

THE GENETICAL ANALYSIS OF COVARIANCE STRUCTURE

N. G. MARTIN and L. J. EAVES

Department of Genetics, University of Birmingham, Birmingham B15 2TT, England

Received 9.ix.76

SUMMARY

The analysis of covariance structures (Jöreskog, 1973) is adapted to the simultaneous maximum likelihood estimation of genetical and environmental factor loadings and specific variances. The goodness of fit is tested by chi square and standard errors of parameter estimates can be obtained.

Any linear model used in univariate genetical analyses can be extended to the multivariate case. Most biological hypotheses about the relationships between variables can be specified by a variety of factor models. Individual parameters can be given fixed values or set to zero and hypotheses concerning the congruence of genetical and environmental correlations can be tested.

The method is illustrated with published twin data on cognitive abilities.

1. INTRODUCTION

THE techniques of factor analysis have been used extensively in the behavioural sciences to simplify the representation of the relationships between multiple variables. Geneticists, rightly, are sceptical about the use of such methods in genetical research. There are several reasons for this. Firstly, factor analysis has generally been confined to the analysis of phenotypic variation and the interpretation in simple genetical or environmental terms of factors defined in this way is not legitimate. Secondly, factor analysis has usually, though not without exception, been an exploratory device. Consequently, factor-analytic studies have often consisted of a *post-hoc* examination of tables of factor loadings and studies designed in advance to test particular hypotheses about the structure of covariation have been atypical. This approach does not commend itself to geneticists for whom predictive statements, not rationalisations, are the ideal of progressive research. Thirdly, in spite of the general development of the statistical and computational methods appropriate to factor analysis towards efficient estimation and tests of significance, most factor studies remain content with techniques which yield estimates of unknown properties and even the crudest statistical tests are seldom attempted.

Any problems of analysing phenotypic covariation apply, *a fortiori*, to attempts to seek a factorial representation of genetical and environmental covariation separately. The idea that genetical and environmental components of covariation might be separated by a technique analogous to factor analysis was formulated long ago in the context of behavioural genetics (*e.g.* Loehlin and Vandenberg, 1968) but there have been few significant advances in the practical problem of providing a general approach to estimation and hypothesis testing in this context.

The main difficulty is that the experimental designs necessary for resolving the different sources of gene action and environmental influence are not appropriate for factor analysis because the information required

to estimate the components of a genotype-environmental model usually come from more than one type of family grouping, *e.g.* identical and fraternal twins. Although there have been many suggestions for approximations, *e.g.* using components estimated from variation within pairs of identical and fraternal twins, such approaches suffer from several important deficiencies. Foremost among these is the fact that such estimates may well be biased by effects for the existence of which no statistical test is conducted. Furthermore, such estimates are inefficient because they ignore information contained, for example, in the variation between families. At best, most existing methods throw away data and lead to estimates which are unnecessarily imprecise. At worst, such methods are actually misleading because they do not provide any test of the assumptions which are made in deriving the estimates.

Although several authors (*e.g.* Jinks and Fulker, 1970; Eaves, 1975; Eaves and Eysenck, 1975) have stressed the inadequacy of the *ad hoc* approach characteristic of behavioural genetics in the univariate case, such remains the approach of most behavioural geneticists to their data. Unfortunately, all the problems of inefficient and misleading univariate methods are exacerbated by their importation into the domain of multivariate genetical analysis because crude methods of genetical analysis have been combined with no less crude methods of factor analysis. Eaves and Gale (1974) attempted to redress the balance somewhat by focusing the attention of a multivariate genetical analysis on the genotype-environmental model which was appropriate to the data before attempting to explain the genetic covariation in factorial terms. They suggested analysing the data in two stages; the first to determine the genetical and environmental model appropriate to the multivariate data without reference to the factor structure and the second attempting to test a simple factor model for the genetical covariation.

Their approach was based on weighted least squares, and therefore, allowed various tests of hypotheses (including a test of goodness of fit) to be conducted. However, their approach shared a common failing with much of the previous work since the attempt at factor analysis was confined to components of variance and covariance which need not have formed a positive definite covariance matrix. It seems likely, therefore, that this approach would not always yield satisfactory parameter estimates. Secondly, as with every other approach, the method involved obtaining a genetical covariance matrix first and then attempting factorisation. In our view the structure of the trait covariation should be no less a matter of advance specification than the basic formulation of the biometrical-genetical model. Thus, the psychological constructs implicit in the choice of variables should be tested at the same time as the causal model being examined in the particular constellation of relatives chosen to be included in the study. We thus need an approach which allows us to formulate both our biometrical-genetical model and our model for the structure of the trait covariation in advance of the analysis, to provide parameter estimates which make the best use of the available data and, wherever possible, facilitate those tests of significance which form the basis of further research.

Generally, geneticists are familiar with the methods of model-fitting but less familiar with the principles of factor analysis. It will help the development of our general model if we outline some of these here, for the classical

case of the analysis of phenotypic variation and covariation of v correlated variables. Our initial data summary is simply the $v \times v$ symmetric, positive definite covariance matrix S . Factor analysis seeks to represent the variation in each variable in terms of a contribution which is specific to that variable (the "specific variance") and contributions from $f < v$ sources (the "common factors") which contribute to variation in more than one variable at the same time. The contribution of each common factor to a particular variable may be assessed from the factor loading which can be standardised to represent the correlation between the variable in question and the hypothesised factor. In a behavioural context the scales of measurement are often arbitrary so psychologists usually prefer to work with correlation matrices rather than covariance matrices. We shall consider only the analysis of covariance matrices, however, because no appropriate scaling of the genetical and environmental components can be suggested in advance of the analysis of the causes of variation. The processes of scaling and estimation are interdependent.

The factors may be constrained to be uncorrelated (orthogonal factors) or a solution may be sought in which the factors are themselves assumed to be correlated ("oblique" factors). For a covariance matrix S , we may represent the covariance matrix predicted in terms of the factor model thus:

$$\Sigma = \Lambda\Phi\Lambda' + \Psi^2$$

We are using, as far as possible, the notation of Jöreskog (1973) whose contributions to factor analysis will form a major component in the formulation of our approach. Λ is a $v \times f$ matrix of factor loadings. The square $f \times f$ matrix Φ contains the correlations between the hypothesised factors and Ψ^2 is a diagonal matrix containing the v specific variances.

Formulated in this way the factor model is indeterminate. An infinite number of solutions for Λ and Φ may be found which give the same Σ . Constraints have to be introduced to ensure a unique solution. Normally these constraints are arbitrary, stemming from mathematical necessity rather than from any psychological considerations about the anticipated form of the factors. Most approaches begin by insisting that the factors be orthogonal (*i.e.* that $\Phi = I$). This is insufficient to ensure a unique solution so further constraints are required. The classical approach to maximum-likelihood factor analysis requires that the loadings satisfy the constraints that the matrix $\Lambda'\Psi^{-2}\Lambda$ be diagonal. Because these and other constraints imposed in obtaining a suitable definition of the factors are arbitrary, factor analysts have usually attempted to improve the interpretation of factors by "rotation". This involves transforming the loadings according to a criterion which might be a better reflection of the investigator's expectations for the outcome of his study. Thurstone suggested (*e.g.* 1945) that variables might fall into groups having high loadings on relatively few factors and small or zero loadings on the rest. This was the notion of "simple structure". Several numerical criteria have been suggested which might be used as a basis for rotation to simple structure. In the course of rotation an investigator may choose to retain the criterion of orthogonality, or he may prefer to allow the factors to become oblique if the interpretation is clarified by doing so. Although significance tests may have been used in determining the number of common factors and specifics (and even this is still the

exception rather than the rule in practical applications) the choice of final rotated solution has rarely rested upon any statistical test of the model. In a psychological context Jöreskog (1973) has indicated that the investigator who is bold enough to commit himself in advance to a precise model for the covariance structure (specifying, for example, which loadings are expected to take zero values) is in a position to proceed directly to the maximum-likelihood solution providing he is able to specify sufficient constraints on the parameters to ensure a determinate solution. Indeed, such hypotheses are much stronger than those which merely involve a guess about the likely number of common factors. In addition, it is possible to construct a statistical test of the worthwhileness of relaxing the constraint of orthogonality upon a set of loadings.

The work of Jöreskog has been developed in many directions, for example to the comparison of factor structure in several populations, but to our knowledge no biometrical-genetical application of the approach has yet been devised. This paper is intended to show how the method can be adapted to the type of problem with which geneticists are often concerned. We hope that this will form the basis of a definitive method for approaching multivariate problems in biometrical genetics. As far as we can see, with a little ingenuity, virtually any model which geneticists might formulate for the relationships between variables could be tested using the approach we outline. Thus, although our own primary interest is in behavioural genetics, we feel that the approach may be worth considering by anyone concerned with the genetical analysis of covariation.

Although we have chosen to illustrate the method with twin data relating to a small set of human abilities we should not be misled into believing that the method is specific to experimental designs involving only twins or to variables for which a single common factor constitutes the simplest explanation. The method is general. It will be most valuable when data are available for many kinds of family grouping and when the variables studied are carefully chosen for biological or psychological reasons.

Our treatment may be outlined as follows:

- i. The general form, following Jöreskog (1973) of the covariance structure of a single population.
- ii. The adaptation of this model for the analysis of genotypic and environmental covariance.
- iii. The numerical and statistical procedure for estimating the parameters and for testing the model.
- iv. The method for obtaining the covariances of the estimates.
- v. The exposition of our treatment on published twin data relating to human abilities.

Our aim is to provide an adequate method within which all the factors of biological interest can be assessed: the components of gene action; the influence of the mating system; the sources of environmental variation and genotype-environmental covariance. All these aspects of individual variation may come within the scope of a multivariate analysis given that the components of a particular model are estimable with the data available. The approach allows us to estimate such effects by methods which are statistically efficient and enables us to decide which of them can be regarded as beyond resolution with a given study.

2. THE GENERAL MODEL

Following Jöreskog (1973) we consider a single observed covariance matrix, \mathbf{S} , ($v \times v$) of v variables and having N d.f. The expected covariance matrix Σ ($v \times v$), is a function of various parameter matrices and has the form:

$$\Sigma = \mathbf{B}(\Lambda\Phi\Lambda' + \Psi^2)\mathbf{B}' + \Theta^2$$

For our purpose, Λ ($v \times f$) is a matrix of factor loadings of the v variables on f factors. Φ ($f \times f$) is the (positive definite) matrix of correlations between the factors. Ψ^2 ($v \times v$) and Θ^2 ($v \times v$) are diagonal matrices containing the specific variances.

The matrix \mathbf{B} allows the factors to be scaled and combined in a variety of ways. We shall discuss its application in our context shortly. Jöreskog has illustrated the application of this model to a variety of problems. In our case however, some extension is needed because the data we shall be analysing will typically consist not just of one covariance matrix, but several. These will be the covariance matrices obtained from the analysis of many different groups of relatives. The classical twin study, for example, will yield covariance matrices between and within pairs for identical and non-identical twins, making four in all.

In the univariate case a genetical analysis attempts to partition the phenotypic variation into contributions from several sources. Usually, the contributions are written in the form of a linear model in which the i th observed variance (or mean square) from a family study is predicted thus:

$$\sigma_i^2 = \sum_{j=1}^p c_{ij}x_j$$

c_{ij} is the coefficient of the j th parameter of the genotype-environmental model (the j th "source" x_j) in the expectation of σ_i^2 . Most of the coefficients can be obtained from genetical theory.

For the multivariate case in which p sources contribute to the variance and covariance of v variables we have, generally,

$$\Sigma_i = \sum_{j=1}^p c_{ij}[\mathbf{B}_j(\Lambda_j\Phi_j\Lambda_j' + \Psi_j^2)\mathbf{B}_j' + \Theta_j^2]$$

Thus corresponding to each source of variation in the univariate case (*e.g.* cultural, additive genetical, non-additive genetical etc.) there are corresponding parameter matrices \mathbf{B}_j , Λ_j , Φ_j , Ψ_j and Θ_j in the multivariate case. In our application \mathbf{B}_j will be used to permit factor loadings contributing to one source (*e.g.* genetical covariation) to be simple scalar functions of those arising due, for example, to environmental differences. We are thus able to test hypotheses about the similarity of genetical and environmental factors in this approach. For this reason, and because we will usually have no reason to scale genetical and environmental specifics relative to one another, we shall usually put $\Psi_j^2 = \mathbf{0}$ and denote the various specific sources of variation by Θ_j^2 .

3. A SPECIFIC FORMULATION OF A MODEL FOR TWIN DATA

As a simple example we may consider how we might formulate a model for the covariance structure of data on identical and fraternal twins on the

assumptions that gene action is additive, mating is random and that members of the same family are not made more alike by shared environmental experiences. We stress that this is only an example, albeit of a model which appears to be appropriate for individual differences for certain dimensions of human personality (Eaves and Eysenck, 1975, 1976, 1977). Ideally we would like more extensive data on other kinds of family, including fostered individuals and further ancestral and collateral relationships. Such studies might expose additional sources of variation which are not subjected to very powerful scrutiny in the usual twin study. Eaves and Eysenck (*e.g.* 1975) have discussed the assumptions in more detail and considered the univariate consequences of their failure.

For a single variable we may write our expectation for the total phenotypic variance in terms of our hypothesised model thus:

$$\sigma_p^2 = \frac{1}{2}D_R + E_1$$

Mather and Jinks (1971) define D_R in terms of the effects and frequencies of the alleles at many loci. It represents the additive portion of the genetical variance. The contribution of environmental differences specific to individuals and not shared by members of the same family is denoted by E_1 . This, we emphasise, is the simplest possible model for the joint action of genetical and environmental influences. The model can be complicated in a variety of ways to include, for example, the contribution of non-additive genetical effects such as dominance, assortative mating, shared environmental effects (E_2), genotype-environmental interaction and genotype-environmental covariance (see *e.g.* Jinks and Fulker, 1970; Eaves, 1975; 1976*a, b*; Eaves *et al.* 1977). However, before we attempt explanations in terms of more intricate effects we must satisfy ourselves by the appropriate scaling tests that a simpler model for individual differences is unacceptable.

This simple model can be generalised to the multivariate case. We may write for the phenotypic covariance matrix, given the above assumptions about gene action, mating system, and environmental causes:

$$\Sigma_p = \frac{1}{2}(\Lambda\Lambda' + D^2) + HH' + E^2$$

We use Λ to represent the loadings of the variables on the additive genetical factors and H to represent the loadings of the same variables on the environmental factors. D^2 and E^2 respectively are used for the contributions of genetical and environmental influences *specific* to the particular variables. We may define the genotypic covariance matrix (given our assumptions of additivity and random mating) as $\frac{1}{2}(\Lambda\Lambda' + D^2)$. Similarly, the environmental covariance matrix is $HH' + E^2$. Thus, provided we can estimate Λ , H , D and E , we are able to obtain genetic correlations and a variety of other summary statistics. These will be illustrated later.

If we only have the phenotypic covariance matrix there is no possibility of any analysis into genetical and environmental factors. Such analysis requires data on more than one kind of family grouping. In the case of identical and fraternal twins the data may be summarised in terms of four matrices. For each of the two different groups of twins we have two distinct matrices of mean products, between and within pairs, which can easily be obtained from a multivariate analysis of variance of the twin pairs. The expected contributions of the genetical and environmental factors to the different matrices can be obtained from genetical theory and are given

below in the univariate case of the four mean squares from a classical twin study for D_R , E_1 and also E_2 , the between families environmental component which we are not considering in the present simple model:

	D_R	E_1	E_2
MS_{MZB}	1	1	2
MS_{MZW}	0	1	0
MS_{DZB}	$\frac{3}{4}$	1	2
MS_{DZW}	$\frac{1}{4}$	1	0

The fact that the contributions of the different sources are not the same for all the raw matrices provides the basis for estimation of the separate genetical and environmental factors.

For the mean squares and mean products derived from pairs of twins we have, by analogy with the univariate case for our simplest model with only two sources of variation (see *e.g.* Eaves and Eysenck, 1977)

$$\begin{aligned}\Sigma_{BMZ} &= \Delta\Delta' + D^2 + HH' + E^2 \\ \Sigma_{WMZ} &= HH' + E^2 \\ \Sigma_{BDZ} &= \frac{3}{4}(\Delta\Delta' + D^2) + HH' + E^2 \\ \Sigma_{WDZ} &= \frac{1}{4}(\Delta\Delta' + D^2) + HH' + E^2\end{aligned}$$

The matrices Σ_{BMZ} , Σ_{WMZ} , Σ_{BDZ} , Σ_{WDZ} represent the mean products between (B) and within (W) pairs of monozygotic (MZ) and dizygotic (DZ) twins respectively. Obviously the model can be easily extended to include the statistics derived from other family groupings, or to incorporate further subtleties of genetic and environmental action and interaction if the data should warrant such complications.

We may now relate the terms of our model for twin data to the general model above. We have four matrices ($m = 4$) and our model involves two sources of variation in the univariate case ($p = 2$). We may now put

$$\begin{aligned}\Lambda_1 &= \Delta; && \text{the matrix of additive genetical loadings} \\ \Lambda_2 &= H; && \text{the matrix of within family environmental loadings} \\ \Theta_1^2 &= D^2; && \text{the diagonal matrix of specific additive genetical variances} \\ \Theta_2^2 &= E^2; && \text{the diagonal matrix of specific within family environmental variances.}\end{aligned}$$

We shall set Ψ_1^2 and Ψ_2^2 to zero. Indeed, if we have sufficient reason (either because of our advance expectations, or because certain parameters have been shown to be non-significant) we may set certain of the loadings to zero at this stage. If we wish, we can introduce matrices of factor correlations, Φ_1 and Φ_2 to denote the correlations of genetical and environmental correlations. These may be fixed to be identity matrices for the orthogonal solution, or the off-diagonal elements may be fixed to specify a desired degree of obliquity, or allowed to vary (within the constraints that the Φ 's should be positive definite) to obtain the oblique solution which yields the best fit. The advantages of fixing certain loadings in advance for the multiple factor solution are that we can test directly particular hypotheses

about the factor structure. Our coefficients (c_{ij}) are the coefficients of the genetical and environmental factors in the expectations of the four matrices. Thus, c_{31} is the coefficient of $(\Delta\Delta^1 + D^2)$ in the expectation of Σ_{BDZ} for example, (i.e. $\frac{3}{4}$) and c_{22} is the coefficient of $(HH' + E^2)$ in the expectation of Σ_{WMZ} . It is quite likely that we shall want to test the hypothesis that the genetical loadings (for example) are simply scaled versions of the environmental loadings. This would imply that the genetical and environmental structures are identical, apart from specific factors, and that genetical and environmental factors are affecting the same aspects of the organism in a consistent manner. Thus, to incorporate such a constraint in our model we simply insist that $\Delta = bH$ by defining $\Lambda_1 = B_1 \Lambda_2$ where $B_1 = bI$ and writing $B_2 = I$ in the expectations above.

4. ESTIMATING THE PARAMETERS AND TESTING THE MODEL

In Jöreskog's treatment the maximum likelihood estimates of the parameters of the model are obtained by maximising the log-likelihood numerically, which is given by

$$\log L = -\frac{1}{2}N[\log |\Sigma| + \text{tr}(\Sigma\mathbf{S}^{-1})]$$

The constant term is omitted. Maximising $\log L$ is equivalent to minimising

$$F = \frac{1}{2}[\log |\Sigma| + \text{tr}(\Sigma\mathbf{S}^{-1})]$$

In our problem, given that the observations are multivariate normal, we may write the log-likelihood of obtaining the m observed independent mean products matrices (S_i) as

$$\log L = -\frac{1}{2} \sum_{i=1}^m N_i[\log |\Sigma_i| + \text{tr}(S_i\Sigma_i^{-1})]$$

(once again, omitting the constant term).

We actually minimise $-\log L$, or

$$F = \frac{1}{2} \sum_{i=1}^m N_i[\log |\Sigma_i| + \text{tr}(S_i\Sigma_i^{-1})]$$

For a given model we require the parameter estimates which minimise F . There are many ways of doing this, some of which have been implemented by Jöreskog for a variety of different applications of his model. We chose, however, to develop our own program for genetical applications using sub-routines written by the Numerical Algorithms Group (1974). From their many routines for minimisation we chose EO4HAF for constrained minimisation as the basis of our numerical method, in conjunction with several accurate routines written for the matrix operations necessary for function evaluation etc. EO4HAF has the advantage of allowing certain flexibility in the choice of method for minimisation in that the user can specify whether first or second derivatives need be evaluated and whether this should be done numerically by difference or by substitution in coded formulae. In view of the complexity of the functions, in most cases we chose to base the minimisation on the Powell 64 method which did not require any differentiation but relied solely on evaluation of the functions themselves for a variety of parameter values. We found that this method gave satisfactory results for our example problem and for a number of similar problems

still unpublished. EO4HAF uses a penalty function technique (Lootsma, 1972) for constraining estimates to lie within a particular region. Classical factor analysis often ran into difficulties because the maximum likelihood solution required that one or more specific variances take negative values. By conducting the analysis in terms of Θ rather than Θ^2 for the specifics we were able to reduce the number of constraints specified and therefore to minimise the amount of computer time required to find the best solution.

For our relatively simple case of a single common genetical factor and a single environmental factor we were able to avoid the specification of any constraints, though we envisage that these would be necessary to obtain multiple factor solutions, for example to constrain the Φ 's to be positive definite.

In some cases we found that minimisation by the Powell 64 method failed when matrices became singular during the course of minimisation. We managed to get around this problem by using another option of the subroutine EO4HAF which employs numerical approximations to the first derivatives of F with respect to the parameter estimates.

Given that we have obtained the maximum-likelihood estimates of our parameters we may test the hypothesis that a less restricting model does not significantly improve the fit by computing $(2L_0 - L_1)$, which is distributed as χ^2 when L_1 is the log likelihood obtained under the restricted hypothesis (H_1) and L_0 is the log likelihood obtained under the less demanding hypothesis (H_0). The H_0 we shall adopt in practice is one which assumes that as many free parameters are required to explain the data as there are independent mean squares and mean products in the first place *i.e.* $\Sigma_i = S_i$ for every i —a "perfect fit" solution. In this case we have simply

$$L_0 = -\frac{1}{2} \sum_{i=1}^m N_i [\log |S_i| + v]$$

When we have m matrices and v variables, the χ^2 has $\frac{1}{2}mv(v+1) - t$ where t is the total number of free parameters estimated under H_1 .

5. COVARIANCES OF THE ESTIMATES

In order that the parameter estimates of a satisfactory model might be interpreted more rigorously we would like their covariance matrix. This is the inverse of the matrix of second derivatives of the log likelihood with respect to the maximum likelihood estimates of the free parameters. Jöreskog (1973) gives second derivatives of the log likelihood for problems involving models for single covariance matrices. We followed his approach, constructing first the matrix of second derivatives with respect to all the parameters, fixed and free, then striking out the rows and columns corresponding to fixed parameters, and finally combining the information on those parameters which are constrained to be equal.

To evaluate the second derivatives we need to know the form of the first derivatives. For any given source j , of the genetical model, we have parameter matrices B_j , A_j etc. and the first derivatives of F with respect to these are:—

$$\partial F / \partial B_j = \sum_{i=1}^m N_i [c_{ij} \Omega_i B_j (A_j \Phi_j A_j' + \Psi_j^2)]$$

$$\partial F/\partial \Lambda_j = \sum_{i=1}^m N_i [c_{ij} \mathbf{B}'_j \Omega_i \mathbf{B}_j \Lambda_j \Phi_j]$$

$$\partial F/\partial \Phi_j = \sum_{i=1}^m N_i [c_{ij} \frac{1}{2} \Lambda_j \mathbf{B}'_j \Omega_i \mathbf{B}_j \Lambda_j]$$

$$\partial F/\partial \Psi_j = \sum_{i=1}^m N_i [c_{ij} \mathbf{B}'_j \Omega_i \mathbf{B}_j \Psi_j]$$

$$\partial F/\partial \Theta_j = \sum_{i=1}^m N_i [c_{ij} \Omega_i \Theta_j]$$

$$\text{where } \Omega_i = \Sigma_i^{-1} (\Sigma_i - S_i) \Sigma_i^{-1}$$

We now obtain the second derivatives by using the following theorem, proved by Jöreskog (1973):

Let the elements of Σ_i be functions of two parameter matrices $\mathbf{M}_k = (\mu_{kef})$ and $\mathbf{N}_i = (v_{igh})$ and let

$$F(\mathbf{M}_k, \mathbf{N}_i) = \frac{1}{2} \sum_{i=1}^m N_i [\log |\Sigma_i| + \text{tr} (S_i \Sigma_i^{-1})].$$

Then if

$$\partial F/\partial \mathbf{M}_k = \sum_{i=1}^m N_i (\mathbf{U}_i \Omega_i \mathbf{X}) \quad \text{and} \quad \partial F/\partial \mathbf{N}_i = \sum_{i=1}^m N_i (\mathbf{V}_i \Omega_i \mathbf{Z})$$

where \mathbf{U}_i , \mathbf{X} , \mathbf{V}_i , \mathbf{Z} are independent of the S_i and Ω_i is defined above, we have asymptotically

$$E(\partial^2 F/\partial \mu_{kef} \partial v_{igh}) = \sum_{i=1}^m N_i [(\mathbf{U}_i \Sigma_i^{-1} \mathbf{V}_i)_{eg} (\mathbf{X}' \Sigma_i^{-1} \mathbf{Z})_{fh} + (\mathbf{U}_i \Sigma_i^{-1} \mathbf{Z})_{eh} (\mathbf{X}' \Sigma_i^{-1} \mathbf{V}_i)_{fg}]$$

Since \mathbf{U}_i and \mathbf{V}_i both contain the scalar constant c_{ij} from the genotype-environmental model, we shall have to calculate these matrices afresh for each expected covariance matrix Σ_i . However, if we define

$$\mathbf{W} = \frac{1}{c_{ij}} \mathbf{U}_i$$

$$\mathbf{Y} = \frac{1}{c_{ij}} \mathbf{V}_i$$

then we may write

$$\begin{aligned} E(\partial^2 F/\partial \mu_{kef} \partial v_{igh}) \\ = \sum_{i=1}^m N_i c_{ik} c_{il} [(\mathbf{W} \Sigma_i^{-1} \mathbf{Y})_{eg} (\mathbf{X}' \Sigma_i^{-1} \mathbf{Z})_{fh} + (\mathbf{W} \Sigma_i^{-1} \mathbf{Z})_{eh} (\mathbf{X}' \Sigma_i^{-1} \mathbf{Y})_{fg}] \end{aligned}$$

which facilitates computation.

We can now calculate the whole information matrix for all elements of the k sets of matrices \mathbf{B}_j , Λ_j , Φ_j , Ψ_j and Θ_j where there are k parameters in the genetical model. However, it is obvious that much computer time can be saved by not computing those elements

$$E(\partial^2 F/\partial \mu_{kef} \partial v_{igh})$$

where either μ_{kef} or v_{igh} has been fixed. This is equivalent to striking out from the total information matrix those rows and columns corresponding to fixed parameters.

We now wish to reduce the information matrix still further by combining rows and columns of information about the same free parameter. Jöreskog suggests the following procedure. Let the remaining free parameters form a vector $\gamma' = (\gamma_1, \gamma_2, \dots, \gamma_l)$ and among these there are some distinct parameters $\pi_1, \pi_2, \dots, \pi_m$. Let $k_{ig} = 1$ if $\gamma_i = \pi_g$ and $k_{ig} = 0$ otherwise and let $\mathbf{K} = (k_{ig}), i = 1, 2, \dots, l; g = 1, 2, \dots, m$. Then we have

$$\partial F / \partial \pi = \mathbf{K}' \partial F / \partial \gamma$$

and

$$E(\partial^2 F / \partial \pi \partial \pi') = \mathbf{K}' E(\partial^2 F / \partial \gamma \partial \gamma') \mathbf{K}.$$

The elements of the information matrix on the right-hand side are obtained as described above. The reduced information matrix on the left-hand side can be inverted to give the variance-covariance matrix of the free parameter estimates.

The expressions given above for the first and second derivatives could be programmed to assist the minimisation of F by other methods and such options are available with the routine we employ.

6. AN EXAMPLE

To illustrate the method we shall use some twin data given by Loehlin and Vandenberg (1968) for the covariation of five of Thurstone's Primary Mental Abilities—Numerical ability (N), Verbal comprehension (V), Spatial ability (S), Word fluency (W) and Reasoning ability (R). These data have already been used to illustrate the method of Eaves and Gale (1974) who discuss the limitations of the data and of the classical twin study. A comparison with their analysis will make clear the conceptual and analytical advantages of the present approach.

Measurements were made on 123 pairs of MZ twins and 75 pairs of DZ twins, the members of each pair having been raised together. Loehlin and Vandenberg (1968) discuss in detail the structure of their sample and conclude that their MZ and DZ twins can be regarded as sub-samples from the same population. They give the between-pair and within-pair mean products matrices for both MZ and DZ twins as Appendices A-D of their paper. Each 5×5 matrix contains 15 unique statistics, providing a total of 60 d.f.

We shall first see whether we can adequately account for the data with the simple model developed in Section 3 above. This contains factor loadings and specific variances corresponding to only D_R and E_1 . We thus estimate five λ_i 's and five θ_i 's for each of the two sources of variation, being 20 parameters in all. We place no bounds on these during the estimation procedure since the expected values are functions of $\hat{\Lambda}\hat{\Lambda}'$ and $\hat{\Theta}^2$ so sign is immaterial.

The log likelihood on this hypothesis is -7311.14 whereas the log likelihood on the hypothesis of a perfect fit solution is -7274.89 . Twice the difference between the log likelihood values, $2(L_0 - L_1)$ yields a chi-square of 72.50 for 40 d.f. (since we have 60 unique statistics which we have

tried to summarise with reference to 20 independent parameters). This model fails badly ($P < 0.001$). We also tried to fit the corresponding model with factor loadings and specifics corresponding to E_2 and E_1 only *i.e.* omitting genetical factors, but this model gives an even worse failure ($\chi_{40}^2 = 127.50$). The maximum likelihood parameter estimates from these two models are given in tables 1 and 2.

Clearly we must try a more complex model if we are to provide a satisfactory explanation of the data. If we extend the model to include factor loadings and specific variances corresponding to all three sources of variation, E_1 , E_2 and D_R , we now have a total of 30 free parameters. The log likelihood on this model is -7291.39 corresponding to $\chi_{30}^2 = 33.0$, an excellent fit.

TABLE 1

Parameter estimates of the E_1D_R model. The λ_{E_1} 's are loadings on the E_1 factor and θ_{E_1} 's are square roots of the specific E_1 variances

	λ_{E_1}	θ_{E_1}	λ_{D_R}	θ_{D_R}
<i>N</i>	0.528	19.201	49.447	34.617
<i>V</i>	10.265	7.315	36.991	22.114
<i>S</i>	1.210	20.717	32.451	54.677
<i>W</i>	5.751	12.830	23.187	22.507
<i>R</i>	6.614	8.588	25.082	15.163

$$\chi_{40}^2 = 72.5$$

TABLE 2

Parameter estimates of the E_1E_2 model

	λ_{E_1}	θ_{E_1}	λ_{E_2}	θ_{E_2}
<i>N</i>	11.359	25.515	31.115	23.523
<i>V</i>	10.420	10.753	27.227	13.759
<i>S</i>	8.824	24.856	19.384	37.662
<i>W</i>	9.799	14.376	14.581	14.155
<i>R</i>	8.462	8.589	17.419	10.863

$$\chi_{40}^2 = 127.5$$

The maximum likelihood estimates of the 30 parameters and their standard errors, calculated as in Section 5, are shown in table 3.

We should be cautious in accepting that only additive genetic and the two types of environmental effect are present, however. (See *e.g.* Eaves, 1975). When this kind of model is fitted in the classical twin study the E_2 estimate tends to be a conglomeration of certain other between-family effects which may be present. If assortative mating is present for example, and the population is in equilibrium then any supposed E_2 is really

$$E_2 + \frac{1}{2} \frac{A}{1-A} \cdot D_R$$

where A is the correlation between the breeding values of spouses, E_2 is the "true" environmental variance between families, and $\frac{1}{2}D_R$ is the "true" additive genetical variance obtained by randomly mating a population with the same gene frequencies.

This immediately suggests a possible reduction of our 30 parameter model. Suppose, for a given variable i the component of variance corresponding to the E_2 factor loading, ϵ_i^2 , were a simple function of an A_i and the corresponding component of variance arising from the additive genetical factor loading, δ_i^2 . Then if there were no contributions of between-families environmental effects ("real" E_2) to the communality of the variables, we should expect A_i to be the same for all i and ϵ_i/δ_i to be a constant.

TABLE 3

Estimates, standard errors and significance of parameters of full $E_1E_2D_R$ model

	$\hat{\lambda}_{E_1}$	s.e.	c	$\hat{\theta}_{E_1}$	s.e.	c
N	1.159	1.989	0.58	19.216	1.218	15.78
V	10.201	1.525	6.69	7.821	1.746	4.48
S	2.123	2.132	1.00	21.094	1.346	15.67
W	6.580	1.376	4.78	12.706	0.918	13.85
R	6.629	1.150	5.77	8.853	0.845	10.48
	$\hat{\lambda}_{E_2}$	s.e.	c	$\hat{\theta}_{E_2}$	s.e.	c
N	19.463	5.237	3.72	1.761	70.668	0.02
V	26.662	3.132	8.51	2×10^{-8}	2×10^7	0.00
S	11.918	4.887	2.44	27.504	5.664	4.86
W	12.453	2.434	5.12	1×10^{-8}	3×10^6	0.00
R	14.720	2.333	6.31	11.414	2.012	5.67
	$\hat{\lambda}_{D_R}$	s.e.	c	$\hat{\theta}_{D_R}$	s.e.	c
N	50.897	7.710	6.60	20.972	18.219	1.15
V	16.830	4.461	3.77	17.114	4.167	4.11
W	31.027	6.352	4.88	37.647	8.339	4.51
R	14.686	3.433	4.28	22.926	4.552	5.04
R	14.451	3.140	4.60	1.628	26.125	0.06

$$\chi_{30}^2 = 33.0$$

We can test this hypothesis by constraining the E_2 and D_R factor loadings to be related by a scalar constant such that $\epsilon_i/\delta_i = b$, for all i . In terms of Jöreskog's model, if E_2 is the second basic source of variation then:

$$B_2 = \begin{pmatrix} b & 0 & 0 & 0 & 0 \\ 0 & b & 0 & 0 & 0 \\ 0 & 0 & b & 0 & 0 \\ 0 & 0 & 0 & b & 0 \\ 0 & 0 & 0 & 0 & b \end{pmatrix} = bI$$

From our 30 parameter model we can thus remove five factor loadings and substitute a single scalar parameter b , leaving 26 parameters in all. The log likelihood in this model was -7301.09 , corresponding to a chi-square of 52.4 on 34 d.f. ($0.01 < P < 0.02$). The estimates of this model are shown in table 4. The failure of this model indicates that, while there may be assortative mating based upon the common genetical variance components, this alone cannot explain the so-called common E_2 components and that in addition there must be some "cultural" common factor whose factor pattern need not follow that of the genetical loadings.

TABLE 4

Parameter estimates of $E_1E_2D_R$ model with loadings on E_2 and D_R factors related by a constant, b

	$\hat{\lambda}_{E_1}$	$\hat{\theta}_{E_1}$	$\hat{\lambda}_{E_2}^*$	$\hat{\theta}_{E_2}$	$\hat{\lambda}_{D_R}$	$\hat{\theta}_{D_R}$
<i>N</i>	0.893	19.349	25.210	1×10^{-5}	35.948	35.025
<i>V</i>	10.127	7.447	18.998	8.226	27.090	18.903
<i>S</i>	1.292	21.104	16.821	25.983	23.985	41.523
<i>W</i>	5.927	12.862	11.731	5×10^{-6}	16.727	22.574
<i>R</i>	6.341	8.913	12.874	10.754	18.358	2×10^{-6}

$$\chi_{34}^2 = 52.4$$

* The $\hat{\lambda}_{E_2}$'s are the product of $\hat{b} = 0.7013$ and the corresponding $\hat{\lambda}_{D_R}$

However, as an exercise it is interesting to calculate the degree of assortative mating which would be implied if the model fitted. We have

$$\epsilon_i^2 = \frac{1}{2} \frac{A}{1-A} \delta_i^2$$

and

$$\epsilon_i^2 = b^2 \delta_i^2$$

$$\frac{A}{1-A} = 2b^2$$

Our maximum likelihood estimate of b is 0.7013 giving an estimate of $A = 0.4958$. This is somewhat higher than other estimates of A for *IQ* (Eaves, 1973, 1975) probably reflecting once again the presence of a genuine E_2 common factor.

Attempting in a similar way to relate the loadings on either the E_2 or D_R factors to the loadings on the E_1 factor would clearly be inappropriate since two of the variables make only specific contributions to E_1 (table 3).

Since this attempt to reduce the model failed, we must accept the 30 parameter model as the most appropriate although we may strike out any non-significant parameters to obtain our final solution. We calculate the covariance matrix of the estimates and use a c -test to assess the significance of each parameter (table 3). The factor loadings can either be positive or negative so a two-tailed test should be used for these ($c < 1.96$) but the specific variances are constrained to be positive so only a one-tailed test should be used for the $\hat{\theta}_i$'s ($c < 1.65$).

If we now fix to zero all non-significant parameters our reduced model contains 23 free parameters and our final $\chi_{37}^2 = 35.01$. Thus the seven non-significant parameters account for a chi-square of only 2.01, although the increased significance of the remaining parameters partly reflects our *post-hoc* reduction of the model. The final values of the parameters and their recalculated standard errors are given in table 5.

We are now in a position to summarise the contribution of each of the six sources of variation to the total variance of each variable. This is done in table 6, along with the sub-totals of variation attributable to E_1 , E_2 and to $\frac{1}{2}D_R$ which would be the heritability if we could ignore assortative mating and non-additive genetical effects. From these proportions we can

calculate any other desired ratios such as the heritabilities of the general or specific variation and the relative contributions of common and specific factors to genetical and environmental variation.

A breakdown of the total variation such as this must be one of the end products of a genetical analysis of covariance structures. Our analysis

TABLE 5

Estimates, standard errors and significance of parameters of reduced $E_1E_2D_R$ model with non-significant parameters fixed to zero

	λ_{E_1}	s.e.	c	θ_{E_1}	s.e.	c
N	0.0	—	—	19.343	1.183	16.95
V	10.582	1.579	6.70	7.281	2.069	3.52
S	0.0	—	—	21.202	1.338	15.85
W	6.481	1.352	4.80	12.795	0.885	14.46
R	6.410	1.103	5.81	9.052	0.730	12.41
	λ_{E_2}	s.e.	c	θ_{E_2}	s.e.	c
N	17.768	5.156	3.45	0.0	—	—
V	26.560	2.502	10.61	0.0	—	—
S	12.990	4.389	2.96	27.316	5.685	4.81
W	12.534	2.189	5.73	0.0	—	—
R	14.761	1.983	7.44	11.527	1.041	11.07
	λ_{D_R}	s.e.	c	θ_{D_R}	s.e.	c
N	56.142	3.653	15.37	0.0	—	—
V	17.211	3.992	4.31	17.063	3.433	4.97
S	28.914	4.694	6.16	39.003	7.720	5.05
W	14.279	2.793	5.11	23.016	1.752	13.135
R	14.263	2.559	5.57	0.0	—	—

TABLE 6

Contributions of general and specific components to total variances of five PMA variables

	E_1			E_2			D_R			Grand Total
	$\lambda_{E_1}^2$	$\theta_{E_1}^2$	Total	$\lambda_{E_2}^2$	$\theta_{E_2}^2$	Total	$\lambda_{D_R}^2$	$\theta_{D_R}^2$	Total	
N	—	374.14	374.14	315.60	—	315.60	1575.97	—	1575.97	2265.71
	—	0.165	0.165	0.139	—	0.139	0.696	—	0.696	1.000
111.98	53.02	165.00	705.44	—	—	705.44	148.11	145.57	293.68	1164.12
0.096	0.046	0.142	0.606	—	—	0.606	0.127	0.125	0.252	1.000
—	449.52	449.52	168.72	746.19	914.91	914.91	418.01	760.59	1178.60	2543.03
	—	0.177	0.177	0.066	0.294	0.360	0.164	0.299	0.463	1.000
W	42.00	163.72	205.72	157.09	—	157.09	101.95	264.86	366.81	729.62
0.058	0.224	0.282	0.215	—	—	0.215	0.140	0.363	0.503	1.000
R	41.09	81.95	123.04	217.88	132.87	350.75	101.70	—	101.70	575.49
0.071	0.142	0.213	0.379	0.231	0.610	0.610	0.177	—	0.177	1.000

shows that although there is common genetical variation contributing to variance in each of the five sub-tests, in three of the sub-tests specific genetical variation is equally or more important. While others (Nichols, 1965; Eaves and Gale, 1974; Martin, 1975) have had to be content with detecting the presence of specific genetical variation in different ability traits we have been able to provide maximum likelihood estimates of these and all the other components.

It can also be seen that although there is a hint of a small E_1 factor corresponding to the verbal traits, nearly all the E_1 variance is specific and this is what we should expect of a source which comprises error and environmental experiences specific to individuals.

Clearly the most unsatisfactory aspect of this particular example is our inability to separate the common factor variance due to assortative mating from a genuine cultural factor. However, this is a failing of the classical twin design rather than our analytical method and a more elaborate constellation of relatives would allow us to make this separation.

A further aim of the genetical analysis of covariation might be to compare and contrast the patterns of correlation of the different sources of variation. Clearly, if all the variation from a given source is due to a common factor then all the correlations between variables for that source will be unity. If all the variation from a source is specific, then all the intercorrelations will be zero. Table 7 shows the intercorrelations of the five variables for

TABLE 7
Correlations for E_1 , E_2 and D_R sources of variation

	E_1					E_2					D_R						
	N	V	S	W	R	N	V	S	W	R	N	V	S	W	R		
N	1.00	0.00	0.00	0.00	0.00	N	1.00	1.00	0.43	1.00	0.79	N	1.00	0.71	0.60	0.53	1.00
V	—	1.00	0.00	0.37	0.48	V	—	1.00	0.69	1.00	0.79	V	—	1.00	0.42	0.37	0.71
S	—	—	1.00	0.00	0.00	S	—	—	1.00	0.43	0.34	S	—	—	1.00	0.31	0.60
W	—	—	—	1.00	0.25	W	—	—	—	1.00	0.79	W	—	—	—	1.00	0.53
R	—	—	—	—	1.00	R	—	—	—	—	1.00	R	—	—	—	—	1.00

the three sources of variation. For instance, we obtain the maximum likelihood estimate of the genetical correlation matrix, given that our model is correct, by calculating $(\Delta\Delta' + D^2)$ and scaling each off-diagonal term by the square root of the products of the corresponding diagonal terms:

$$\text{e.g. } r_{DRi,j} = \frac{\delta_i \delta_j}{\sqrt{(\delta_i^2 + d_i^2)(\delta_j^2 + d_j^2)}}$$

As expected, all the E_1 intercorrelations are zero or low, while those for E_2 and D_R are higher reflecting the greater importance of common factors in those two sources.

7. CONCLUSIONS

We should not allow a specific application to obscure the generality of our approach. Given an adequate set of family groupings or generations and a sufficiently strong biological and psychological theory it is possible to formulate and test a model of individual differences which embodies the biological and cultural sources of variation and specifies precisely the way these are expected to affect a number of correlated traits.

Although we have been content to fit a single common factor, because this was inherent in our choice of measurements, the approach can be extended to the estimation of additional correlated or uncorrelated factors as long as appropriate constraints are specified or fixed values are assigned to certain of the factor loadings (Jöreskog, 1973).

In our case we have shown that the multivariate structure of five different ability measures is consistent with a causal explanation in terms of additive gene action and within and between families environmental effects. Although we are unable to determine how much of the latter source may be

due to assortative mating the failure of the model which attempted to estimate all the E_2 factor loadings as a simple multiple of the D_R loadings indicated that there must be some common cultural factor.

As well as yielding maximum likelihood estimates of the factor loadings and specific variances, with all their desirable properties, the approach enables us to test the adequacy of the fitted model and provides us with standard errors of the parameter estimates so that the margin of error attached to the individual estimates can be assessed.

With data on relatives other than twins it would be possible to study in still greater subtlety the mechanisms underlying the multivariate structure of individual differences. With organisms other than man, of course, this is not difficult. Providing the investigator possesses the ingenuity to write the appropriate model and collect the right data the possibilities for the causal analysis of trait covariation in quantitative genetical terms seem extensive.

Acknowledgements.—This work is part of a research programme in psychogenetics supported by the Medical Research Council.

8. REFERENCES

- EAVES, L. J. 1973. Assortative mating and intelligence: an analysis of pedigree data. *Heredity*, *30*, 199-210.
- EAVES, L. J. 1975. Testing models of variation in intelligence. *Heredity*, *34*, 132-136.
- EAVES, L. J. 1976a. A model for sibling effects in man. *Heredity*, *36*, 205-214.
- EAVES, L. J. 1976b. The effect of cultural transmission on continuous variation. *Heredity*, *37*, 41-57.
- EAVES, L. J., AND EYSENCK, H. J. 1975. The nature of extraversion; a genetical analysis. *Journal of Personality and Social Psychology*, *32*, 102-112.
- EAVES, L. J., AND EYSENCK, H. J. 1976. Genetic and environmental components of inconsistency and unrepeatability in twins' responses to a neuroticism questionnaire. *Behavior Genetics*, *6*, 145-60.
- EAVES, L. J., AND EYSENCK, H. J. 1977. A genetic model for psychoticism. *Behaviour Research and Therapy Monograph Supplement*, Vol. 1 (in press).
- EAVES, L. J., AND GALE, J. S. 1974. A method for analysing the genetic basis of covariation. *Behavior Genetics*, *4*, 253-67.
- EAVES, L. J., LAST, K. A., MARTIN, N. G., AND JINKS, J. L. 1977. A progressive approach to non-additivity and genotype-environmental covariance in the analysis of human differences. *British Journal of Mathematical and Statistical Psychology* (in press).
- JINKS, J. L., AND FULKER, D. W. 1970. A comparison of the biometrical genetical, MAVA and classical approaches to the analysis of human behaviour. *Psychological Bulletin*, *73*, 311-49.
- JÖRESKOG, K. G. 1973. Analysis of covariance structures. In *Multivariate Analysis III*, ed. P. R. Krishnaiah. Academic Press, N.Y.
- LOEHLIN, J. C., AND VANDENBERG, S. G. 1968. Genetic and environmental components in the covariation of cognitive abilities: an additive model. In *Progress in Human Behaviour Genetics*, ed. S. G. Vandenberg, pp. 261-278. Johns Hopkins, Baltimore.
- LOOTSMA, F. A. 1972. A survey of methods for solving constrained minimisation problems via unconstrained minimisation. In *Numerical Methods for Non-linear Optimisation*, ed. F. A. Lootsma. Academic Press, London.
- MARTIN, N. G. 1975. The inheritance of scholastic abilities in a sample of twins. II. Genetical analysis of examination results. *Annals of Human Genetics*, *39*, 219-29.
- MATHER, K., AND JINKS, J. L. 1971. *Biometrical Genetics*. Chapman and Hall, London.
- NICHOLS, R. C. 1965. The National Merit Twin Study. In *Methods and Goals in Human Behaviour Genetics*, ed. S. G. Vandenberg. Academic Press, N.Y.
- NUMERICAL ALGORITHMS GROUP 1974. EO4HAF in N.A.G. Library Manual. Mark IV. N.A.G. Central Office, Oxford University, Oxford.
- THURSTONE, L. L. 1945. *Multiple Factor Analysis*. University of Chicago Press, Chicago.