

COMMET 01302

Section II. Systems and programs

TWINAN90: a FORTRAN program for conducting ANOVA-based and likelihood-based analyses of twin data

Christopher J. Williams ^{a,c}, Joe C. Christian ^a and James A. Norton, Jr. ^{a,b}

Departments of ^a Medical Genetics, ^b Psychiatry, and ^c Division of Biostatistics, Department of Medicine, Indiana University School of Medicine, Indianapolis, IN, USA

We discuss the program, TWINAN90, which can perform several different types of analysis of twin data. TWINAN90 incorporates the ANOVA-based twin analyses from the TWINAN twin analysis program, and also includes maximum likelihood estimation of parameters from three path models. Another feature of TWINAN90 is the optional output of a pedigree file which can be read by the quantitative genetics package FISHER. The diagnostic features of the program make TWINAN90 useful also for preliminary analyses prior to the use of more sophisticated modeling procedures which are available in packages such as LISREL and FISHER. An annotated printout from TWINAN90 is presented to illustrate the statistical analyses performed in the program.

Twin analysis; Variance components; Maximum likelihood; Path analysis; Method of moments

1. Introduction

Three commonly used approaches for the analysis of twin data are (i) methods based on the analysis of variance of twin pairs, (ii) maximum likelihood methods based upon covariance matrices from the sample, and (iii) maximum likelihood methods based on individual pedigrees (twin pairs). Analysis of variance based approaches have been described by Kempthorne and Osborne [1], Haseman and Elston [2], and Christian et al. [3], and are currently available in computer programs such as TWINAN [4]. Maximum likelihood analyses of twin data based on sample covariance matrices are described in Rao et al. [5], and are now widely used via the computer package LISREL [6]. Maximum likelihood analyses based on

individual pedigrees are described in Haseman and Elston [2], and Lange et al. [7], and are now available in the quantitative genetics package, FISHER [8]. One of the problems for individuals attempting to understand and interpret these various approaches to twin analysis has been that they each require knowledge of different software. These differing software requirements and consequent difficulty in performing more than one type of analysis has tended to isolate many users from learning about (and using) alternate methods of twin analysis.

We hope to facilitate greater awareness of alternate methods of twin analysis with the release of the computer program TWINAN90. TWINAN90 includes ANOVA-based analyses from the TWINAN twin analysis program, and adds three major new features. The first major feature is expanded diagnostic information for assessing the adequacy of the assumed normal distribution of twin values. This includes a Kolmogorov-Smirnov goodness-of-fit test to the as-

Correspondence: Dr. Christopher J. Williams. Present address: Department of Mathematics and Statistics, The University of Idaho, Moscow, ID 83843, USA.

sumed normal distribution [9, section 9.6] for each zygosity group, and calculation of a suggested Box-Cox data transformation [10] if the goodness-of-fit test for a zygosity has a p value that is less than 0.05.

The second major feature in TWINAN90 is maximum likelihood estimation of parameters for three models, based on the sample covariance matrices for the two twin zygositys. This output is similar to analyses that can be performed using the LISREL computer package, with the only difference being that TWINAN90 obtains parameter estimates for the squares of path model parameters. These parameters are then directly interpretable as variance component estimates. This section of the program will be useful both (i) to introduce users with little path analysis experience to maximum likelihood estimators of genetic variance, and (ii) to provide users with path analysis experience concise analyses of some simple path models that can guide the formulation of more complex models. More elaborate analyses are also possible with the TWINAN90 program, as the program can be run on subsets of the data to test for parameter homogeneity across strata.

The third major feature in TWINAN90 is the optional output of a pedigree file that can be read by the FISHER package. The variable analyzed by TWINAN90 and two other variables can be put into a new file, to be read by FISHER. Among the potential advantages in performing estimation with the FISHER package is more direct assessment of the effect of measured variables on genetic and environmental sources of variation, as in studies of the effect of social contact on twin similarity [11].

2. Statistical methods

2.1. ANOVA-based twin analyses

The data from a univariate twin sample consist of n_1 monozygous (MZ) pairs $\mathbf{y}_{1k} = (y_{11k}, y_{12k})$ and n_2 dizygous (DZ) pairs $\mathbf{y}_{2k} = (y_{21k}, y_{22k})$. A model for the data may be written as:

$$y_{ijk} = \mu_i + a_{ijk} + d_{ijk} + \epsilon_{ijk}, \quad (2.1)$$

where μ_i is the population mean for the i th zygosity, and a , d , and ϵ are the additive genetic, dominance genetic, and residual environmental effects on the k th twin pair. Additive genetic effects combine independently over different alleles at the same genetic locus and over different loci. Dominance genetic effects arise due to allelic interactions at the same locus. The model (2.1) assumes that interactions between genes at different loci are negligible compared to a and d effects. Note that the levels of a and d are unobservable random variables, unlike typical random effects ANOVA models.

By performing separate nested analyses of variance on the two zygositys, estimators of the additive genetic variance σ_a^2 can be developed using the method of moments:

$$\hat{\sigma}_{a(WP)}^2 = 2(MS_{WDZ} - MS_{WMZ}), \quad (2.2)$$

and

$$\hat{\sigma}_{a(AC)}^2 = (MS_{AMZ} - MS_{ADZ} + MS_{WDZ} - MS_{WMZ}), \quad (2.3)$$

where MS_{WDZ} and MS_{WMZ} are the within-pairs mean squares for the DZ and MZ groups, respectively, and MS_{AMZ} and MS_{ADZ} are the among-pairs mean squares for the MZ and DZ groups, respectively. These two estimators are the within-pair and among-component estimators of genetic variance from Christian et al. [3], multiplied by 2 to make them consistent estimators of σ_a^2 .

Kempthorne and Osborne [1], Haseman and Elston [2], and Christian, Kang, and Norton [3], provide details on the development of ANOVA-based analyses of twin data. Several assumptions must be made in order for $\hat{\sigma}_{a(WP)}^2$ and $\hat{\sigma}_{a(AC)}^2$ to be consistent estimators of σ_a^2 ; in particular, the dominance genetic variance σ_d^2 must be zero. Among other requirements are the assumption of equal means for the two twin zygositys, and, for hypothesis tests and confidence intervals, that the twin variable has a bivariate normal distribution. Both refs. [2] and [3] discuss twin model assumptions, with the latter recommending the use of

$\hat{\sigma}_{a(\text{WP})}^2$ when a preliminary test of the equality of the total variances of the two zygositys is not rejected. Christian et al. [3] recommend the use of $\hat{\sigma}_{a(\text{AC})}^2$ when the equal variance hypothesis is rejected.

In addition to obtaining estimators of genetic variance, an estimator of heritability h^2 , which is the proportion of variance of a variable that is due to additive genetic effects, can be developed by dividing an estimator of genetic variance by an estimator of the total variance of the variable. Thus an estimator of heritability based on $\hat{\sigma}_{a(\text{WP})}^2$ is

$$\hat{h}_{\text{WP}}^2 = \frac{\hat{\sigma}_{a(\text{WP})}^2}{\frac{1}{4}(MS_{\text{WDZ}} + MS_{\text{WMZ}} + MS_{\text{ADZ}} + MS_{\text{AMZ}})} \quad (2.4)$$

The estimator \hat{h}_{AC}^2 is defined similarly.

2.2. Likelihood-based twin analyses

Maximum likelihood estimation of parameters in the twin model, first discussed in [2], was extended to more complex models in Rao et al. [5] and Lange et al. [7]. The approach described in [5] expresses the genetic model via a path diagram, which then yields a set of equations relating observed covariances or correlations to model parameters. For a given set of data, maximum likelihood estimates of model parameters can then be obtained. For these analyses, the likelihood function of the sample is based on the covariance matrices of the data. Thus if a sample of n_1 MZ pairs has covariance matrix \mathbf{S}_1 and a sample of n_2 DZ pairs has covariance matrix \mathbf{S}_2 , the log likelihood of the sample is:

$$\begin{aligned} \log L(\theta) &= -\left(\frac{n_1 - 1}{2}\right) \left(\log |\boldsymbol{\Omega}_1(\theta)| + \text{tr}[\mathbf{S}_1 \boldsymbol{\Omega}_1^{-1}(\theta)] \right) \\ &\quad - \left(\frac{n_2 - 1}{2}\right) \left(\log |\boldsymbol{\Omega}_2(\theta)| + \text{tr}[\mathbf{S}_2 \boldsymbol{\Omega}_2^{-1}(\theta)] \right) \\ &\quad + K, \end{aligned} \quad (2.5)$$

$$\boldsymbol{\Omega}_1 = \begin{pmatrix} \sigma_a^2 + \sigma_d^2 + \sigma_\epsilon^2 & \sigma_a^2 + \sigma_d^2 \\ \sigma_a^2 + \sigma_d^2 & \sigma_a^2 + \sigma_d^2 + \sigma_\epsilon^2 \end{pmatrix}$$

$$\boldsymbol{\Omega}_2 = \begin{pmatrix} \sigma_a^2 + \sigma_d^2 + \sigma_\epsilon^2 & \frac{1}{2}\sigma_a^2 + \frac{1}{4}\sigma_d^2 \\ \frac{1}{2}\sigma_a^2 + \frac{1}{4}\sigma_d^2 & \sigma_a^2 + \sigma_d^2 + \sigma_\epsilon^2 \end{pmatrix}$$

Fig. 1. The matrices $\boldsymbol{\Omega}_1$ and $\boldsymbol{\Omega}_2$ for the ADE model.

where $\boldsymbol{\Omega}_i(\theta)$ is the population phenotypic covariance matrix of zygosity i , and is a function of the model parameters θ . For example, $\theta = (\sigma_a^2, \sigma_c^2, \sigma_\epsilon^2)$ for the ACE model discussed below. Also, $|\mathbf{X}|$ is the determinant of the matrix \mathbf{X} , $\text{tr}(\mathbf{X})$ is the trace of the matrix \mathbf{X} , and K is a constant not depending on θ .

Three models that are typically of interest in likelihood-based univariate twin analyses are called the ADE, ACE, and the AE models. The ADE model is derived from the linear model (2.1), and is represented by the matrices $\boldsymbol{\Omega}_1$ and $\boldsymbol{\Omega}_2$ listed in Fig. 1. The ACE model replaces the σ_d^2 term in the ADE model with a σ_c^2 term that represents variation due to common twin environment. The coefficient of σ_c^2 in the covariance between twins is one in both the MZ and DZ groups. Thus, for the ACE model, the off-diagonal term of $\boldsymbol{\Omega}_1$ is $\sigma_a^2 + \sigma_c^2$, while the off-diagonal term in $\boldsymbol{\Omega}_2$ is $\frac{1}{2}\sigma_a^2 + \sigma_c^2$. The AE model has only the two parameters, σ_a^2 and σ_ϵ^2 .

Maximum likelihood estimation for very general path models can be performed with the computer package LISREL [6]. A description of the use of LISREL for analyzing twin data can be found in Heath et al., [12], Neale et al. [13], and in several other articles in the December 1989 issue of Behavior Genetics.

Lange et al. [7] discuss likelihood estimation for quantitative genetics models where the likelihood function is defined on individual pedigrees.

For a sample of n_1 MZ twin pairs and n_2 DZ twin pairs, the log likelihood defined on individual twin pairs is:

$$\begin{aligned}
 \log L(\theta) &= -\left(\frac{n_1}{2}\right) \log |\boldsymbol{\Omega}_1(\theta)| \\
 &\quad - \frac{1}{2} \sum_{i=1}^{n_1} (\mathbf{y}_{1i} - \mu)' \boldsymbol{\Omega}_1^{-1}(\theta) (\mathbf{y}_{1i} - \mu) \\
 &\quad - \left(\frac{n_2}{2}\right) \log |\boldsymbol{\Omega}_2(\theta)| \\
 &\quad - \frac{1}{2} \sum_{i=1}^{n_2} (\mathbf{y}_{2i} - \mu)' \boldsymbol{\Omega}_2^{-1}(\theta) (\mathbf{y}_{2i} - \mu) + K.
 \end{aligned} \tag{2.6}$$

The parameter μ is the population mean of the variable being analyzed.

A wide range of models of the form (2.6) can be analyzed with the quantitative genetics package FISHER [8]. Possible advantages of using a maximum likelihood analysis based on individual pedigrees include direct modeling of measured variables such as social contact [11], and easier identification of outliers and application of robust variance estimation [14]. Robust methods of parameter estimation can also be conveniently applied through use of the multivariate t distribution [15], or for more general weight functions [16,17].

3. Program guide

3.1. General information

TWINAN90 is written in FORTRAN-77 in double precision arithmetic. The source code consists of approximately 1250 lines, and requires 44 kbytes of disk storage space. The program can be run on mainframe computers or microcomputers. TWINAN90 as currently written can analyze data sets consisting of up to 500 twin pairs from each zygosity group. The program can easily be modified, however, to handle larger data sets.

Several algorithms appearing in the subroutines are adapted or obtained from the excellent reference book by Press et al. [18]. The source code of these algorithms is protected by copyright, © 1986 by Numerical Recipes Software, and are not included with the TWINAN90 source code, although the routines are included in the TWINAN90 executable file. The calls to these subroutines are well marked so that individuals wishing to modify TWINAN90 can refer to Press et al. [18] to develop their own source code.

Routines in TWINAN90 for likelihood estimation of genetic parameters calculate derivatives using methods discussed in Lange et al. [7]. Computation of parameter estimates is accomplished via an iteratively reweighted least squares procedure, as discussed in Williams et al. [17]. Maximum likelihood estimation of a suggested Box-Cox transformation, when required, is accomplished by a likelihood grid search.

3.2. Input

Input for TWINAN90 is illustrated via a sample run displayed in Fig. 2. After the program is called, it requests the name of the data file. The first line of the data file must list the number of MZ twin pairs, followed by the number of DZ twin pairs, with at least one space between them to allow the numbers to be read in free format. Also, the data file must be sorted with all MZ twin pairs occurring first, followed by the DZ twin pairs. The members of a twin pair must occur sequentially for each twin pair, all with the same input format.

The program then requests a variable name that can be up to 16 characters long. Next, the names to be assigned to the two default output files are requested. Each of these filenames can also be up to 16 characters long. Next the user is prompted to provide the FORTRAN input format for the data. The input format is used to read the variables from both members of the twin pair. In the example (Fig. 2), the data is arranged to have one person per line, so the format (16X, F7.2/16X, F7.2) reads the first twin's value on the first line, then moves down to the second line to read the second twin's value.

```

twinan90
ENTER THE NAME OF THE DATA FILE.
bone.dat
ENTER THE NAME OF THE VARIABLE.
WID18
ENTER THE NAME OF THE OUTPUT FILE.
wid18.res
ENTER THE NAME OF THE DIAGNOSTIC OUTPUT FILE.
wid18.chk
ENTER THE INPUT FORMAT (A FORTRAN FORMAT).
(16x,f7.2/16x,f7.2)
ENTER THE POWER TRANSFORMATION PARAMETER.
ENTER 1. IF YOU DO NOT WISH TO TRANSFORM YOUR DATA.
1.
ENTER A VALUE TO ADD TO EACH DATA POINT.
40.
ENTER THE VALUE FOR MISSING DATA.
-99.00
ENTER 1. IF YOU WISH TO OBTAIN
PATH MODEL RESULTS, OTHERWISE ENTER 0.
1.
ENTER 1. IF YOU WISH TO PRODUCE
AN OUTPUT PEDIGREE FILE, OTHERWISE ENTER 0.
1.
THE SAMPLE SIZES ARE 91.00000000 AND 31.00000000
ENTER THE NAME OF THE OUTPUT FILE.
wid18.ped
ENTER THE NUMBER OF VARIABLES TO BE OUTPUT (1-3).
1.
ENTER THE FORTRAN FORMAT STATEMENT FOR EACH TWIN
PAIR. (INCLUDE AN "A1" CHARACTER IN THE FORMAT
FOR EACH PERSON, PRECEDING THE VARIABLES. THIS
VARIABLE CAN BE USED TO IDENTIFY TWIN GENDER.
AN EXAMPLE IS (1X,A1.5X,F6.2/1X,A1.5X,F6.2) FOR
OUTPUTING ONE VARIABLE. THE "A1" MUST BE INCLUDED
EVEN IF IT IS A DUMMY VARIABLE, AND NO INFORMATION
IS AVAILABLE ON TWIN GENDER.)
(4x,a1,11x,f7.2/4x,a1,11x,f7.2)
.
. (program runs)
.
IF YOU WOULD LIKE TO RUN THE PROGRAM AGAIN,
ENTER A 1, OTHERWISE ENTER A 0.
0

```

Fig. 2. A sample run of the TWINAN90 program.

The next three values requested by TWINAN 90 are a transformation parameter from the Box-Cox family of transformations, a constant to add to each data point to make all values positive, and a code that is assigned to missing values. The option of adding a constant to the data is useful because the subroutine that calculates a suggested transformation parameter requires that all values be positive. Also, the subroutine is not run if a transformation parameter other than 1.0 is entered.

The last two options in the TWINAN90 sample run allow likelihood estimates of the ADE, ACE, and AE models to be computed, and allow an output pedigree file to be produced that can be read by the quantitative genetics package FISHER. If an output pedigree file is requested by the user, then TWINAN90 prompts for a filename and a FORTRAN format statement for

reading up to three variables and an identification variable per person. Note that the 'A1' code must be entered whether or not an identification variable exists. After running, TWINAN90 returns to ask if another run is requested.

3.3. Output

TWINAN90 always produces two output files: a file containing results of the ANOVA-based analyses and if requested, likelihood-based analyses; and a second 'diagnostic' file that is useful for ensuring that the program ran satisfactorily. We will describe the information in the diagnostic file in this section, and will describe output from the primary output file in the following section with an example. Information on the optional output pedigree file and on the FISHER quantitative genetics programs can be found by consulting documentation for FISHER [8].

The diagnostic file contains three types of information: a listing of the data that was analyzed, a listing of the likelihood grid search for a suggested transformation of the data (if required), and the results from the iterative process of computing the maximum likelihood estimates for the ADE, ACE, and AE models (if requested).

The file first lists data that was analyzed, both before and after transformation. If one member of a twin pair has a value that has been user-designated as missing (see Section 3.2), then the pair is not analyzed or listed in the diagnostic file. Thus, the file can be checked to ensure that the data were read correctly and that no missing values were included in the analyses.

Second, if the data from a zygosity group failed the goodness-of-fit test, the likelihood L_λ of the Box-Cox transformation λ is listed for values of λ from -1 to 2 in increments of 0.1 , along with the value of λ having the largest likelihood value for that zygosity group. The value of λ for a zygosity group is set equal to 1 if the Kolmogorov-Smirnov test has a p value ≥ 0.05 for that zygosity.

The third and last section of the diagnostic file lists the results of the iterations leading to the maximum likelihood estimates for the ADE, ACE, and AE models, if requested by the user. For each iteration, the updated parameter estimates

are listed, along with the first derivatives of the loglikelihood, the current loglikelihood value, and the criterion value that is used to determine when to stop the likelihood search. The criterion value is equal to the Euclidean norm of the vector difference between the previous and current parameter values, divided by the Euclidean norm of the current vector of parameter estimates. The iterative process stops when the criterion is less than 1×10^{-12} , or when 500 iterations have been

completed, whichever occurs first. It is important to check the diagnostic file to insure that the likelihood search ended at a maximum, with derivatives equal to 0. Otherwise, if a maximum likelihood value was not found in 500 iterations, the program will print the parameter estimates for the model for iteration 500. Fortunately, in our experience this problem usually only occurs when the data set is very small (less than 5 twin pairs for one of the zygosity groups).

TABLE 1

TWINAN90 output, first page (See text for details)

TWINAN90 TWIN ANALYSIS PROGRAM FOR THE VARIABLE WID18				
NO. OF TWIN PAIRS: MZ=91 DZ=31				
ASSESSMENT OF MODEL ASSUMPTIONS, DESCRIPTIVE STATISTICS				
	INDIVIDUALS		TWIN ABSOLUTE DIFFERENCES	
	MZ	DZ	MZ	DZ
MEAN	2.5695	2.5929	0.0690	0.1204
VAR	0.0330	0.0323	0.0030	0.0111
SKEW;SE	0.07; 0.18	0.10; 0.31		
EX KURT;SE	-0.64; 0.36	0.07; 0.62		
KOLMOGOROV-SMIRNOV TEST OF NORMALITY (VALID FOR N>25 TWIN PAIRS)				
MZ: D"=0.57; P>0.150 DZ: D"=0.66; P>0.150				
TEST FOR DIFFERENCE IN MZ AND DZ MEANS:				
MZ-DZ=-0.0234 T"=-0.68 DF=56.1 P=0.500				
ANALYSIS OF VARIANCE				
	MZ	DF	DZ	DF
AMONG MEAN SQUARES:	0.0625	90	0.0525	30
WITHIN MEAN SQUARES:	0.0039	91	0.0128	31
SUM OF MEAN SQUARES:	0.0663	101.2	0.0653	43.9
TEST OF TWIN MODEL:				
SUM OF MEAN SQUARES RATIO: F"=0.9846 APPROX P=0.979				
THE EQUAL VARIANCE HYPOTHESIS IS NOT REJECTED.				
ESTIMATES OF GENETIC VARIANCE				
	ESTIMATE	VARIANCE	F-RATIO	PROB
WITHIN PAIR (WP):	0.008882	0.000010	3.28	0.000
AMONG COMPONENT (AC):	0.009394	0.000067	1.33	0.167
INTRACCLASS CORRELATIONS: MZ=0.8825; P=0.000 DZ=0.6087; P=0.000				
AVERAGE ABS DIFFERENCE TEST: T"=2.83 DF=37.6 P=0.008				
HERITABILITY ESTIMATES				
	FORMULA	ESTIMATE	PROBABILITY	
WITHIN PAIR:	$2 * (WDZ - WMZ) / ((SMZ + SDZ) / 4)$	0.53963	SEE PROB (WP)	
AMONG COMPONENT:	$(2 * AC) / ((SMZ + SDZ) / 4)$	0.57071	SEE PROB (AC)	
INTRACCLASS CORRELATION:	$2 * (RMZ - RDZ)$	0.54754	0.001	

4. Application

4.1. Model diagnostics and ANOVA-based estimates

An example of output from the new program is shown in Tables I and II for a study of the genetic effects on bone mass in adult women [19]. The variable, named WID18, is a measure of average bone width along a section of the distal radius. The output is divided into two pages, with descriptive information, tests of model assumptions, and ANOVA-based analyses on the first page, with likelihood-based analyses on the second page.

On page 1 of the output, after listing the

variable name and numbers of MZ and DZ pairs that were included in the analyses, the mean, variance, skewness, and excess kurtosis (kurtosis minus its expected value) of y_{ijk} are displayed for the grouped twin pairs of each zygosity. The standard errors of the skewness and excess kurtosis estimates are printed to allow a check of whether the observed values are consistent with those expected for a normal distribution. Next, a Kolmogorov-Smirnov goodness-of-fit test is applied separately to the grouped pairs for each of the two zygositys. Since the calculations for the Kolmogorov-Smirnov test use means and variances that are estimated from the data, p values for the test are based on the results of Stephens [20].

TABLE 2

TWINAN90 output, second page (See text for details)

TWINAN90 PATH ANALYSIS OUTPUT FOR THE VARIABLE WID18					
NO. OF TWIN PAIRS: MZ=91 DZ=31					
FOR THE COVARIANCE MATRICES					
MZ=	0.0331	0.0291	DZ=	0.0245	0.0197
	0.0291	0.0333		0.0197	0.0403
MODEL 1 MAXIMUM LIKELIHOOD (ML) ESTIMATES					
(ADD GEN VAR, DOM GEN VAR, ERR VAR)=	(0.0519,	-0.0230,	0.0041)		
WITH STANDARD ERRORS OF	0.0162	0.0142	0.0006		
THE ASYMPTOTIC CORRELATION MATRIX=	1.0000	-0.9760	0.0706		
	-0.9760	1.0000	-0.1080		
	0.0706	-0.1080	1.0000		
GOODNESS-OF-FIT LIKELIHOOD RATIO STATISTIC=2.990 ON 3 DF, P=0.394					
MODEL 2 ML ESTIMATES					
(ADD GEN VAR, SHARED ENV VAR, ERR VAR)=	(0.0174,	0.0115,	0.0041)		
WITH STANDARD ERRORS OF	0.0066	0.0071	0.0006		
THE ASYMPTOTIC CORRELATION MATRIX=	1.0000	-0.8424	-0.1770		
	-0.8424	1.0000	0.1080		
	-0.1770	0.1080	1.0000		
GOODNESS-OF-FIT LIKELIHOOD RATIO STATISTIC=2.990 ON 3 DF, P=0.394					
MODEL 3 ML ESTIMATES					
(ADD GEN VAR, ERR VAR)=	(0.0280,	0.0040)			
WITH STANDARD ERRORS OF	0.0036	0.0006			
THE ASYMPTOTIC CORRELATION MATRIX=	1.0000	-0.1243			
	-0.1243	1.0000			
GOODNESS-OF-FIT LIKELIHOOD RATIO STATISTIC=4.930 ON 4 DF, P=0.294					
MODEL 3 VS MODEL 1 LIKELIHOOD RATIO STATISTIC=1.940 ON 1 DF, P=0.164					
MODEL 3 VS MODEL 2 LIKELIHOOD RATIO STATISTIC=1.940 ON 1 DF, P=0.164					
THE OUTPUT PEDIGREE FILE IS NAMED WID18.PED					

If the data from either zygosity fails the Kolmogorov-Smirnov goodness-of-fit test ($p < 0.05$), then a subroutine is called to suggest a Box-Cox transformation for that group. The transformation parameter λ is equal to one for a zygosity if its p value is ≥ 0.05 . The average of the λ values for the two zygosity groups is printed below the Kolmogorov-Smirnov test results if data from at least one zygosity fails the test.

Next, the assumption of equal means of the two zygosity groups is tested by a method discussed in Christian and Norton [21]. The mean and variance of the twin absolute differences, $|X_{i1} - X_{i2}|$, for each zygosity are also printed in this section.

For the bone width data, the model assumptions appear to be met, with both the MZ and DZ data sets having reasonably high p values for the null hypothesis of normality, according to the Kolmogorov-Smirnov test results. This is also reflected by the skewness and excess kurtosis values, which are within two standard errors of 0 for both zygosity groups. Additionally, the difference between the MZ and DZ means is not significantly different from zero. Thus all of the twin model assumptions that are tested are met.

The next section of the output lists the analyses of variance results. The results of the test of equality of the MZ and DZ total variances is then printed at the bottom of the section under the heading 'Test of the Twin Model'. Following the rejection rule suggested by Christian et al. [3], the significance level for the equal variance test is conservatively set at 0.2. The output of TWINAN90 includes the two-tailed p value for the test, followed by a message that either (i) the equal variance hypothesis is not rejected, (ii) the equal variance hypothesis is rejected, with DZ variance greater than MZ variance, or (iii) the equal variance hypothesis is rejected, with MZ variance greater than DZ variance.

The following section of output lists estimates of genetic variance, including the within-pair and among-component estimates of genetic variance. The MZ and DZ intraclass correlations are also listed, as they are often interpreted as measures of genetic variance. The final line of output in this section is a test for genetic variance based on

the average absolute difference between twins [22]. This test has been found to perform well under a variety of failures of the assumptions of the twin model when the equal variance assumption is correct [23]. For this test the difference between zygosity groups of the average absolute twin differences is standardized to yield an approximate t-test for the null hypothesis of no genetic variance.

The final section on the first page presents three estimates of heritability. The first two heritability estimates are derived from the within-pair (WP) and among-component (AC) estimates of genetic variance, respectively. The final estimate of heritability is based on the intraclass correlations between twins of the two twin zygosity groups. The p values listed with the within-pair and among-component heritability estimates are taken from the corresponding genetic variance estimators, on the presumption that significant genetic variance implies nonzero heritability.

The results for the bone width data are fairly straightforward, with the trait in question appearing highly heritable. The test of equal variances of the two twin zygosity groups is not rejected, thus the within-pair and average difference tests are preferred as tests of genetic variance. Both of these tests reject the null hypothesis of no genetic variance, indicating that the trait is heritable. Although attention typically centers on the within-pair and intraclass correlation estimates of heritability when the equal variance assumption is not rejected, for these data all three ANOVA-based heritability estimates are quite similar, showing a high heritability for the trait of distal radius width. The comparative inefficiency of the AC estimator, however, is exemplified by its lack of significance.

4.2. Maximum likelihood analyses

The second page of output, in Table 2, contains the results of the likelihood-based analyses. The first section on page 2 is identical to the first section on page 1, and identifies the variable and numbers of twin pairs being analyzed. Next, the MZ and DZ covariance matrices, which the likelihood functions are based on, are printed. Twin

covariances for each zygosity are computed from intraclass correlations, which are invariant under the ordering of twins within a twin pair.

The next section of the output presents maximum likelihood estimates for the ADE model. These estimates are equal to the squared estimates of path parameters that would be obtained from the LISREL package. The parameter estimates from TWINAN90 are thus directly interpretable as components of variance. The standard errors for estimates from the model are related to those obtained from a LISREL analysis via the equation:

$$s.e.(\hat{\theta}^2) = 2\hat{\theta}s.e.(\hat{\theta}),$$

where $\hat{\theta}$ is the LISREL estimate and $\hat{\theta}^2$ is the TWINAN90 estimate. The asymptotic correlation matrix and goodness-of-fit statistic in TWINAN90 are the same as those from LISREL.

The next section contains estimates for the ACE model. The following section lists the parameter estimates for the AE model. The last statistical results on page 2 are likelihood-ratio statistics for testing improvement in fit of the ADE or ACE models over the reduced AE model. In practice, only one of these tests will be of interest, as determined by the values of r_{MZ} and r_{DZ} . The ACE model will be preferable to the ADE model on pragmatic grounds when the DZ intraclass correlation exceeds half the MZ intraclass correlation, since an effect due to shared environment will tend to raise r_{DZ} above half of r_{MZ} . Otherwise, the ADE model is preferable, because a dominance effect will tend to lower r_{DZ} below half of r_{MZ} . Finally, below the last section of output of the maximum likelihood analyses, is a statement reporting whether or not an output pedigree file was created.

For the distal bone radius width example, the ACE model is more relevant than the ADE model as $r_{DZ} > r_{MZ}/2$. The likelihood ratio statistic testing for improvement in fit of ACE over the AE model is 1.94 on 1 degree of freedom ($p = 0.164$), which indicates that the AE model is the most parsimonious. Although the proportion of variance explained by additive genetic variance in the AE model (0.88) is much higher than the

ANOVA-based estimates of heritability, it is interesting to note that the proportion explained by additive genetic variance in the ACE model (0.53) is much closer to the ANOVA-based estimates.

Although this example is relatively straightforward to interpret, in many cases the output from a twin analysis is much more complex. Transformations are often required to satisfy the normality assumptions, and in some instances the equal variance hypothesis is rejected. When the equal variance hypothesis is rejected, then, as discussed above, close attention must be paid to the choice of estimator of genetic variance and heritability. Often heritability estimates are not in agreement, and occasionally heritability estimates fall outside the [0,1] range to which they are theoretically restricted. Examples of differing heritability estimates, and of possible interpretations in these cases, are discussed in Christian et al. [24].

5. Availability

The TWINAN90 program is available from the author at no charge. Please send a request for the program to the following address: Dr Christopher J. Williams, Attn: TWINAN90, Department of Mathematics and Statistics, The University of Idaho, Moscow, ID 83843, U.S.A.

Acknowledgements

The authors thank Dr C.J. Brown and Dr J. Arnold for their valuable comments. This work was supported in part by PHS Grants AA 07611, HL 46674, and AG 05793.

References

- [1] O. Kempthorne, R. Osborne, The interpretation of twin data, *Am. J. Hum. Genet.* 13 (1961) 320–339.
- [2] J.K. Haseman, R.C. Elston, The estimation of genetic variance from twin data, *Behav. Genet.* 1 (1970) 11–19.
- [3] J.C. Christian, K.W. Kang and J.A. Norton Jr., Choice of an estimate of genetic variance from twin data, *Am. J. Hum. Genet.* 26 (1974) 154–161.

- [4] K.W. Kang, TWINAN: Twin data analysis program for microcomputers, *Acta Genet. Med. Gemellol.* 34 (1985) 113–114.
- [5] D.R. Rao, N.E. Morton, S. Yee, Analysis of familial resemblance. II. A linear model for familial correlation, *Am. J. Hum. Genet.* 26 (1974) 331–359.
- [6] K.G. Jöreskog and D. Sörbom, LISREL VI, Scientific Software, Mooresville, Ind. (1986).
- [7] K. Lange, J. Westlake and M.A. Spence, Extensions to pedigree analysis. II. Variance components by the scoring method, *Ann. Hum. Genet.* 39 (1976) 485–491.
- [8] K.L. Lange, D. Weeks and M. Boehnke, Programs for pedigree analysis: MENDEL, FISHER, and dGENE, *Genet. Epidemiol.* 5 (1988) 471–472.
- [9] P.J. Bickel and K.A. Doksum, *Mathematical Statistics: Basic Ideas and Selected Topics* (Holden-Day, Inc., San Francisco, 1977).
- [10] G.E.P. Box and D.R. Cox, An analysis of transformations, *J. Roy. Stat. Soc., Ser. B* 26 (1964) 211–243.
- [11] J.L. Hopper and P.R. Culross, Covariation between family members as a function of cohabitation history, *Behav. Genet.* 13 (1983) 459–471.
- [12] A.C. Heath, M.C. Neale, J.K. Hewitt, L.J. Eaves and D.W. Fulker, Testing structural equation models for twin data using LISREL, *Behav. Genet.* 19 (1989) 9–35.
- [13] M.C. Neale, A.C. Heath, J.K. Hewitt, L.J. Eaves and D.W. Fulker, Fitting genetic models with LISREL: Hypothesis testing, *Behav. Genet.* 19 (1989) 37–49.
- [14] T.H. Beaty, S.G. Self, K.Y. Liang, M.A. Connolly, G.A. Chase and P.O. Kwiterovich, Use of robust variance components models to analyse triglyceride data in families, *Ann. Hum. Genet.* 49 (1985) 315–328.
- [15] K.L. Lange, R.J.A. Doolittle and J. Taylor, Robust statistical modeling using the *t* distribution, *J. Am. Stat. Assoc.* 84 (1989) 881–897.
- [16] D. Pregibon, Resistant fits for some commonly used logistic models with medical applications, *Biometrics* 38 (1982) 485–498.
- [17] C.J. Williams, W.W. Anderson and J. Arnold, Generalized linear modeling methods for selection component experiments, *Theor. Pop. Biol.* 37 (1990) 389–423.
- [18] W.H. Press, B.P. Flannery, S.A. Teukolsky and W.T. Vetterling, *Numerical Recipes, The Art of Scientific Computing* (Cambridge University Press, Cambridge, 1986).
- [19] C.W. Slemenda, J.C. Christian, C.J. Williams, J.A. Norton, Jr. and C.C. Johnston Jr., Genetic determinants of bone mass in adult women: a reevaluation of the twin model and the potential importance of gene interaction on heritability estimates, *J. Bone Min. Res.* 6 (1991) 561–567.
- [20] M. Stephens, E.D.F. statistics for goodness of fit, *J. Am. Stat. Assoc.* 69 (1974) 730–737.
- [21] J.C. Christian and J.A. Norton Jr., A proposed test of the difference between the means of monozygotic and dizygotic twins, *Acta Genet. Med. Gemellol.* 26 (1977) 49–53.
- [22] L.A. Corey, K.W. Kang, J.C. Christian J.A. Norton Jr., R.E. Harris and W.E. Nance, Effects of chorion type on variation in cord blood cholesterol of monozygotic twins, *Am. J. Hum. Genet.* 28 (1976) 433–441.
- [23] J.E. Bailey-Wilson, The effects of failures of assumptions on several tests used for genetic analysis (Ph.D. thesis, Indiana University, 1980).
- [24] J.C. Christian, N.O. Borhani, W.P. Castelli, R. Fabsitz, J.A. Norton Jr., T. Reed, R. Rosenman, P.D. Wood and P.L. Yu, Plasma cholesterol variation in the National Heart, Lung and Blood Institute Twin Study, *Genet. Epidemiol.* 4 (1987) 433–446.