

## Genetics in Psychosomatic Medicine: Research Designs and Statistical Approaches

JEANNE M. McCAFFERY, PhD, HAROLD SNIEDER, PhD, YANBIN DONG, MD, PhD, AND ECO DE GEUS, PhD

It has become increasingly clear that genetic factors influence many of the behaviors and disease endpoints of interest to psychosomatic medicine researchers. There has been increasing interest in incorporating genetic variation markers into psychosomatic research. In this Statistical Corner article, we build on the valuable experiences gained during two workshops for “starters in the field” at the American Psychosomatic Society and the Society for Psychophysiological Research to review two common genetically informative research designs for human studies: twin and genetic association studies. We outline statistical techniques for each and, for genetic association studies, address special topics, including the treatment of race and ethnicity, gene  $\times$  gene and gene  $\times$  environment interaction, haplotype analysis, and power and sample size. Finally, we discuss the issue of nonreplication and interpretation of results derived from genetic association studies. We hope this overview of twin and genetic association designs will support and stimulate thoughtful applications of genetic approaches within psychosomatic medicine. **Key words:** statistics, genetics, twin studies.

**MZ** = monozygotic; **DZ** = dizygotic; **SES** = socioeconomic status; **SEM** = structural equation modeling; **SBP** = systolic blood pressure; **HWE** = Hardy-Weinberg equilibrium; **SNP** = single nucleotide polymorphism; **VNTR** = varying number of tandem repeats; **LD** = linkage disequilibrium; **TDT** = transmission disequilibrium test.

### INTRODUCTION

Recently researchers in psychosomatic medicine have witnessed an increasing awareness of the importance of genetic factors in stress and health relationships (1–6). Twin and family studies have confirmed a clear-cut genetic contribution to cardiovascular disease (7–11) and its major risk factors (12–21). In addition, tracking of these risk factors over time (22–25) as well as their tendency to cluster, as in the metabolic syndrome (26–29), is largely due to genetic factors. This finding is relevant to psychosomatic medicine. By lumping together subjects who are genetically susceptible to the effects of psychosocial stressors with those subjects who are not susceptible, previous studies may have underestimated the significance of negative health effects in the former susceptible group. Future research, therefore, should strive to include genetic variation as a potential source of individual variance in psychosomatic risk factors.

We aim to review two common research designs in genetics. The starting point of genetic research on any risk factor is the establishment of significant heritability. The twin study has been the work horse of such heritability estimation and we will start by reviewing its principles. Because most researchers in this field are expected to use candidate gene association approaches, the largest part of this paper will consider the statistical methods for this type of association. Throughout, we based this paper on the valuable experiences gained during two workshops for “start-

ers in the field” at the American Psychosomatic Society (30) and the Society for Psychophysiological Research (31). Although we expect that some statistical approaches may be familiar to the readers of *Psychosomatic Medicine*, some genetic terminology may not be. A glossary of genetic terms is available at <http://www.genome.gov/glossary.cfm>.

### Twin Studies

Perhaps one of the most robust clinical observations in psychiatry and cardiology is that disease tends to “run in the family”. However, familial resemblance for a trait cannot automatically be attributed to genes. In family studies, the genetic relatedness is confounded with the shared environment of the family members. This includes potentially important sources of interindividual variance like culture, socioeconomic status (SES), neighborhood, school, sports club, peers, family diet, and parental rearing style and attitudes. A unique experiment of nature has provided the solution to separating genetic and shared environmental influences: the existence of monozygotic (MZ) and dizygotic (DZ) twins.

Because MZ twins reared together share part of their environment and 100% of their genes (32) except for some rare exceptions, any resemblance between them is attributed to these two sources of covariance. The extent to which MZ twins do not resemble each other is ascribed to so-called unique or nonshared environmental factors like differential jobs or lifestyle, accidents or other life events, and in childhood, differential treatment by the parents, and nonshared peers. Unique environment also includes measurement error. Resemblance between DZ twins reared together is ascribed to the sharing of both environment and genes. DZ twins share on average 50% of their segregating genes; any resemblance between them attributable to genetic influences will be less than for MZ pairs. The extent to which DZ twins do not resemble each other is due to unique environmental factors and nonshared genetic influences.

Based on molecular genetic theory, we can further divide the genetic variance in two separate parts: a) additive and b) dominant genetic variance. Genetic effects at a single locus are called additive when the effect of one parental allele is added to the effect of the other parental allele. Genetic effects are called dominant when they deviate from purely additive

---

From the Weight Control and Diabetes Research Center (J.M.M.), Brown Medical School and The Miriam Hospital, Providence, RI; Georgia Prevention Institute (H.S., Y.D.) and Department of Pediatrics, Medical College of Georgia, Augusta, GA; Twin Research and Genetic Epidemiology Unit (H.S.), St. Thomas' Hospital, London, UK; and the Department of Biological Psychology (E.D.G.), Vrije Universiteit, Amsterdam, The Netherlands.

Address correspondence and reprint requests to Jeanne M. McCaffery, Weight Control and Diabetes Research Center, 196 Richmond Street, Providence, RI. E-mail: [Jeanne\\_McCaffery@brown.edu](mailto:Jeanne_McCaffery@brown.edu)

Received for publication December 19, 2005; revision received September 5, 2006.

DOI: 10.1097/PSY.0b013e31802f5dd4

## GENETICALLY INFORMATIVE DESIGNS

effects, e.g., when the two alleles of the locus interact. The total additive and dominance variance estimated in twin studies reflects the additive and dominant effects summed over all contributing loci. The total variance in any trait can arise from the four components identified above: a) unique environmental factors (“E”), b) shared or common environmental factors (“C”), c) additive (“A”) genetic factors, and d) dominant (“D”) genetic factors. For simplicity, we will first consider the case where there is no interaction or correlation among these four components. The value of a trait is then defined as  $P = A + D + C + E$ , where  $P$  is a quantitative trait;  $A$  and  $D$  are the effects of additive and dominant genetic factors; and  $C$  and  $E$  are the effects of common and unique environmental factors (with  $E$  also including the residual variance due to measurement error). The variance ( $V$ ) in trait  $P$  then becomes  $V_P = V_A + V_D + V_C + V_E$ , and the MZ and DZ twin covariances become  $\text{Cov}(\text{MZ}) = V_A + V_D + V_C$ , and  $\text{Cov}(\text{DZ}) = 0.50V_A + 0.25V_D + V_C$ , respectively (33,34).

From the pattern of MZ and DZ twin correlations, we can obtain a first crude estimate of these variance components. However, we cannot estimate common environmental influences and dominant genetic influences at the same time. Therefore, we first test for evidence of dominance, which would yield MZ correlations that are much larger than twice the DZ correlation (e.g.,  $r_{\text{MZ}} = 0.42$ ,  $r_{\text{DZ}} = 0.10$ ). If there is no evidence for dominance, the contribution of additive genetic influences to the total variance in a trait can be estimated as twice the difference between the MZ and DZ correlations ( $V_A/V_P = 2(r_{\text{MZ}} - r_{\text{DZ}})$ ). For instance, typical MZ and DZ correlations for resting systolic blood pressure (SBP) are 0.52 and 0.26 (17); therefore, the percentage of SBP variance explained by the additive genetic influences is estimated at 52%. An estimate of the proportional contribution of the shared environmental influences to the total phenotypic variance is given by subtracting the MZ correlation from twice the DZ correlation ( $V_C/V_P = 2r_{\text{DZ}} - r_{\text{MZ}}$ ). The proportional contribution of the unique environmental influences can be obtained by subtracting the MZ correlation from unit correlation ( $V_E/V_P = 1 - r_{\text{MZ}}$ ). If, for instance, the MZ correlation for exercise behavior of adolescents is 0.8 and the DZ correlation is 0.6, estimates of the relative contribution  $V_A$ ,  $V_C$ , and  $V_E$  to total variance are 40%, 40%, and 20%, respectively (35). If there is evidence for genetic dominance (i.e., the MZ correlation is larger than twice the DZ correlation), the estimate for the proportional contribution of additive genetic influences changes to  $V_A/V_P = (4r_{\text{DZ}} - r_{\text{MZ}})$ . An estimate of the proportional contribution of the dominant genetic influences is then obtained by subtracting four times the DZ correlation from twice the MZ correlation ( $V_D/V_P = 2r_{\text{MZ}} - 4r_{\text{DZ}}$ ).

These are rules of thumb only. They are based on a model that has no interaction terms (e.g.,  $A \times E = 0$ ) and assumes that mating is random, and that the genetic and environmental factors are uncorrelated in the population (e.g.,  $\text{Cov}(A, C) = 0$ ). If these assumptions do not hold, these intuitively simple rules may yield incorrect estimates. Interaction across multiple loci

(gene-gene interaction or epistasis), for instance, will reduce the DZ correlation and inflate the estimate of genetic dominance. Interaction of genetic and unique environmental influences will inflate the contribution of the unique environment and underestimate genetic influences, whereas interaction of genetic and shared environmental factors will inflate the contribution of genetic influences (36,37). Incorrect estimates may also arise when genetic and environmental factors are correlated, for instance, because people actively seek environments that fit their temperament and skills, or because parents pass on their genes as well as a specific environment to their offspring (vertical cultural transmission). Finally, phenotypic assortment, which is nonrandom mate selection based on shared traits (e.g., education, religion, lifestyle choices), increases both MZ and DZ twin correlations that lead to an inflated estimate of the contribution of shared environment.

The other major assumption of the classical twin study is the “Equal Environments Assumption” that MZ twin pairs experience the same degree of environmental similarity as DZ twin pairs. If this is not the case and MZ twin pairs are exposed to more similar environments than DZ pairs, then any excess similarity between MZ pairs compared with DZ pairs may result from environmental rather than genetic factors. Several empirical findings argue in favor of the validity of the equal environment assumption (38,39). For instance, heritability estimates obtained from twin-adoption studies (where the MZ twins are raised in entirely different families) closely resemble those from ordinary twin studies. Also, studies of parents with misclassified twins (the parents always thought the twin was MZ but they turned out to be DZ and vice versa) have not shown any consistent effect of perceived zygosity on twin similarity for a range of personality traits.

Structural equation modeling (SEM) of twin variance-covariance data has several advantages over merely comparing the MZ and DZ correlations (34,40,41). SEM allows the comparison of the fit of alternative models (e.g., ACE versus AE) with the observed data and provides confidence intervals around the estimates for  $V_A$ ,  $V_C/V_D$ , and  $V_E$ . In SEM, the relationship between several latent unobserved and observed variables is summarized by a series of structural equations. In a genetic analysis, these equations relate the observed trait to latent genetic and environmental variables (i.e., the additive and dominant effects of genes and common and unique environmental influences). From these equations, it is possible to derive the variance-covariance matrix implied by the model through covariance algebra (42). The variances and covariances for the basic twin model can be represented by linear structural equations of the total phenotypic variance ( $V_P$ ) of both MZ and DZ twins ( $V_P = V_A + V_D + V_C + V_E$ ), the MZ covariance ( $\text{Cov}[\text{MZ}] = V_A + V_D + V_C$ ), and the DZ covariance ( $\text{Cov}[\text{DZ}] = 0.50V_A + 0.25V_D + V_C$ ). As stated earlier, since we have four unknowns and only three observations, at most only one of  $V_C$  and  $V_D$  can be estimated. This is not to say that  $V_C$  and  $V_D$  cannot both contribute to the phenotypic variance of a trait but rather they cannot be esti-

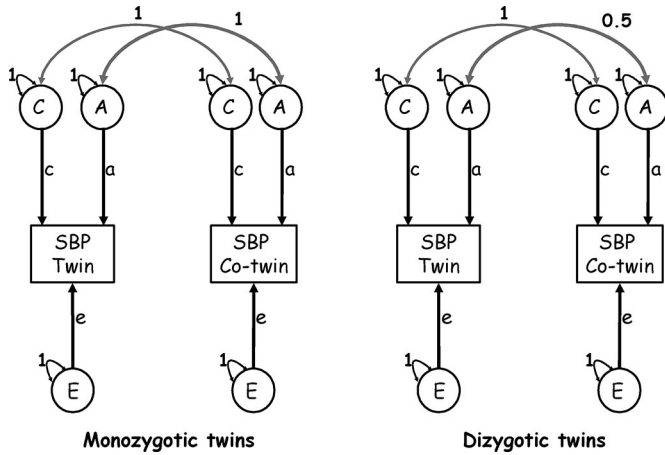


Figure 1. Variance decomposition in twin studies using a path model. Following standard path tracing rules, the expectation for the variance and covariances obtains as: Variance (SBP) =  $(a^2 * 1 * a) + (c^2 * 1 * c) + (e^2 * 1 * e) = a^2 + c^2 + e^2$ ; monozygotic covariance (SBP) =  $(a^2 * 1 * a) + (c^2 * 1 * c) = a^2 + c^2$ ; dizygotic covariance (SBP) =  $(a^2 * 0.5 * a) + (c^2 * 1 * c) = 0.5a^2 + c^2$ . Path coefficients  $a$ ,  $c$ ,  $e$  can be estimated by a maximum likelihood estimation procedure that optimally fits the expected variance and covariances to the observed variance in SBP and the observed cross-twin SBP covariances in MZ and DZ twin pairs. SBP = systolic blood pressure.

mated simultaneously with data from twins alone. Consequently, when the correlation between MZ twins is less than twice the DZ correlation, we estimate  $V_C$  and assume that genetic dominance is absent; conversely, when the MZ correlation is more than twice the DZ correlation, we estimate  $V_D$  and assume that  $V_C$  is zero.

Structural equation models may be represented diagrammatically using path diagrams, which can be helpful in understanding complex multivariate designs. A first simple univariate example relevant to psychosomatic medicine is depicted in Figure 1. SBP is measured at rest in MZ and DZ twin pairs. Our model specifies one latent genetic factor, one latent shared environmental factor, and one latent unique environmental factor, all with a variance of 1. In the example, dominance is assumed not to influence SBP and all the genetic variance is assumed to be additive; this seems to be the case in reality as well (17). Path coefficients “ $a$ ,” “ $c$ ,” and “ $e$ ” represent the factor loadings of SBP on the latent factors. As seen from biometrical theory,  $a^2 = V_A$ ,  $c^2 = V_C$ , and  $e^2 = V_E$  (43). In structural equation modeling, parameter estimates for these path coefficients are obtained by using a fitting function, which quantifies the difference between the observed variance-covariance matrix and the variance-covariance matrix implied by the model. These functions provide a measure of how likely the data are under the specified model for the causes of familial resemblance. They also provide the significance of each of the model parameters (e.g.,  $a^2$ ,  $c^2$ , and  $e^2$ ). The relative contribution of the genetic factor to the total variance in resting SBP, also known as the heritability ( $h^2$ ), now obtains as the ratio of  $a^2 / (a^2 + e^2 + c^2)$ .

One huge advantage of structural equation modeling is that it can easily be expanded to the multivariate case, enabling us to examine if two traits are correlated through common ge-

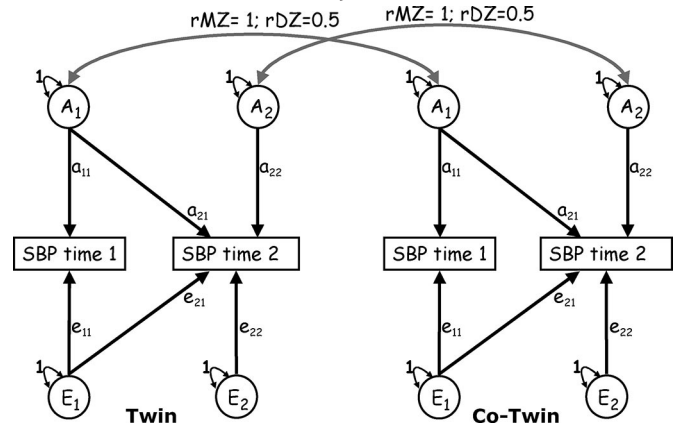


Figure 2. Bivariate twin model of the tracking of systolic blood pressure (SBP) across time. From the bivariate path diagram, we can compute the genetic correlation between the two time points ( $r_g$ ) as the genetic covariance divided by the square root of the genetic variances of both traits ( $a_{11} * a_{21} / \sqrt{(a_{11}^2 * (a_{21}^2 + a_{22}^2))}$ ). The unique environmental correlation obtains as  $e_{11} * e_{21} / \sqrt{(e_{11}^2 * (e_{21}^2 + e_{22}^2))}$ .

netic or through common environmental effects. A typical example in our field would be to detect the nature of the well-known tracking of SBP level across time, which in adulthood is about 0.55 over 5- to 10-year periods. This tracking may reflect the effects of an underlying genetic factor affecting SBP across time points, but it may also reflect the effects of chronic stress or other persistent unique environmental factor. Figure 2 depicts a bivariate twin model that can test this and various other hypotheses. In the example, we assume that only two sources of variance explain individual differences in SBP at the two time points: additive genetic and unique environmental factors; again this seems to be true in reality (23,44). If coefficient  $a_{22}$  can be set to zero without a significant loss of fit, only a single genetic factor influences SBP at both time points; i.e., there is no additional contribution of genetic factors at time 2 that is not already seen at time 1. If coefficient  $e_{21}$  can be set to zero, then the unique environmental factors causing variance in SBP at time points 1 and 2 are uncorrelated. If coefficient  $a_{21}$  is significant and  $e_{21}$  is not, this means that the tracking of SBP over time is caused entirely by underlying genetic factors. Such a structure was found across multiple time points in Dutch twin samples (23) whereas in Australian and American twins both genetic and environmental factors contributed to temporal stability of SBP (44,45).

Multivariate structural equation models of twin data can also be used to analyze the interaction between siblings, the genetic and environmental correlation between different traits, and the direction of causation between variables (34,43,46). It is also easy to extend the classical twin design by including other informative relationships in the analysis including siblings (47,48), parents of twins (25,49,50), the offspring of MZ and DZ twins (51,52), and the spouses of twins (53–55). These designs can quantify the effects of phenotypic assortment and vertical cultural transmission, which the classical twin study cannot do. Finally, if important aspects of the environment are



GENETICALLY INFORMATIVE DESIGNS

measured, the presence and extent of gene × environment interaction can be tested (36,37).

An example of a twin model incorporating gene-environment interaction is given in Figure 3, where we additionally control for possible gene-environment correlation. Regular exercise is known to be associated with lowered SBP (56). However, the extent of SBP reduction after an identical exercise program shows large differences between individuals; family studies (98,99) have suggested these differences to be partly heritable (57). This suggests that subjects with different genetic make-up can differ in their sensitivity to the beneficial effects of exercise. To account for this gene-exercise interaction, the path loadings on SBP in Figure 3 are weighted for the exercise status (which is “yes” = 1/“no” = 0) of the twins. If a model with nonzero β weights for the genetic factors fits the observed data better than a model with zero β weights for the genetic factors, we have formal evidence of gene-exercise interaction. Some complexity is introduced to the model by allowing part of the association between exercise behavior and blood pressure to derive from genes that independently influence both traits (i.e., the latent genetic factor  $A_c$ ). This phenomenon is known as “pleiotropy” and may play a role in many traits that can be considered “environmental” modifiers of risk factors, on the one hand (e.g., lifestyle, SES, chronic stress), but may themselves be heritable. When there is evidence of potential gene-environment correlation, i.e., when the “environmental factor” itself shows heritability, as is the case for exercise behavior (35), allowing for gene-environment correlation as in Figure 3 is prudent.

In short, twin studies provide a first necessary step in genetic research by establishing that genes contribute to the observed population variation in psychosomatic risk factors and by estimating the size of this genetic contribution relative to other factors that create resemblance within families. Twin

studies do not identify the actual genes. This effort requires molecular genetic research on the actual genetic variation.

Molecular Genetics

That DZ twins share, on average, 50% of their genetic material refers exclusively to the part of the genes in which people can differ. Any one person’s deoxyribonucleic acid (DNA) is 99.9% the same as any other person’s DNA. The 0.01% difference in the sequence of DNA among individuals is the source of all genetic variation. Variation in a single gene is responsible for some disorders, such as cystic fibrosis and sickle cell disease. Variation in multiple genes, environmental factors, gene by gene interactions, and gene by environment interactions are thought to account for complex traits, including most traits of interest in psychosomatic medicine.

A gene consists of two units of information, the alleles. One allele is inherited from the father and one from the mother. Together they constitute the genotype, which may be homozygous (same allele from both parents) or heterozygous (different allele from each of the parents). Under a simple Mendelian inheritance model and random mating assumption, lack of selection according to genotype, and absence of mutation or migration, the frequencies of the genotypes in the population are perfectly predicted by the frequencies of the two alleles, which is referred to as Hardy-Weinberg equilibrium (HWE) (58). As an example, consider a gene with two alleles, denoted “short” (s) with frequency p and “long” (l) with frequency q. Let the least frequent, or minor, allele s take up 40% of all alleles in the population ( $p = .4$ ). The three potential genotypes, ss, ls and ll, have expected frequencies, namely,  $p^2$  (.16),  $2pq$  (0.48), and  $q^2$  (0.36). A  $\chi^2$  test for HWE compares these expected genotype frequencies with the observed genotype frequencies; a significant  $\chi^2$  test indicates that HWE does not hold. Many of the association analyses discussed below require HWE to hold.

Large-scale genetic variation includes loss or gain of chromosomes or breakage and rejoining of chromatids. This variation is abnormal and often leads to profound developmental problems. Smaller-scale genetic variation is at the level of a single allele and contributes to most of the normal variation in the population. Smaller-scale genetic variation can be classified into three groups: a) single nucleotide polymorphisms (SNPs), b) insertion/deletion polymorphisms, and c) varying number of tandem repeats (VNTR). Deletion occurs when one or more nucleotides are eliminated from a sequence, whereas insertion occurs when one or more nucleotides are inserted into the sequence. VNTRs (which include very short repeats or microsatellites) are short identical segments of DNA aligned head to tail in a repeating fashion. The number of repeated segments at a locus varies between individuals. An SNP is defined as a single base substitution. SNPs are the most abundant form of DNA variation in the human genome with approximately 7 million common SNPs with a minor allele frequency of at least 5% across the entire human population (59–62).

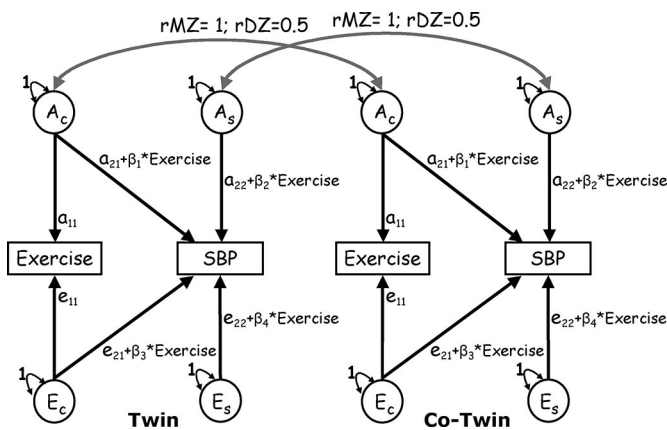


Figure 3. Twin model testing gene-environment interaction in the presence of gene-environment correlation.  $A_c$  is a hypothetical set of pleiotropic genes that lower SBP and increase the drive to exercise regularly.  $A_s$  are genes specific to SBP. Regular exercise is simultaneously allowed to act as an environmental modulator of the genetic effects on SBP ( $\beta_1, \beta_2$ ) as well as of the unique environmental effects on SBP ( $\beta_3, \beta_4$ ). Significance of the gene-exercise interactions can be tested by comparing this full model to models with  $\beta_1$  or  $\beta_2$  set to zero. SBP = systolic blood pressure.

### Candidate Gene Association Studies

Candidate gene association studies test if a particular allele in a candidate gene and a trait co-occur above chance level, given the frequency of the allele and the distribution of the trait in the population (63). In these studies, selection of candidate genes a priori is required. The selection of genes may be based on the biological role of the gene in a causative pathway (physiological candidate) or a location close to a peak from a linkage, or genetic mapping, study (positional candidate). Ideally, the gene fits both criteria. In a direct association study, one or more putatively functional variants are genotyped and serve as the independent variable predicting a dependent variable, the trait of interest. It is presumed that the selected variant is causative in the trait of interest although, in practice, association may be attributable to linkage disequilibrium (LD) with another functional site nearby. Genetic variants should be prioritized by apparent functional significance or location within coding, promoter, or splice regions. These typically include SNPs, VNTRs, and insertion/deletion polymorphisms.

An example of a direct association study is examining the role of variants within  $\alpha$ - and  $\beta$ -adrenergic receptor genes as predictors of blood pressure level. The adrenergic receptor genes are good biological candidates due to their location on the heart ( $\beta_1$ ), in the vasculature ( $\alpha_1$ ,  $\beta_2$ ), or within the central nervous system ( $\alpha_{2a}$ ), and their involvement in cardiovascular regulation. In addition, individual variants within the genes have been shown to be functional. For example, receptors with the C  $\rightarrow$  G SNP at base pair (bp) 1165 within ADRB1 ( $\beta_1$ -adrenoreceptor gene), resulting in an amino acid substitution from arginine to glycine at position 389, show increased adenylyl cyclase activity in the presence of an agonist. In a study of young adult twins, the genotypes at this SNP were examined in relation to blood pressure at rest and in response to a combined mental arithmetic and Stroop task (3). After statistically controlling for age, sex, and body mass index, participants carrying any G allele at base pair (bp) 1165 in ADRB1 exhibited increased resting SBP (GG/GC =  $115.52 \pm 8.47$  versus CC =  $112.94 \pm 10.14$  mm Hg), DBP (GG/GC =  $61.88 \pm 6.32$  versus CC =  $59.64 \pm 7.16$  mm Hg), and a larger DBP response (GG/GC =  $6.97 \pm 6.94$  versus CC =  $4.85 \pm 6.87$   $\Delta$ mm Hg) to mental challenge as compared with CC genotypes (CG and GG groups were combined due to the small sample size for GG homozygotes).

There are numerous online resources with information about candidate genes and variation in or near the genes of interest. The National Center for Biotechnology Information home page, available at <http://www.ncbi.nlm.nih.gov/>, includes resources such as Online Medelian Inheritance in Man (OMIM), dbSNP, the Genome Database and Pubmed. Other excellent resources include Ensembl available at <http://www.ensembl.org/>, the Genome Browser from the University of California, Santa Cruz available at <http://genome.ucsc.edu/>, the International Hapmap project available at <http://www.hapmap.org/>, the SNP Consor-

tium (TSC) available at <http://snp.cshl.org/>, the SeattleSNPs variation discovery resource available at <http://pga.gs.washington.edu/>, and SNPper, a Web-based application to automate the tasks of extracting SNPs from public databases available at <http://SNPper.chip.org/>.

The statistical approach to association studies depends on the research design (63). Common research designs for association studies include cohort designs and case-control designs. Within these designs, special topics with statistical implications include treatment of race and ethnicity, gene  $\times$  gene and gene  $\times$  environment interaction, and use of haplotypes and power.

### Cohort Studies: Quantitative Traits

Single diallelic polymorphisms, such as SNPs, may be analyzed using general linear modeling. For individual SNPs, genotype (e.g., GG, CG, CC for a G to C substitution) typically serves as the independent variable. In the absence of knowledge about whether alleles at a given site function in an additive, dominant, or recessive manner (as is the case for many of the polymorphisms of interest in psychosomatic medicine), the three possible genotypes should be treated as independent groups. This would translate to a between-subjects group factor with the number of levels ( $k$ ) equal to the number of genotypes and  $k-1$  degrees of freedom ( $df$ ) (i.e., 2  $df$ ). Evidence for apparent dominance of one allele over another may be detected through posthoc group contrasts (e.g., GG = CG > CC). Covariates and additional predictors of the dependent variable may also be incorporated.

Within a regression framework, the most general model for genetic effects at a single locus includes a term for linear effects of a given allele and an additional parameter for the deviation from this linear effect, i.e., a dominance term (63). For the linear term, genotypes (e.g., GG, CG, and CC) are assumed to function in an additive manner and are coded as 0, 1, and 2, reflecting dose of the C allele. The associated  $\beta$  weight is the additive effect of the C allele. This linear model alone predicts that the mean of the heterozygotes (CG) will be located at the midpoint between the two types of homozygotes (GG, CC); however, in practice, this may or may not be the case. Deviation of the mean of the heterozygotes from the midpoint between the means of the homozygotes suggests that one allele is dominant over the other. To quantify this effect, an additional, dominance term, is necessary. Specifically, genotypes GG, CG and CC may be coded 0, 1, and 0 with the associated  $\beta$  weight reflecting deviation of the heterozygotes from the midpoint of the two homozygous groups, as would be predicted by the linear term alone. The general regression framework for a diallelic locus is given by  $P = \alpha + \beta_a A + \beta_d D + e$ , where  $P$  is a quantitative trait;  $\alpha$  is the baseline mean of  $P$ ;  $A$  and  $D$  are dummy variables reflecting coding for linear and nonlinear effects of the underlying genotype at a single locus; and  $e$  is a residual error term assumed to be normally distributed.

For polymorphisms with more than two alleles (e.g., microsatellites), genotypes may be treated individually, although

## GENETICALLY INFORMATIVE DESIGNS

there will be little power to examine the effects of the more rare alleles. Alternatively, alleles may be ranked according to function based on *in vitro* assays (64). If there are no functional data available and several rare genotypes, it may be necessary to limit analyses to the most common genotypes to preserve statistical power.

### *Case-Control Studies: Disease Traits*

Case-control genetic association studies are typically comprised of a group of cases with a trait of interest and well-matched controls. Ideally, the cases and controls should represent “identical” subsamples from a single population differing only on the trait of interest (65). Statistical analyses compare allele frequencies or genotypes across cases and controls. In well-matched samples, differences in genotypes across cases and controls may be tested using  $\chi^2$  tests. Alternatively, the risk of having the disorder may be modeled using logistic regression with a 2 *df* test. Within this approach, the log odds of expressing the disease trait is modeled as a function of the additive effects of the dose of one of the alleles (e.g., 0, 1, or 2 copies of the C allele for genotypes GG, CG, and CC, respectively) and a dominance term representing deviance from this additive pattern (e.g., genotypes GG, CG, and CC coded as 0, 1, and 0). For the additive term, the log odds of disease expression for heterozygotes is midway between the log odds of the two homozygous groups. The dominance term quantifies the extent to which the log odds for heterozygotes differs from the additive prediction. The general logistic regression framework for a diallelic locus is given by  $\ln(P/1 - P) = \alpha + \beta_a A + \beta_d D + e$ , where P is the binary expression of a phenotype;  $\alpha$  is the baseline log odds of P; A and D are dummy variables reflecting coding for linear and nonlinear effects of the underlying genotype at a single locus; and *e* is a residual error term assumed to be normally distributed. The natural log raised to the power of the additive  $\beta$  weight ( $e^{\beta_a}$  or  $\text{Exp}(\beta_a)$ ) reflects the change in odds of expression of the phenotype based on a unit increase in allele dose. The GG genotype becomes the reference group (0 allele) and the effect of genotype is quantified by determining if there is a significant change in the probability of the expression of the phenotype for each additional C allele (CG = 1 additional allele and CC genotypes = 2 additional alleles). The natural log raised to the power of the dominance  $\beta$  weight reflects the deviation of heterozygotes from the midpoint of the log odds for the two homozygous groups. For a binary genotype (i.e., GG versus CG or CC), natural log raised to the power of the additive  $\beta$  weight would be an odds ratio.

### *Treatment of Race and Ethnicity*

In cohort-based and case-control analysis of unrelated individuals, spurious genetic association may result due to differences in allele frequencies and the trait of interest in subgroups within the larger population, often reflecting racial or ethnic groups (population stratification). The classic example of population stratification is a hypothetical association between chopstick use and any genetic marker that differs markedly

between Asian and Caucasian populations in a larger population with substantial representation of both ethnicities, such as San Francisco, California (66). It has been argued that there have been relatively few documented instances of bias due to population stratification reported in the literature and that population-based studies are largely robust to this type of bias (67). However, recent empirical tests do find evidence of stratification effects, particularly among populations that have recently been mixed from two or more distinct parental populations (genetic admixture), including African Americans and Hispanic Americans (68).

Population stratification is essentially a problem of sample matching, occurring primarily when the genetic background of the cases differs from that of controls (67). Accordingly, it is possible that matching cases and controls on self-reported race in homogeneous populations (such as European Americans) will mitigate concerns about population stratification. Two methods are available to control for stratification using markers throughout the genome. In structure assessment (69–74), genetic markers, either anonymous markers or markers that differ substantially among ethnic groups, are used to predict membership in homogeneous subgroups within a stratified population. Once identified, genetic associations may be conducted within these subgroups to ensure a similar genetic background of cases and controls. A second method, genomic control (75–79), uses anonymous genetic markers to estimate the degree of inflation of the  $\chi^2$  statistic due to population stratification and yields a correction factor to account for these background genetic effects in genetic association studies. With the rapid reduction in genotyping costs and further development of these methods (69), it is likely that the threat of population stratification will be routinely controlled in cohort and case-control genetic association studies using these types of techniques.

Another good method to ensure genetic matching is to conduct genetic studies within families. Tests using within-family controls to control for population stratification are collectively known as transmission disequilibrium tests (TDTs). The classic TDT requires information on trios, i.e., parents and an affected offspring. The principal idea is that the allele associated with disease will be transmitted more often to an affected offspring (80). The TDT compares the actual and expected probabilities of transmission of the allele (an offspring has an expected chance of 0.5 of receiving a specific allele from either the mother or the father). Overtransmission can only occur if the marker and disease locus are linked. However, power of the TDT is less than for an association test based on cases and controls because only heterozygote parents provide information about preferential allele transmission. After the introduction of the classic TDT by Spielman and colleagues (80), the TDT has undergone many developments and has, for example, been adapted for quantitative traits and nuclear families of any size (81–83) as well as for haplotypes (84).



**Gene × Gene and Gene × Environment Interaction**

From a genetics perspective, nearly all psychosomatic traits are considered “complex,” meaning that the causal pathways are likely to involve multiple genes of small effect, environmental factors, and gene × gene and gene × environment interaction (85). Genetic interaction within a given locus is termed genetic dominance. Interaction between two loci is termed epistasis. However, a distinction between epistasis referring to a statistical interaction and that referring to a physical interaction of gene products is warranted, as the presence of statistical interaction does not necessarily imply an underlying biological interaction (86,87). Similarly, statistical gene-environment interactions should be interpreted with caution as the mathematical model may again have no obvious biological interpretation (88).

Modeling statistical gene × gene or gene × environment interaction may be accomplished by incorporating two genetic predictors or one genetic and one environmental predictor into linear or logistic regression in standard statistical packages and testing for their interaction (87,89). The choice of scale becomes important because factors that are additive with respect to an outcome in one scale may exhibit interaction if a transformed scale is used. For linear regression with two genetic predictors, the general regression model is given by:

$$P = \alpha + \beta_{a1}A_1 + \beta_{d1}D_1 + \beta_{a2}A_2 + \beta_{d2}D_2 + \beta_{a1a2}A_1A_2 + \beta_{a1d2}A_1D_2 + \beta_{a2d1}A_2D_1 + \beta_{d1d2}D_1D_2 + e$$

where P is a quantitative trait; α is the baseline mean of P; A<sub>1</sub>, A<sub>2</sub>, D<sub>1</sub>, and D<sub>2</sub> are dummy variables coding for the additive and dominance effects of the underlying genotype for sites 1 and 2; and e is a residual error term assumed to be normally distributed. Statistical epistasis implies that at least one of the interaction coefficients differs significantly from zero.

For gene × environment interaction, at least one genetic and one environmental predictor are included in the regression equation plus the interaction of the additive and dominance term with the environmental predictor. Statistical interaction implies that either of the interaction terms differs significantly from zero. The general regression framework for a gene × environment interaction for a continuous trait is given by:

$$P = \alpha + \beta_a A + \beta_d D + \beta_e E + \beta_{ae} AE + \beta_{de} DE + e$$

where P is a quantitative trait; α is the baseline mean of P; A and D are dummy variables coding for linear and nonlinear effects of the underlying genotype; E is a measured environmental factor; and e is a residual error term assumed to be normally distributed. Assuming no genetic dominance or associated interactions, this equation reduces to:

$$P = \alpha + \beta_a A + \beta_e E + \beta_{ae} AE + e$$

Finally, if the genotype is correlated to the environmental risk factor (e.g., genetic susceptibility to aggression and pa-

**HAPLOTYPES:**

snp1:	A	snp1:	A	snp1:	A	snp1:	A	snp1:	T	snp1:	T	snp1:	T	snp1:	T
snp2:	T	snp2:	T	snp2:	A	snp2:	A	snp2:	T	snp2:	T	snp2:	A	snp2:	A
snp3:	C	snp3:	G	snp3:	C	snp3:	G	snp3:	C	snp3:	G	snp3:	C	snp3:	G

**GENOTYPES:**

		mother							
		ATC	ATG	AAC	AAG	TTC	TTG	TAC	TAG
father	ATC	<b>AATTCC</b>	AATTGC	AAATCC	AAATGC	TATTCC	TATTGC	TAATCC	TAATGC
	ATG	<b>AATTCC</b>	<b>AATTGG</b>	AAATCG	AAATGG	TATTCC	TATTGG	TATACG	TAATGG
	AAC	<b>AATACC</b>	AATAGC	<b>AAAACC</b>	AAAAGC	TATACC	TATAGC	TAAACC	TAAAGC
	AAG	<b>AATACC</b>	<b>AATAGG</b>	<b>AAAACG</b>	<b>AAAAGG</b>	TATACC	TATAGG	TAAACG	TAAAGG
	TTC	ATTTCC	ATTTGC	ATATCC	ATATGC	<b>TTTTCC</b>	TTTTGC	TTATCC	TTATGC
	TTG	ATTTCC	<b>ATTTGG</b>	ATATCG	ATATGG	<b>TTTTCG</b>	<b>TTTTGG</b>	TTATCG	TTATGG
	TAC	ATTACC	ATTAGC	<b>ATAACC</b>	ATAAGC	<b>TTTACC</b>	TTTAGC	<b>TTAACC</b>	TTAAGC
	TAG	ATTACC	ATTAGG	ATAACG	<b>ATAAGG</b>	TTTACG	<b>TTTAGG</b>	<b>TTAACG</b>	<b>TTAAGG</b>

Figure 4. Three-SNP haplotypes and genotypes. The upper part of the figure shows the eight haplotypes that can be formed from three SNPs that are in full linkage equilibrium (the general rule is 2<sup>n</sup> with n = the number of SNPs). These eight haplotypes give rise to the 36 different genotypes (the general rule is ((m + 1)\*m)/2 with m the number of haplotypes). Because genotyping does not discriminate between paternal or maternal alleles, the upper part of the matrix (blue) has identical genotypes as the lower part of the matrix. The genotypes printed boldface unambiguously translate to haplotypes. All other genotypes printed in black can derive from multiple combinations of haplotypes. If the SNPs are in linkage disequilibrium, only a few of these possible haplotypes and genotypes will be observed in the population. SNP = single nucleotide polymorphism.

rental maltreatment), the interpretation of the statistical interaction is not straightforward (90). In addition, observational studies can be associated with substantially less power than well-designed experiments to detect interaction effects (91), suggesting that controlled interventions may be a useful alternative to observational studies in detecting gene × environment interaction effects. An example would be to test whether certain candidate genes in the sympathetic nervous system (e.g., ADRB2, or the β<sub>2</sub>-adrenoreceptor gene) may explain part of the large individual variability in the beneficial effects of exercise on blood pressure.

**Haplotype Analysis**

The primary disadvantage of characterizing a single variant per gene is that there may be additional variants within the gene that are relevant to the trait of interest but are not captured by variation at a single marker. Hence, there has been increasing interest in using haplotypes, rather than single markers, as the unit of analysis in association studies (92). A haplotype refers to multiple SNPs along a short region of a chromosome (e.g., within a gene) that occur in a block pattern (Figure 4). There are three good reasons to perform haplotype analysis as part of candidate gene association studies: a) a haplotype might be in higher LD with the causal locus than any of the individual markers, b) interactions among the individual markers might form a functional haplotype, and c) haplotype analysis reduces the number of multiple tests of individual SNP analysis. A common problem of all statistical methods that use haplotype information is linkage phase ambiguity; i.e., it is unknown which alleles are located on the maternal chromosome and which are located on the paternal

## GENETICALLY INFORMATIVE DESIGNS

chromosome. As our genotyping analyses yield only the full genotypes, not the parental alleles separately, we do not know from which haplotype (maternal or paternal) the alleles originated. When multiple members in a family are genotyped, preferably including the parents, the two haplotypes constituting each genotype can often be determined from the Mendelian principles of gene segregation within a pedigree. Alternatively, statistical algorithms can be used to reconstruct haplotypes in unrelated individuals using the frequency and correlation of the SNPs in the population. The reliability of such algorithms seems to be good for multiple diallelic markers, such as SNPs (93,94), although there is some power loss for the association tests as a result of the haplotype phase uncertainty.

Because the SNPs in a haplotype are strongly associated (in LD) with each other, it is possible to test for the association of a haplotype with a trait or disease by genotyping only a few SNPs ("haplotype tagging" SNPs or simply "tagging" SNPs) within the haplotype. Tagging SNPs are first selected in a subset of the sample or in samples of the same ethnicity from freely available web resources, such as the HapMap. This reduces the cost of genotyping in the full sample, yet it ensures reasonably good coverage of common variation throughout the gene. The tagging SNPs are then examined for association with the trait of interest in the total sample and the effects of unassayed SNPs would then be detected through LD with tagging SNPs (95). The International Hapmap project (available at <http://www.hapmap.org>) has characterized >4 million SNP markers on a genome-wide scale in three ethnic groups (Caucasians, Africans, and Asians), greatly facilitating the use of tagging SNPs in association studies.

### *Power and Sample Size Considerations*

Although many power calculations required for genetic association studies may be derived from texts well known to behavioral researchers (96,97), excellent on-line resources specific to power and sample size calculations for genetic association studies also exist, e.g., Quanto (98) available at <http://hydra.usc.edu/gxe> and the Genetic Power Calculator (99) available at <http://statgen.iop.kcl.ac.uk/gpc/>. In molecular genetic studies of quantitative traits, assuming a simple additive model, the effect size of a locus is a function of mean trait differences between homozygotes (e.g., the CC versus GG genotype) and allele frequency (100). Differing modes of inheritance (additive, dominant, and recessive) will also influence the effect size and have resulting effects on power and sample size.

As psychosomatic traits are likely to be influenced by multiple genes and interactions of small effect, the effect size for each is generally expected to be small. Sample sizes required to detect gene main effects and gene  $\times$  environment interaction with sufficient statistical power in this context are relatively large. Although previous studies (98,99) have suggested that association can be detected even in modestly sized samples, standard power calculations show that up to 1000 participants are required to detect gene main effects and

approximately 1500 to 2000 participants are required to detect gene  $\times$  environment interaction with small to medium effect sizes. The required sample size will be even larger if one of the alleles is rare (e.g., less than 5% to 10%) or a large number of markers is typed and the statistical criterion, typically set at  $\alpha = 0.05$  for two-tailed tests, must be adjusted for multiple comparisons. One method to adjust for multiple comparisons is to use techniques that control the false discovery rate (FDR), i.e., the proportion of significant findings (or discoveries) that are false-positives (101,102).

### *Integration*

Although statistical analysis of genetic association may, in many cases, be conducted using well-known methods, the strength of the interpretation of results is grounded in the study design. Nonsignificant results may be attributable to Type II error (at 80% power, there remains a 20% chance that you will falsely accept the null hypothesis) or experimental biases like genotyping error or overmatching of controls (65). Many prior studies have lacked sufficient statistical power (at least 80%) to detect the small effects expected, particularly if gene-gene, gene-environment, or genetic heterogeneity (i.e., more than one genetic variant can produce the same outcome) effects are involved. In addition, negative results could be due to inadequate coverage of a gene, for example, in studies of single variants. Nonsignificant results may also be attributable to a true lack of etiological relationship (65). Given the importance of nonreplications in the literature, calls have been made for convenient formats to publish negative results (65,89,103).

A significant result may indicate that a causal relationship between genotype and trait has been identified. However, because there are several other potential explanations of significant results, this type of interpretation should be used with caution. A common cause of false-positive results is increased Type I error due to multiple statistical tests. This problem will only increase with the availability of high throughput methods, which can easily generate millions of genotypes. The optimal correction method for multiple comparisons depends on the number of markers and phenotypes studied. For example, correction is mandatory for whole genome studies that use large numbers of random markers but may not always be necessary for candidate gene studies in which the prior probability of a true discovery is likely to be higher (104) and in which gene-wide significance levels can be used (105). The issue is further complicated by the correlation between SNPs (LD) and between phenotypes, making it difficult to assess the number of independent tests. Recently Manly and colleagues (104) showed that correction techniques for multiple comparisons based on the original Bonferroni are generally too conservative. New procedures based on FDR effectively control the proportion of false discoveries without sacrificing the power to discover.

Population stratification is another source of false-positive results. Although it is less likely that spurious association due



to population stratification will occur among seemingly homogeneous populations, such as European Americans, this type of bias remains a concern particularly among recently admixed populations, such as African Americans and Hispanic Americans, and populations of mixed racial and ethnic composition. It is also possible that the genetic marker showing significant association is not the causal variant per se but is co-inherited (in LD) with a causal variant. Many of these threats to the interpretation of the results may be mitigated with careful study design, such as appropriate correction for multiple comparisons, incorporation of genetic markers to characterize population substructure, and haplotyping to characterize variation throughout a candidate gene. Most importantly, to minimize the probability that an observed association is a false-positive, significant findings must be replicated in independent samples. Many of these issues may be novel for persons considering genetic research for the first time and consultation on study designs with geneticists, statistical geneticists, or genetic epidemiologists is always recommended.

Finally, although this review focused on methods for candidate gene association studies, it should be noted that for most complex traits, our knowledge of underlying causative pathways is likely incomplete. Limiting the search for contributing genetic variation to known candidate genes only will likely prevent the identification of potentially novel pathways that contribute to psychosomatic traits. Thus, the candidate gene approach should ideally capitalize on knowledge generated with genome-wide searches, using techniques such as linkage analysis (106) and genome-wide association (107).

REFERENCES

1. Jeanmonod P, von Kanel R, Maly FE, Fischer JE. Elevated plasma C-reactive protein in chronically distressed subjects who carry the A allele of the TNF-alpha-308 G/A polymorphism. *Psychosom Med* 2004; 66:501-6.
2. McCaffery JM, Bleil M, Pogue-Geile MF, Ferrell RE, Manuck SB. Allelic variation in the serotonin transporter gene-linked polymorphic region (5-HTTLPR) and cardiovascular reactivity in young adult male and female twins of European-American descent. *Psychosom Med* 2003;65:721-8.
3. McCaffery JM, Pogue-Geile M, Ferrell R, Petro N, Manuck SB. Variability within alpha- and beta-adrenoreceptor genes as predictors of cardiovascular function at rest and in response to mental challenge. *J Hypertens* 2002;20: 1105-14.
4. Raynor DA, Pogue-Geile MF, Kamarck TW, McCaffery JM, Manuck SB. Covariation of psychosocial characteristics associated with cardiovascular disease: genetic and environmental influences. *Psychosom Med* 2002;64:191-203; discussion 204-5.
5. Scherrer JF, Xian, Hong, Bucholz, Kathleen K, Eisen, Seth A, Lyons, Michael J, Goldberg, Jack, Tsuang, Ming, True, William R. A twin study of depression symptoms, hypertension, and heart disease in middle-aged men. *Psychosom Med* 2003;65:548-57.
6. Wang X, Trivedi R, Treiber F, Snieder H. Genetic and environmental influences on anger expression, John Henryism, and stressful life events: the Georgia cardiovascular twin study. *Psychosom Med* 2005;67:16-23.
7. Bak S, Gaist D, Sindrup SH, Skytthe A, Christensen K. Genetic liability in stroke: a long-term follow-up study of Danish twins. *Stroke* 2002; 33:769-74.
8. Flossmann E, Schulz UG, Rothwell PM. Systematic review of methods and results of studies of the genetic epidemiology of ischemic stroke. *Stroke* 2004;35:212-27.
9. Swan L, Birnie DH, Inglis G, Connell JM, Hillis WS. The determination

of carotid intima medial thickness in adults—a population-based twin study. *Atherosclerosis* 2003;166:137-41.

10. Zdravkovic S, Wienke A, Pedersen NL, Marenberg ME, Yashin AI, De Faire U. Heritability of death from coronary heart disease: a 36-year follow-up of 20 966 Swedish twins. *J Intern Med* 2002;252:247-54.
11. Zdravkovic S, Wienke A, Pedersen NL, Marenberg ME, Yashin AI, de Faire U. Genetic influences on CHD-death and the impact of known risk factors: comparison of two frailty models. *Behav Genet* 2004;34: 585-92.
12. Beekman M, Heijmans BT, Martin NG, Pedersen NL, Whitfield JB, DeFaire U, van Baal GC, Snieder H, Vogler GP, Slagboom PE, Boomsma DI. Heritabilities of apolipoprotein and lipid levels in three countries. *Twin Res* 2002;5:87-97.
13. de Lange M, Snieder H, Ariens RA, Spector TD, Grant PJ. The genetics of haemostasis: a twin study. *Lancet* 2001;357:101-5.
14. de Maat MP, Bladbjerg EM, Hjelmberg JB, Bathum L, Jespersen J, Christensen K. Genetic influence on inflammation variables in the elderly. *Arterioscler Thromb Vasc Biol* 2004;24:2168-73.
15. Dunn EJ, Ariens RA, de Lange M, Snieder H, Turney JH, Spector TD, Grant PJ. Genetics of fibrin clot structure: a twin study. *Blood* 2004; 103:1735-40.
16. Edwards KL, Newman B, Mayer E, Selby JV, Krauss RM, Austin MA. Heritability of factors of the insulin resistance syndrome in women twins. *Genet Epidemiol* 1997;14:241-53.
17. Evans A, Van Baal GC, McCarron P, DeLange M, Soerensen TI, De Geus EJ, Kyvik K, Pedersen NL, Spector TD, Andrew T, Patterson C, Whitfield JB, Zhu G, Martin NG, Kaprio J, Boomsma DI. The genetics of coronary heart disease: the contribution of twin studies. *Twin Res* 2003;6:432-41.
18. Kupper N, Willemsen G, Riese H, Posthuma D, Boomsma DI, de Geus EJ. Heritability of daytime ambulatory blood pressure in an extended twin design. *Hypertension* 2005;45:80-5.
19. Kupper NH, Willemsen G, van den Berg M, de Boer D, Posthuma D, Boomsma DI, de Geus EJ. Heritability of ambulatory heart rate variability. *Circulation* 2004;110:2792-6.
20. Peetz D, Victor A, Adams P, Erbes H, Hafner G, Lackner KJ, Hoehler T. Genetic and environmental influences on the fibrinolytic system: a twin study. *Thromb Haemost* 2004;92:344-51.
21. Retterstol L, Eikvar L, Berg K. A twin study of C-reactive protein compared to other risk factors for coronary heart disease. *Atherosclerosis* 2003;169:279-82.
22. Colletto GM, Cardon LR, Fulker DW. A genetic and environmental time series analysis of blood pressure in male twins. *Genet Epidemiol* 1993;10:533-8.
23. Hottenga JJ, Boomsma DI, Kupper N, Posthuma D, Snieder H, Willemsen G, de Geus EJ. Heritability and stability of resting blood pressure. *Twin Res Hum Genet* 2005;8:499-508.
24. Iliadou A, Lichtenstein P, Morgenstern R, Forsberg L, Svensson R, de Faire U, Martin NG, Pedersen NL. Repeated blood pressure measurements in a sample of Swedish twins: heritabilities and associations with polymorphisms in the renin-angiotensin-aldosterone system. *J Hypertens* 2002;20:1543-50.
25. Snieder H, van Doornen L, Boomsma DI. Developmental genetic trends in blood pressure levels and blood pressure reactivity to stress. In: Turner J, Cardon L, Hewitt J, editors. *Behavior genetic approaches in behavioral medicine*. New York: Plenum Press; 1995.
26. de Lange M, Snieder H, Ariens RA, Andrew T, Grant PJ, Spector TD. The relation between insulin resistance and hemostasis: pleiotropic genes and common environment. *Twin Res* 2003;6:152-61.
27. Hong Y, Pedersen NL, Brismar K, de Faire U. Genetic and environmental architecture of the features of the insulin-resistance syndrome. *Am J Hum Genet* 1997;60:143-52.
28. Poulsen P, Vaag A, Kyvik K, Beck-Nielsen H. Genetic versus environmental aetiology of the metabolic syndrome among male and female twins. *Diabetologia* 2001;44:537-43.
29. Kissebah AH, Sonnenberg GE, Mykalebust J, Goldstein M, Broman K, James RG, Marks JA, Krakower GR, Jacob HJ, Weber J, Martin L, Blangero J, Comuzzie AG. Quantitative trait loci on chromosomes 3 and 17 influence phenotypes of the metabolic syndrome. *Proc Natl Acad Sci USA* 2000;97:14478-83.
30. McCaffery JM, De Geus E, Snieder H, Dong Y, Flory JD. *Genetic epidemiology of psychosomatic traits*. Orlando, FL: American Psychosomatic Society; 2004.

## GENETICALLY INFORMATIVE DESIGNS

31. McCaffery JM, De Geus E, Snieder H, Dong Y, Flory JD. Genetics in the psychophysiological laboratory. Sante Fe, NM: Society for Psychophysiological Research; 2004.
32. Martin N, Boomsma D, Machin G. A twin-pronged attack on complex traits. *Nat Genet* 1997;17:387–92.
33. Mather K, Jinks JL. Biometrical genetics. London: Chapman & Hall; 1971.
34. Neale MC, Cardon LR. Methodology for genetic studies of twins and families. Dordrecht, The Netherlands: Kluwer Academic Publishers; 1992.
35. Stubbe JH, Boomsma DI, De Geus EJ. Sports participation during adolescence: a shift from environmental to genetic factors. *Med Sci Sports Exerc* 2005;37:563–70.
36. Purcell S. Variance components models for gene-environment interaction in twin analysis. *Twin Res* 2002;5:554–71.
37. Purcell S, Koenen KC. Environmental mediation and the twin design. *Behav Genet* 2005;35:491–8.
38. Bouchard TJ Jr. Genes, environment, and personality. *Science* 1994; 264:1700–1.
39. Kendler KS, Neale MC, Kessler RC, Heath AC, Eaves LJ. A test of the equal-environment assumption in twin studies of psychiatric illness. *Behav Genet* 1993;23:21–7.
40. Eaves L. The genetic analysis of continuous variation: a comparison of experimental designs applicable to human data. *Br J Math Stat Psychol* 1969;22:131–47.
41. Jinks JL, Fulker DW. Comparison of the biometrical genetical, MAVA, and classical approaches to the analysis of human behavior. *Psychol Bull* 1970;73:311–49.
42. Bollen K. Structural equations with latent variables. New York: John Wiley; 1989.
43. Posthuma D, Beem AL, de Geus EJ, van Baal GC, van Hjelmborg JB, Iachine I, Boomsma DI. Theory and practice in quantitative genetics. *Twin Res* 2003;6:361–76.
44. Hottenga JJ, Whitfield JB, de Geus EJ, Boomsma DI, Martin NG. Heritability and stability of resting blood pressure in Australian twins. *Twin Res Hum Genet* 2006;9:205–9.
45. Kupper N, Ge D, Treiber FA, Snieder H. Emergence of novel genetic effects on blood pressure and hemodynamics in adolescence: the Georgia cardiovascular twin study. *Hypertension* 2006;47:948–54.
46. Evans DM, Gillespie NA, Martin NG. Biometrical genetics. *Biol Psychol* 2002;61:33–51.
47. Nance WE, Corey LA. Genetic models for the analysis of data from the families of identical twins. *Genetics* 1976;83:811–26.
48. Posthuma D, Boomsma DI. A note on the statistical power in extended twin designs. *Behav Genet* 2000;30:147–58.
49. Eaves LJ, Last KA, Young PA, Martin NG. Model-fitting approaches to the analysis of human behaviour. *Heredity* 1978;41:249–320.
50. Snieder H, van Doornen LJ, Boomsma DI. The age dependency of gene expression for plasma lipids, lipoproteins, and apolipoproteins. *Am J Hum Genet* 1997;60:638–50.
51. Eaves LJ, Silberg JL, Maes HH. Revisiting the children of twins: can they be used to resolve the environmental effects of dyadic parental treatment on child behavior? *Twin Res Hum Genet* 2005;8:283–90.
52. Haley CS, Last K. The advantages of analysing human variation using twins and twin half-sibs and cousins. *Heredity* 1981;47:221–36.
53. Eaves L, Heath A, Martin N, Maes H, Neale M, Kendler K, Kirk K, Corey L. Comparing the biological and cultural inheritance of personality and social attitudes in the Virginia 30,000 study of twins and their relatives. *Twin Res* 1999;2:62–80.
54. Eaves L. The use of twins in the analysis of assortative mating. *Heredity* 1979;43:399–409.
55. Heath AC, Kendler KS, Eaves LJ, Markell D. The resolution of cultural and biological inheritance: informativeness of different relationships. *Behav Genet* 1985;15:439–65.
56. de Geus EJ, van Doornen LJ, Orlebeke JF. Regular exercise and aerobic fitness in relation to psychological make-up and physiological stress reactivity. *Psychosom Med* 1993;55:347–63.
57. Rice T, An P, Gagnon J, Leon AS, Skinner JS, Wilmore JH, Bouchard C, Rao DC. Heritability of HR and BP response to exercise training in the HERITAGE family study. *Med Sci Sports Exerc* 2002;34:972–9.
58. Salanti G, Sanderson S, Higgins JP. Obstacles and opportunities in meta-analysis of genetic association studies. *Genet Med* 2005;7:13–20.
59. Altshuler D, Brooks LD, Chakravarti A, Collins FS, Daly MJ, Donnelly P. A haplotype map of the human genome. *Nature* 2005;437:1299–320.
60. Hugot JP, Chamaillard M, Zouali H, Lesage S, Cezard JP, Belaiche J, Almer S, Tysk C, O'Morain CA, Gassull M, Binder V, Finkel Y, Cortot A, Modigliani R, Laurent-Puig P, Gower-Rousseau C, Macry J, Colombel JF, Sahbatou M, Thomas G. Association of NOD2 leucine-rich repeat variants with susceptibility to Crohn's disease. *Nature* 2001;411: 599–603.
61. Kruglyak L, Nickerson DA. Variation is the spice of life. *Nat Genet* 2001;27:234–6.
62. Romualdi C, Balding D, Nasidze IS, Risch G, Robichaux M, Sherry ST, Stoneking M, Batzer MA, Barbujani G. Patterns of human diversity, within and among continents, inferred from biallelic DNA polymorphisms. *Genome Res* 2002;12:602–12.
63. Cordell HJ, Clayton DG. Genetic association studies. *Lancet* 2005;366: 1121–31.
64. Manuck SB, Flory JD, Ferrell RE, Mann JJ, Muldoon MF. A regulatory polymorphism of the monoamine oxidase-A gene may be associated with variability in aggression, impulsivity, and central nervous system serotonergic responsivity. *Psychiatry Res* 2000;95:9–23.
65. Sullivan PF, Eaves LJ, Kendler KS, Neale MC. Genetic case-control association studies in neuropsychiatry. *Arch Gen Psychiatry* 2001;58: 1015–24.
66. Lander ES, Schork NJ. Genetic dissection of complex traits. *Science* 1994;265:2037–48.
67. Cardon LR, Palmer LJ. Population stratification and spurious allelic association. *Lancet* 2003;361:598–604.
68. Freedman ML, Reich D, Penney KL, McDonald GJ, Mignault AA, Patterson N, Gabriel SB, Topol EJ, Smoller JW, Pato CN, Pato MT, Petryshen TL, Kolonel LN, Lander ES, Sklar P, Henderson B, Hirschhorn JN, Altshuler D. Assessing the impact of population stratification on genetic association studies. *Nat Genet* 2004;36:388–93.
69. Hoggart CJ, Parra EJ, Shriver MD, Bonilla C, Kittles RA, Clayton DG, McKeigue PM. Control of confounding of genetic associations in stratified populations. *Am J Hum Genet* 2003;72:1492–504.
70. Pritchard JK, Donnelly P. Case-control studies of association in structured or admixed populations. *Theor Popul Biol* 2001;60:227–37.
71. Pritchard JK, Rosenberg NA. Use of unlinked genetic markers to detect population stratification in association studies. *Am J Hum Genet* 1999; 65:220–8.
72. Pritchard JK, Stephens M, Donnelly P. Inference of population structure using multilocus genotype data. *Genetics* 2000;155:945–59.
73. Pritchard JK, Stephens M, Rosenberg NA, Donnelly P. Association mapping in structured populations. *Am J Hum Genet* 2000;67:170–81.
74. Satten GA, Flanders WD, Yang Q. Accounting for unmeasured population substructure in case-control studies of genetic association using a novel latent-class model. *Am J Hum Genet* 2001;68:466–77.
75. Bacanu SA, Devlin B, Roeder K. The power of genomic control. *Am J Hum Genet* 2000;66:1933–44.
76. Bacanu SA, Devlin B, Roeder K. Association studies for quantitative traits in structured populations. *Genet Epidemiol* 2002;22:78–93.
77. Devlin B, Roeder K. Genomic control for association studies. *Biometrics* 1999;55:997–1004.
78. Devlin B, Roeder K, Bacanu SA. Unbiased methods for population-based association studies. *Genet Epidemiol* 2001;21:273–84.
79. Devlin B, Roeder K, Wasserman L. Genomic control, a new approach to genetic-based association studies. *Theor Popul Biol* 2001;60:155–66.
80. Spielman RS, McGinnis RE, Ewens WJ. Transmission test for linkage disequilibrium: the insulin gene region and insulin-dependent diabetes mellitus (IDDM). *Am J Hum Genet* 1993;52:506–16.
81. Abecasis GR, Cardon LR, Cookson WO. A general test of association for quantitative traits in nuclear families. *Am J Hum Genet* 2000;66: 279–92.
82. Horvath S, Xu X, Laird NM. The family based association test method: strategies for studying general genotype-phenotype associations. *Eur J Hum Genet* 2001;9:301–6.
83. van den Oord EJ, Snieder H. Including measured genotypes in statistical models to study the interplay of multiple factors affecting complex traits. *Behav Genet* 2002;32:1–22.
84. Horvath S, Xu X, Lake SL, Silverman EK, Weiss ST, Laird NM. Family-based tests for associating haplotypes with general phenotype data: application to asthma genetics. *Genet Epidemiol* 2004;26:61–9.

85. Strachan T, Read AP. *Human Molecular Genetics 2*. New York: John Wiley & Sons; 1999.
86. Thompson WD. Effect modification and the limits of biological inference from epidemiologic data. *J Clin Epidemiol* 1991;44:221–32.
87. Cordell HJ. Epistasis: what it means, what it doesn't mean, and statistical methods to detect it in humans. *Hum Mol Genet* 2002;11:2463–8.
88. Clayton D, McKeigue PM. Epidemiological methods for studying genes and environmental factors in complex diseases. *Lancet* 2001;358:1356–60.
89. Moffitt TE, Caspi A, Rutter M. Strategy for investigating interactions between measured genes and measured environments. *Arch Gen Psychiatry* 2005;62:473–81.
90. Turkheimer E, D'Onofrio BM, Maes HH, Eaves LJ. Analysis and interpretation of twin studies including measures of the shared environment. *Child Dev* 2005;76:1217–33.
91. McClelland GH, Judd CM. Statistical difficulties of detecting interactions and moderator effects. *Psychol Bull* 1993;114:376–90.
92. Zaykin DV, Westfall PH, Young SS, Karnoub MA, Wagner MJ, Ehm MG. Testing association of statistically inferred haplotypes with discrete and continuous traits in samples of unrelated individuals. *Hum Hered* 2002;53:79–91.
93. Fallin D, Schork NJ. Accuracy of haplotype frequency estimation for biallelic loci, via the expectation-maximization algorithm for unphased diploid genotype data. *Am J Hum Genet* 2000;67:947–59.
94. Stephens M, Donnelly P. A comparison of bayesian methods for haplotype reconstruction from population genotype data. *Am J Hum Genet* 2003;73:1162–9.
95. Carlson CS, Eberle MA, Kruglyak L, Nickerson DA. Mapping complex disease loci in whole-genome association studies. *Nature* 2004;429:446–52.
96. Cohen J. *Statistical power analysis for the behavioral sciences*. Hillsdale, NJ: Lawrence Erlbaum Associates, Publishers; 1988.
97. Kraemer HC, Thiemann S. A strategy to use soft data effectively in randomized controlled clinical trials. *J Consult Clin Psychol* 1989;57:148–54.
98. Gauderman J, Morrison J. *Quanto Version 1.0*; 2005.
99. Purcell S, Cherny SS, Sham PC. Genetic power calculator: design of linkage and association genetic mapping studies of complex traits. *Bioinformatics* 2003;19:149–50.
100. Blangero J. Localization and identification of human quantitative trait loci: king harvest has surely come. *Curr Opin Genet Dev* 2004;14:233–40.
101. Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical powerful approach to multiple testing. *J Royal Stat Soc B* 1995;57:289–300.
102. Storey JD, Tibshirani R. Statistical significance for genomewide studies. *Proc Natl Acad Sci USA* 2003;100:9440–5.
103. Colhoun HM, McKeigue PM, Davey Smith G. Problems of reporting genetic associations with complex outcomes. *Lancet* 2003;361:865–72.
104. Manly KF, Nettleton D, Hwang JT. Genomics, prior probability, and statistical tests of multiple hypotheses. *Genome Res* 2004;14:997–1001.
105. Neale B, Sham P. The Future of Association Studies: Gene-Based Analysis and Replication. *Am J Hum Genet* 2004;75:353–62.
106. Vink JM, Boomsma DI. Gene finding strategies. *Biol Psychol* 2002;61:53–71.
107. Risch N, Merikangas K. The future of genetic studies of complex human diseases. *Science* 1996;273:1516–7.

## ERRATUM

While drafting a related article, the authors of an article published last year discovered an error in the description of how sample size was calculated (Frasure-Smith N, Koszycki D, Swenson JR, Baker B, van Zyl LT, Laliberté M-A, Abramson BL, Lambert J, Gravel G, Lespérance F. Design and rationale for a randomized, controlled trial of interpersonal psychotherapy and citalopram for depression in coronary artery disease (CREATE). *Psychosom Med* 2006;68:87–93). On page 91, in the second paragraph of the Sample Size section, points 5 and 6 should have read as follows: “(5) adjustment for loss of final assessment of not more than 5% ( $n/0.95$ ); (6) adjustment for noncompletion of 12 weeks of treatment of not more than 20% in each group ( $n/((1 - 0.20)2)$ ).”