper prior density. Improper prior densities would have  $\delta \leq 0$ , and the Jeffreys invariant improper density corresponds to  $\delta = 1 - qp$ .

The scale matrix hyperparameter  $Q$  of the prior distribution for  $\Omega$  from equation (4) is given as proportional to a Kronecker form  $\Gamma \otimes \Sigma$ . Here  $\Gamma$  is a  $q \times q$  matrix of covariation across time, whereas  $\Sigma$  is  $p \times p$  and measures covariability across markers. Other forms are of course possible, but here this Kronecker product structure is a key simplifying form, borne out by empirical evidence in Section 5.1. It can be argued that imposing a Kronecker product structure at the hyperparameter level of the prior mean of this covariance matrix, results in robust inference. A small value of  $\delta$  will not strongly impose the structure whereas a large value will effectively impose the structure at the first level of covariance  $\Omega$ . This Kronecker product model applies when we have a multiplicative random-effects model:  $Y_{tr} - \mu_{tr} = \theta_t \times \gamma_r$ , where  $Y_{tr}$  denotes the  $r$ th response at the  $t$ th time point,  $r = 1, \dots, p$ ,  $t = 1, \dots, q$ . Brown *et al.* (1994) in modelling spatial variation across areas rather than time used a full  $pq \times pq$  covariance structure, with a similar Kronecker product form at the hyperparameter level. Like them we will adopt an empirical Bayes approach to estimation of the covariance structure, following (Chen, 1979). We prefer, however, to fix the shape parameter  $\delta$  at a selection of values and examine the ‘likelihood’ of these, so exploring sensitivity to its choice rather than estimating it. This simplifies and speeds up the algorithms and offers more control. It does mean that we will need to retain terms solely in  $\delta$ , or equivalently the degrees of freedom  $\nu = \delta + qp - 1$ , in our relevant likelihoods.

## 4. BAYESIAN INFERENCE

### 4.1 Posterior and marginal distributions

The probability density function of  $Y_d$ ,  $d = 0, \dots, g$  from model (1) is

$$f_Y(Y_d | \alpha_d, B, \Omega) \propto |\Omega|^{-n_d/2} \exp\{-\frac{1}{2} \text{trace}(Y_d - 1\alpha'_d - X_d B)\Omega^{-1}(Y_d - 1\alpha'_d - X_d B)'\}, \quad (5)$$

for  $d = 0, \dots, g$ . Explicit forms for the constant of proportionality and other matrix-variate densities needed below can be found in Brown (1993, Appendix A).

The likelihood marginalized over the vague priors for  $\alpha_d$ ,  $d = 0, \dots, g$ , and  $B$  is shown in Appendix B for given  $\Omega$  as proportional to

$$|\Omega|^{-N/2} \exp\{-\frac{1}{2} \text{trace}(\Omega^{-1}S)\}, \quad (6)$$

with centred  $Y_d^* = Y_d - 1\bar{Y}_d$  and sum of products of residuals

$$S = \sum_d (Y_d^* - X_d^* \hat{B})'(Y_d^* - X_d^* \hat{B}), \quad (7)$$

where  $X_d^* = X_d - 1\bar{X}_d$  and

$$\hat{B} = \left[ \sum X_d^{*'} X_d^* \right]^{-1} \left( \sum X_d^{*'} Y_d \right).$$

In the process *a posteriori* given  $\Omega$ ,  $Y$  it is shown that

$$B - \hat{B} \sim \mathcal{N} \left( \left[ \sum X_d^{*'} X_d^* \right]^{-1}, \Omega \right). \quad (8)$$

The final stage injects prior knowledge from the inverse Wishart prior for  $\Omega$  in (4). We take the scale matrix  $Q = \nu Q^*$  so that the inverse of  $\Omega$  given as  $\Omega^{-1} = \Lambda$  is Wishart( $\nu$ ,  $[\nu Q^*]^{-1}$ ) and has expectation  $(Q^*)^{-1}$ , independent of the degrees of freedom. The form of the prior density of  $\Omega$  is

$$\pi(\Omega) = c(qp, \delta) |\Omega|^{\nu/2} |\Omega|^{-(\nu+qp+1)/2} \exp\{-\frac{1}{2} \text{trace}(\Omega^{-1} Q)\} \quad (9)$$

with  $c(l, \delta) = 2^{-l\nu/2} / \Gamma_l(\nu/2)$  and  $\Gamma_l(\frac{\nu}{2}) \propto \prod_{i=1}^l \Gamma[\frac{\nu}{2} - (i-1)/2]$ . Here we have used the alternative parameters  $\delta, \nu$  interchangeably where  $\delta = \nu - l + 1$  with dimension  $l = qp$ . Multiplying the likelihood (6) by the prior density (9) and integrating out  $\Omega$  gives the integrated or type II likelihood (Good, 1965) of  $Q$  as

$$L(Q; Y) = \frac{c(qp, \delta)}{c(qp, \delta^*)} |\Omega|^{\nu/2} |Q + S|^{-(\nu+N)/2} \quad (10)$$

with  $\delta^* = \delta + N$ . In the limit as  $\nu \rightarrow \infty$  this becomes the standard likelihood for  $Q$  as if  $Q$  were the level I parameter. Explicitly with  $Q = \nu Q^*$

$$L(Q^*; Y) \propto |Q^*|^{-N/2} \exp\{-\frac{1}{2} \text{trace}[(Q^*)^{-1} S]\} \quad (11)$$

where this normal likelihood has the constant of proportionality  $(2\pi)^{-qpN/2}$ . This is straightforward to maximize when  $Q^* = \Gamma \otimes \Sigma$  using a flip-flop algorithm: alternatively maximising over  $\Gamma$  for fixed  $\Sigma$  and then  $\Sigma$  for fixed  $\Gamma$ , as shown in Brown *et al.* (1999b). The same algorithm will be developed below in a slightly different context. The non-asymptotic likelihood (10) does not seem straightforward to maximize under the Kronecker product structure: we rather take an EM algorithm route suggested by Dempster *et al.* (1977) and developed by Chen (1979). The use of a conditional maximization within the M-step renders it an ECM algorithm (Meng and Rubin, 1993).

#### 4.2 The ECM algorithm

Chen (1979) considers parameter  $\nu$  as well as  $\Omega$ , whereas we concentrate on  $\Omega$  alone. Taking  $\Lambda = \Omega^{-1}$ , the complete data is  $(S, \Lambda)$ , the incomplete data is  $S$  alone and  $\Lambda$  is the complete data sufficient statistic for  $Q^*$ . Now the posterior distribution of  $\Omega$  for given  $Q^*$  is inverse Wishart,  $\mathcal{IW}(\delta + N, S + \nu Q^*)$  and  $\Lambda$  is Wishart,  $\Lambda \sim \mathcal{W}[\nu + N, (S + \nu Q^*)^{-1}]$ . Thus for the E-step we have

$$E(\Lambda) = (\nu + N)(S + \nu Q^*)^{-1}. \quad (12)$$

For the M-step a new value of  $Q^*$  is found by maximising the density of  $\Lambda$  given  $Q^*$ , that is the Wishart prior density for  $\Lambda$  with the current value of  $\Lambda$  taken at the E-step (12). This means that we are required to maximize over  $Q^*$

$$H(Q^*) = \log |Q^*| - \text{trace}(\Lambda Q^*). \quad (13)$$

This can be achieved by the same type of flip-flop used in the maximization of (11), although further simplifications are evident in the case of the normal likelihood. Notice that by virtue of the parametrization  $Q = \nu Q^*$ , as used by Chen (1979), the maximization of (13) is decoupled from the prior degrees of freedom  $\nu$ .

We now consider the Kronecker product structure for  $Q^*$ , setting

$$Q^* = \Gamma \otimes \Sigma. \quad (14)$$

The quantity to maximize is

$$H(Q^*) = p \log |\Gamma| + q \log |\Sigma| - \text{trace} [\Lambda(\Gamma \otimes \Sigma)]. \quad (15)$$

The value of  $\Lambda$  is set at the current E-step value (12). Keeping this fixed we need to cycle through alternate maximizations of  $\Gamma$  and  $\Sigma$ . For given current  $\Gamma$  we need to maximize

$$q \log |\Sigma| - \text{trace} (\Sigma A). \quad (16)$$

In (16),  $A = \sum_{i=1}^q \sum_{j=1}^q \gamma_{ij} B_{ij}$ , where  $B_{ij}$  is the  $ij$ th  $p \times p$  submatrix of  $\Lambda$ , partitioned into blocks and  $\Gamma = (\gamma_{ij})$ . The maximum of (16) is achieved at  $\Sigma^{-1} = A/q$  with a consequent increase in  $H(Q^*)$  in (15) due to concavity (Mardia *et al.*, 1979, p. 104). Fixing  $\Sigma$  at this new value we can achieve a further increase by maximising over  $\Gamma$  to give  $\Gamma^{-1} = A^*/p$  where  $A^* = \sum_{i=1}^p \sum_{j=1}^p \sigma_{ij} B_{ij}^*$  and  $B_{ij}^*$  is the  $ij$ th  $q \times q$  submatrix of  $\Lambda$ , permuted to be varying fastest for time rather than variables, in the order implied by  $\Sigma \otimes \Gamma$ . Cycling in this way through successive maximizations of  $\Sigma$  and  $\Gamma$  will converge to a unique maximum of

$$\hat{Q}^* = \hat{\Gamma} \otimes \hat{\Sigma},$$

even though there is no uniqueness about  $\Gamma$  and  $\Sigma$  separately, just their product.

After this cycle of maximizations within the M-step we update the value of  $\Lambda$  through the E-step and repeat the cycle of maximizations within the M-step. We repeat the E- and M-steps until convergence. The monotonic increase in likelihood and convergence is guaranteed by the exponential family formulation and EM theory (Meng and Rubin, 1993). Hybrid algorithms which adopt fewer CM cycles per iteration are also possible, but exploration of alternatives was not warranted given the fast convergence (seconds rather than minutes) of the algorithm.

At designated convergence of the overall E- and M-step process we can evaluate the type II likelihood at the given  $(\nu, \hat{Q}^*)$  in (10). The smallest integer value of  $\nu$  such that  $\delta > 0$  would be  $\nu = qp$ . We may also calculate the likelihood for  $\nu \rightarrow \infty$  from (11) with the appropriate constant of proportionality. A few intervening values of  $\nu$  may also be desirable. The precise value of  $\nu$  may influence the tail probabilities for prediction, and we will supplement its estimation by a type II maximum likelihood estimate, denoted  $\hat{\nu}$ , with the two estimates of  $\nu$  at 2 units of log likelihood below the maximum.

This completes the specification of the posterior distribution of the parameters  $(\alpha_0, \dots, \alpha_g, B, \Omega)$  as a normal-inverse Wishart with estimated hyperparameters  $(\hat{\alpha}_0, \dots, \hat{\alpha}_g, \hat{B}, \hat{Q}^*, \hat{\nu})$ .

### 4.3 Predictive distribution and discrimination

We are first interested in deriving the predictive distribution for a future  $qp$ -dimensional  $Y$ -vector from group  $d$ , and denoted as the  $1 \times qp$  response vector  $Z$ , from the same model as (1), with covariates,  $x$  ( $1 \times c$ ), measured on the same scale as with the training data. The model, analogous to (1), and using the same matrix notation for consistency is

$$Z - \alpha'_d - xB \sim \mathcal{N}(1, \Omega). \quad (17)$$

Conditional on  $(B, \Omega, Y)$ , using the conditional posterior of  $\alpha_d$  in (5) we have

$$Z - (\bar{Y}_d - \bar{X}_d B) - xB \sim \mathcal{N}(1 + 1/n_d, \Omega). \quad (18)$$

Now we average this over the posterior distribution of  $B$  for given  $\Omega, Y$ . This is symbolically accomplished by multiplying (8) by  $(x - \bar{X}_d)$  on the left and adding the resultant distribution to (18) to eliminate  $B$ . This gives

$$Z^* - x^* \hat{B} \sim \mathcal{N}\left(1 + 1/n_d + x^* \left[\sum X_d^{*'} X_d^*\right]^{-1} x^{*'}, \Omega\right), \quad (19)$$

with mean corrected  $x^* = (x - \bar{X}_d)$ ,  $Z^* = Z - \bar{Y}_d$ . The posterior distribution of  $\Omega$  is  $\mathcal{IW}(\hat{\delta} + N, \hat{Q} + S)$  where  $\hat{\delta} = \hat{\nu} - qp + 1$ , and  $\hat{Q} = \hat{\nu} \hat{Q}^*$ . Hence, we finally have the predictive distribution of  $Z$  as multivariate Student- $T$ . In our matrix notation (see Appendix B)

$$Z^* - x^* \hat{B} \sim \mathcal{T}\left(\hat{\delta} + N; 1 + 1/n_d + x^* \left[\sum X_d^{*'} X_d^*\right]^{-1} x^{*'}, \hat{Q} + S\right). \quad (20)$$

We will also be interested in predicting sub-vectors of  $Z$ . With this aim let  $\eta$  be a binary vector of zeros and ones of length  $qp$  so that  $Z^\eta$  identifies (through ones) the chosen elements of the vector  $Z$ . The predictive distribution of the sub-vector identified by  $\eta$  is then

$$Z^{*\eta} - x^* \hat{B}^\eta \sim \mathcal{T}\left(\hat{\delta} + N; 1 + 1/n_d + x^* \left[\sum X_d^{*'} X_d^*\right]^{-1} x^{*'}, \hat{Q}^\eta + S^\eta\right), \quad (21)$$

where  $\eta$  selects columns of  $Z^*$ ,  $B$  and rows and columns of  $\hat{Q}$ ,  $S$ . Typically, we will be interested in the sub-vector that corresponds to all  $p$  marker measurements at a particular time point,  $t$ . Then  $\hat{Q}^\eta + S^\eta = \hat{\nu} \hat{\gamma}_{tt} \hat{\Sigma} + S_{tt}$ , with  $S_{tt}$  the  $p \times p$  sum of products matrix for the  $t$ th time point. Imagine a single time point from the prediction data: that is, an individual whose GH-dose history is unknown provides  $p$  biochemical measurements via a blood sample. We assume such an individual history is like taking  $p$ -markers out of the time trajectory of  $q$  sets of  $p$  measurements. Thus the  $p \times 1$  vector  $Z$  for given time  $t$  and dose  $d$  is as model (21) with the time-specific selection vector  $\eta$ , and whose probability density function is denoted  $f(Z|t, d, Y, x)$ . This future observation has an unknown value of  $t$  and  $d$  and these are taken *a priori* independent and given as

$$\pi(t, d) = \pi(t)\pi(d). \quad (22)$$

Here  $\pi(t)$  should have zero probability on  $t = 0$  which provides baseline measurements confounded with placebo.

The posterior distribution is given as

$$\pi(t, d | Z, Y, x) \propto f(Z | t, d, Y, x) \pi(t, d). \quad (23)$$

If in (21) we let

$$\begin{aligned} Z_{xd} &= Z^{*\eta} - x^* \hat{B}^\eta, \\ h_{xd} &= 1 + 1/n_d + x^* \left[\sum X_d^{*'} X_d^*\right]^{-1} x^{*'}, \\ \delta^* &= \hat{\delta} + N, \\ R_t &= \hat{\nu} \hat{\gamma}_{tt} \hat{\Sigma} + S_{tt} \end{aligned} \quad (24)$$

then this becomes

$$Z_{xd} \sim \mathcal{T}(\delta^*; h_{xd}, R_t). \quad (25)$$

The kernel of the density  $f(Z_{xd} | t, d, Y, x)$  is that of a multivariate Student density. This is then

$$f(Z | t, d, Y, x) \propto (h_{xd})^{\delta^*/2} |R_t|^{-1/2} \{h_{xd} + Z_{xd}(R_t)^{-1} Z'_{xd}\}^{-(\delta^*+p)/2}, \quad (26)$$

a simple function of the squared distance  $Z_{xd}(R_t)^{-1} Z'_{xd}$  in the common metric of  $(R_t)^{-1}$ . When  $\delta$  is large so that the Student kernel becomes that of a normal, then the quadratic dependence on  $Z$  would cancel top and bottom in (23) if  $\gamma_{tt}$  did not vary with time, and discrimination would be implicitly linear. With our general  $\gamma_{tt}$  multiplier the relative odds at a particular time loses the quadratic dependence in this asymptotic case.

The case of multiple time point data on a future individual is also subsumed within the above framework. For two time points the selection vector  $\eta$  would be zero except for the  $p$  variables at each of the two time points. The formulae above apply aside from  $R_t$  in (24) being specified from the more general form given as the  $2p \times 2p$  matrix

$$R_t = \hat{Q}^\eta + S^\eta. \quad (27)$$

The time index is now bivariate and would naturally consist of  ${}^7C_2$  distinct ordered pairs. Since one would typically know the time lag between sampling points for the individual, there is a temptation to restrict attention to the pairs of observations with similar elapsed time within the double-blind trial. This would be unwise however, as elapsed time for a random competitor would be unlikely to correspond to elapsed time in the controlled doping regime of the double-blind trial. It does suggest that one may wish to gather further information from a future additional double-blind trial, featuring several sampled points within the doping phase, and even intermittent doping regimes.

Features of the joint posterior distribution (23) of particular interest are

- (1) The marginal posterior  $\pi(d | Z, Y, x)$  formed by summing (23) over  $t$ . In particular one may require  $\pi(d = 0 | Z)$  to be very small, say less than 1/10 000 for convincing evidence against an athlete. This is the level at which medical investigators and lawyers connected with the project would feel justified in taking legal proceedings.
- (2) The marginal conditional posterior  $\pi(t | Z, Y, x, d > 0)$ , identifying the likely time phase relative to doping of the suspect.

## 5. STRUCTURAL ASSUMPTIONS FOR THE ANALYSIS OF THE GH2000 DATA

### 5.1 Kronecker product

The assumption of a Kronecker product covariance structure was motivated by the multiplicative random effects model of Section 3. For selected pairs of markers it is feasible to compare the 14 by 14 ( $= 2 \times 7$ ) full covariance matrix with that of the Kronecker structure of the  $2 \times 2$  markers covariance by the  $7 \times 7$  times covariance, estimated iteratively by maximum likelihood. The agreement was good. This is not so feasible for the full  $8 \times 7 = 56 \times 56$  matrix due to sheer complexity and the number of variables being greater than the number of observations. As one alternative check we calculated the maximum likelihood estimates of the Kronecker product covariance structure for the GH2000 data. Assuming this structure and multivariate normality, we could generate 20 new sets of data and calculate the eigenvalues of their sample covariance structures. These ordered eigenvalue profiles together with the actual profile are given in Figure 4.

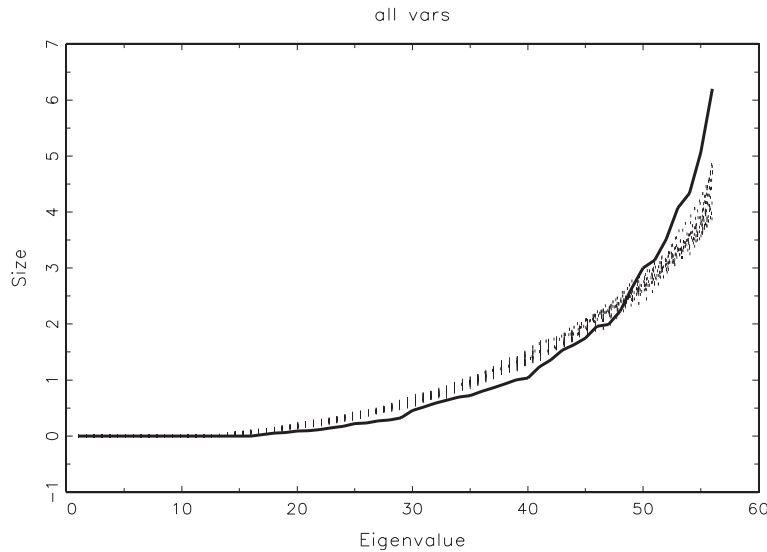


Fig. 4. 20 simulated eigenvalue profiles from a true Kronecker covariance structure with the observed profile.

The observed eigenvalue profile lies outside the envelope of the 20 simulated ones so that the Kronecker product assumption is some way from perfect, with an over-representation of larger eigenvalues, and so warrants the cautious assumption of a Kronecker covariance at the secondary rather than primary level, the parameter  $\delta$  of the inverse Wishart accounting for variability about the Kronecker hyperparameter  $Q$  in (4).

### 5.2 Prior settings

We have assumed vague prior distributions for the location constants,  $\alpha_d, d = 0, \dots, 2$ . The prior for  $\Omega$  in (4) has two hyperparameters,  $\delta$  and  $Q$ . The scalar  $\delta$  is fixed at a grid of values so as to form a posterior distribution for this. Implicitly, this has been assumed to have a constant vague prior of the positive real line. The scale matrix hyperparameter  $Q$  is structured as  $\Gamma \otimes \Sigma$ , as discussed in Section 3. The hyperparameter matrices are estimated by maximum integrated likelihood via the EM algorithm discussed in Section 4.2, implicitly assuming vague third-stage priors for these (Smith, 1973).

The prior probabilities for prediction are taken to be equally likely on the three groups and equally likely on time points 21, 28, 30, 33, 42, 84 days, that is omitting day 0 since this provides baseline measurements for all individuals.

## 6. RESULTS

The type II profile log likelihood of  $\nu$  is plotted in Figure 5. This is formed by choosing a grid of  $\nu$  values and recording the marginal likelihoods maximized over  $\Gamma$  and  $\Sigma$  by the EM algorithm. Changes in  $\nu$  which result in changes in the log likelihood of 2 units from the maximized value of  $\nu = 114$  roughly correspond to a generalized likelihood ratio 95% confidence interval, that is  $\nu$  in (98, 137). This approximates a Bayesian 95% credibility interval for  $\nu$  if the maximized values of  $\Gamma$  and  $\Sigma$  change little over this range, which appears to be the case.

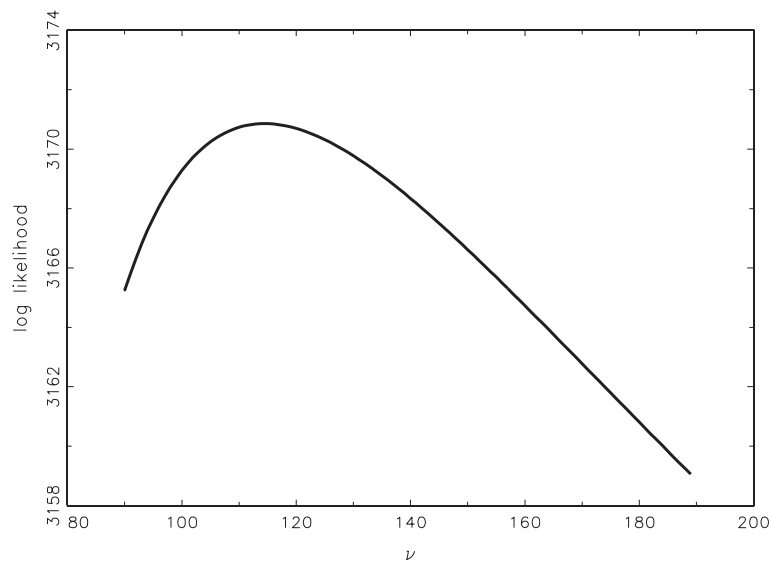


Fig. 5. Profile of the type II likelihood of  $\nu$  at the maximized values of  $\Gamma$ ,  $\Sigma$ .

We summarize the results for prediction of nine validation males from the training data of 43 males with complete records. For simplicity of presentation we amalgamate the two treated groups which would be achieved by adding their predictive probabilities, or equivalently we simply give the probabilities on the placebo group for each individual. The probabilities are predictive probabilities from Section 4.3 for placebo (versus treated) tabulated according to day of sampling, given just the information on that individual at the sampled time point as well as the training data. For each of the nine individuals in the validation sample these placebo probabilities are given for  $\nu = 114, 98, 137$  by row in that order, the top row corresponding to the maximized  $\nu$ . These are given in Table 1. The coding of the nine subjects has as its first digit the country of the participating individual, 1 = Denmark, 2 = Italy, 3 = Sweden, 4 = UK; thus 302 stands for the second individual in Sweden, and we see from the second column that he had a double dose of GH in the double-blind trial.

All those individuals truly on placebo have probabilities of order 1 of being drug-free. Those on dose level 1 have small probabilities of being drug-free when they are submitted at days 21, 28 or 30 and even at day 42 for individual 301. The one individual on double dose is detectable up to day 33 at a threshold of 1 in 10 000.

These probabilities do not vary much over the likely values of  $\nu$  and in this sense the method is robust to the maximum integrated likelihood plug-in.

## 7. DISCUSSION

There are various possible elaborations of the model. The use of  $\Omega_d$  rather than plain  $\Omega$  in model (1) may be appealing. This could feed through to the hyperparameter  $\Sigma$ , the scale matrix for variation amongst markers. What is to be gained may best be judged by a pseudo-Bayes compromise such as regularized discrimination given in Friedman (1989) and underpinned by Brown *et al.* (1999a). Seeming differences in group covariances are often not useful due to the introduction of too many parameters. The Friedman approach uses a hyperparameter to compromise between  $\hat{\Omega}_d$  and  $\hat{\Omega}$  and has a further hyperparameter for

Table 1. Validation sample of nine males, probabilities on placebo for  $v = 114$ , 98(L), 137(U) by row

Subject	Group	Day						
		0	21	28	30	33	42	84
102	1	0.98	$2 \times 10^{-4}$	0.001	0.003	0.05	0.06	–
L	1	0.98	$2 \times 10^{-4}$	0.001	0.004	0.07	0.07	–
U	1	0.97	$1 \times 10^{-4}$	$8 \times 10^{-4}$	0.002	0.04	0.04	–
105	0	1.0	0.97	0.99	0.87	0.66	0.40	–
L	0	1.0	0.97	0.98	0.85	0.64	0.41	–
U	0	1.0	0.97	0.99	0.88	0.68	0.40	–
106	0	1.0	1.0	0.99	–	0.96	0.99	–
L	0	1.0	1.0	0.99	–	0.95	0.98	–
U	0	1.0	1.0	0.99	–	0.96	0.99	–
115	1	0.90	$6 \times 10^{-5}$	$3 \times 10^{-4}$	$6 \times 10^{-4}$	–	0.021	0.24
L	1	0.90	$1 \times 10^{-4}$	$6 \times 10^{-4}$	0.001	–	0.025	0.26
U	1	0.90	$3 \times 10^{-5}$	$1 \times 10^{-4}$	$3 \times 10^{-4}$	–	0.016	0.22
301	1	1.0	$9 \times 10^{-4}$	$3 \times 10^{-6}$	–	$3 \times 10^{-5}$	$5 \times 10^{-4}$	0.16
L	1	1.0	0.001	$4 \times 10^{-6}$	–	$4 \times 10^{-5}$	$6 \times 10^{-4}$	0.16
U	1	1.0	$6 \times 10^{-4}$	$1 \times 10^{-6}$	–	$2 \times 10^{-5}$	$4 \times 10^{-4}$	0.16
302	2	0.71	$3 \times 10^{-7}$	$2 \times 10^{-6}$	–	$7 \times 10^{-5}$	0.001	0.51
L	2	0.67	$4 \times 10^{-7}$	$4 \times 10^{-6}$	–	$9 \times 10^{-5}$	0.001	0.48
U	2	0.75	$2 \times 10^{-7}$	$2 \times 10^{-6}$	–	$5 \times 10^{-5}$	0.001	0.55
310	1	1.0	0.029	$1 \times 10^{-4}$	–	0.46	0.85	0.94
L	1	1.0	0.031	0.002	–	0.47	0.84	0.92
U	1	1.0	0.027	0.001	–	0.46	0.86	0.95
402	0	1.0	1.0	1.0	–	–	1.0	1.0
L	0	1.0	1.0	1.0	–	–	1.0	1.0
U	0	1.0	1.0	1.0	–	–	1.0	1.0
405	1	1.0	$2 \times 10^{-3}$	$8 \times 10^{-5}$	–	$6 \times 10^{-4}$	0.014	0.58
L	1	1.0	$3 \times 10^{-3}$	$1 \times 10^{-4}$	–	$9 \times 10^{-4}$	0.017	0.60
U	1	1.0	0.001	$4 \times 10^{-5}$	–	$4 \times 10^{-4}$	0.012	0.55

regularization.

We have adopted a type II likelihood approach to estimation of the level II hyperparameters  $\Gamma$ ,  $\Sigma$  given by the Kronecker structure in (14). The kernel of the log-integrated likelihood is given by (15). Focusing on one of these hyperparameters, say  $\Sigma$ , the log-integrated likelihood conditional on  $\Gamma$  is given by (16). This is of the form of the log of an inverted Wishart density (Brown, 1993, Appendix A). Thus, had we assumed inverted Wishart prior distributions for these then the posterior distribution would also be inverted Wishart conditional on the other matrix of the pair. This would allow more extensive inference via MCMC in the form of Gibbs sampling.

We have presented a multivariate model for discrimination using longitudinal learning data where prediction is at a single time point. A number of general issues arise to be investigated in future work:

- selection of subsets of responses;
- incorporation of cases with missing information.

Preliminary work looking at subsets of markers using a scoring rule for prediction suggests that subsets do worse than the full eight markers. An algorithm for incorporation of simple structures of missing

information in the case of simple (level I) maximum likelihood, that is when  $\delta \rightarrow \infty$  in our formulation, is given in Brown *et al.* (1999b).

#### ACKNOWLEDGEMENTS

This work was supported by the International Olympic Committee and the European Union under the Biomed 2 Programme, Contract BMH4 CT950678. We are also grateful to Professor Peter Sönksen, the coordinator and driving force of the multicentre project, the seven partners, researchers in Denmark, Italy, Sweden and the UK, and the volunteers who provided the data.

#### APPENDIX A

##### *Matrix-variate distributions*

We shall follow the notation introduced by Dawid (1981) for matrix-variate distributions. This has the advantage of preserving the matrix structures without the need to string by row or column as a vector. It redefines the degrees of freedom as shape parameters for both inverse Wishart and matrix-variate  $T$ , to allow notational invariance under marginalization and very easy symbolic Bayesian manipulations.

With  $U$  a matrix having independent standard normal entries,  $M + \mathcal{N}(\Gamma, \Sigma)$  will stand for a matrix-variate normal distribution of  $V = M + A'UB$  where  $M, A, B$  are fixed matrices satisfying  $A'A = \Gamma$  and  $B'B = \Sigma$ . Thus  $M$  is the matrix mean of  $V$  and  $\gamma_{ii}\Sigma$  and  $\sigma_{jj}\Gamma$  are the covariance matrices of the  $i$ th row and  $j$ th column, respectively, of  $V$ . There is thus a lack of uniqueness in  $\Gamma, \Sigma$ . In this Kronecker product structure they are unique up to an arbitrary scalar multiplier of  $\Gamma$  (which then divides  $\Sigma$ ).

If  $U$  is of order  $n \times p$  with  $n \geq p$ , the notation  $\mathcal{IW}(\delta; \Sigma)$  with  $\delta = n - p + 1$ , will stand for the distribution of  $B'(U'U)^{-1}B$ , an inverse Wishart distribution. The shape parameter  $\delta$  differs from the more conventional degrees of freedom, and may be generalized, using the density function, to take on any positive real value. The matrix-variate  $T$  distribution  $M + \mathcal{T}(\delta; \Gamma, Q)$  is the distribution of  $T$  where  $T$  follows the  $M + \mathcal{N}(\Gamma, \Sigma)$  distribution conditional on  $\Sigma$ , and  $\Sigma \sim \mathcal{IW}(\delta; Q)$ . Corresponding probability density functions are given in Brown (1993), Appendix A.

#### APPENDIX B

##### *Bayes integration*

The prior probability density function of  $(\alpha_d, d = 0, \dots, g, B)$  given  $\Omega$  is the product of

$$\pi(\alpha_d|\Omega) \propto h^{-q/2}|\Omega|^{-1/2} \exp\left\{-\frac{1}{2h}(\alpha - \alpha_{d0})'\Omega^{-1}(\alpha - \alpha_{d0})\right\} \quad (1)$$

for  $d = 0, \dots, g$  and

$$\pi(B|\Omega) \propto |H|^{-qp/2}|\Omega|^{-c/2} \exp[-\frac{1}{2} \text{trace}\{H^{-1}(B - B_0)\Omega^{-1}(B - B_0)'\}]. \quad (2)$$

We first seek to integrate over  $(\alpha_d, d = 0, \dots, g, B)$  for given  $\Omega$ . In this Gaussian setting, to do this we should 'complete the square' in  $\alpha_d$  and  $B$  within the exponentiated terms of likelihood times prior. First focusing on the likelihood for the  $g + 1$  sets of data comprising (5), the exponential term is

$$-\frac{1}{2} \text{trace} \Omega^{-1} \sum_{d=0}^g \{(Y_d - X_d B)'\Omega^{-1}(Y_d - X_d B) - 2(Y_d - X_d B)'\Omega^{-1}\alpha'_d + n_d \alpha'_d \Omega^{-1} \alpha'_d\}, \quad (3)$$

using  $\text{trace}(AC) = \text{trace}(CA)$ .

The  $\alpha_d$ -term of (3) combines the prior (1) to give

$$|\Omega|^{-1/2} \exp\{-\frac{1}{2} \text{trace } \Omega^{-1} \sum \{-2(Y_d - X_d B)' 1\alpha'_d + n_d \alpha_d \alpha'_d + (\alpha_d - \alpha_{d0})(\alpha_d - \alpha_{d0})'/h\}\}. \quad (4)$$

Letting  $h \rightarrow \infty$ , the prior is vague and completing the square we see that given  $B, \Omega, Y$  a posteriori

$$\alpha'_d - (\bar{Y}_d - \bar{X}_d B) \sim \mathcal{N}(1/n_d, \Omega) \quad (5)$$

independently for  $d = 0, \dots, g$ .

Turning to the integration of  $B$  for given  $\Omega$ , the first term of (3) minus the quadratic in  $B$  left from completion of the square in  $\alpha$  is

$$\sum \{(Y_d - X_d B)'(Y_d - X_d B) - n_d(\bar{Y}_d - \bar{X}_d B)'(\bar{Y}_d - \bar{X}_d B)\}.$$

For weak prior knowledge, letting  $H$  become large in (2), and completing the square in  $B$ , this gives a posteriori given  $\Omega, Y$

$$B - \hat{B} \sim \mathcal{N}\left(\left[\sum X_d^{*'} X_d^*\right]^{-1}, \Omega\right) \quad (6)$$

where  $X_d^* = X_d - 1\bar{X}_d$  and

$$\hat{B} = \left[\sum X_d^{*'} X_d^*\right]^{-1} \left(\sum X_d^{*'} Y_d\right).$$

This completion of the square leaves the likelihood marginalized over the vague priors  $\alpha_d, d = 0, \dots, g, B$  for given  $\Omega$  as proportional to

$$\{|\Omega|^{-N/2} \exp\{-\frac{1}{2} \text{trace } (\Omega^{-1} S)\}\} \quad (7)$$

with centred  $Y_d^* = Y_d - 1\bar{Y}_d$  and sum of products of residuals

$$S = \sum_d (Y_d^* - X_d^* \hat{B})'(Y_d^* - X_d^* \hat{B}). \quad (8)$$

REFERENCES

BROWN, P. J. (1993). *Measurement, Regression, and Calibration*. Oxford: Clarendon.

BROWN, P. J., FEARN, T. AND HAQUE, M. S. (1999a). Discrimination with many variables. *Journal of the American Statistical Association* **94**, 1320–1329.

BROWN, P. J., KENWARD, M. G. AND BASSETT, E. E. (1999b). Discrimination with longitudinal data. Technical Report UKC/IMS/99/09. University of Kent.

BROWN, P. J., LE, N. D. AND ZIDEK, J. V. (1994). Multivariate spatial interpolation and exposure to air pollutants. *Canadian Journal of Statistics* **22**, 489–509.

CHEN, CHAN-FU (1979). Bayesian inference for a normal dispersion matrix and its application to stochastic multiple regression analysis. *Journal of the Royal Statistical Society B* **41**, 235–248.

DAWID, A. P. (1981). Some matrix-variate distribution theory: notational considerations and a Bayesian application. *Biometrika* **8**, 265–274.

DEMPSTER, A. P., LAIRD, N. M. AND RUBIN, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm (with discussion). *Journal of the Royal Statistical Society B* **39**, 1–38.

- FRIEDMAN, J. H. (1989). Regularised discriminant analysis. *Journal of the American Statistical Association* **84**, 165–175.
- GEISSER, S. (1964). Posterior odds for multivariate normal classifications. *Journal of the Royal Statistical Society B* **26**, 69–76.
- GOOD, I. J. (1965). *Estimation of Probabilities: an Essay on Modern Bayesian Methods*. Cambridge, MA: MIT Press.
- KRZANOWSKI, W. J. AND MARRIOTT, F. H. C. (1995). *Multivariate Analysis, Part 2: Classification, Covariance Structures and Repeated Measurements. (Kendall's Library of Statistics)*. London: Arnold.
- LOGAN, T. P. AND GUPTA, A. K. (1993). Bayesian discrimination using multiple observations. *Communications in Statistics: Theory and Methodology* **22**, 1735–1754.
- MARDIA, K. V., KENT, J. T. AND BIBBY, J. M. (1979). *Multivariate Analysis*. London: Academic.
- MCLACHLAN, G. J. (1992). *Discriminant Analysis and Statistical Pattern Recognition*. New York: Wiley.
- MENG, X. L. AND RUBIN, D. B. (1993). Maximum likelihood estimation via the ECM algorithm. *Biometrika* **80**, 267–278.
- SMITH, A. F. M. A general Bayesian linear model. *Journal of the Royal Statistical Society B* **35**, 67–75.

[Received 5 August, 2000; first revision 27 October, 2000; second revision 25 January, 2001; accepted for publication 5 February, 2001]