# Inference About Genetic Correlations

## Gregory Carey[1]

*In polygenic systems genetic correlations and the factors and specific genetic variances from genetic correlation matrices are often interpreted in terms of sets of genes common or specific to variables. While these inferences may indeed be true, a genetic correlation is not always sufficient evidence for the inferences. In some cases two variables with all genes in common can have low genetic correlations, and systems with only a few genes in common can have high genetic correlations. The assumptions about genic effects in polygenic systems and their effects on a genetic correlation are explicated and discussed. It is suggested that a distinction be made between* biological *pleiotropism and* statistical *pleiotropism to promote more accurate communication about the genetic associations among traits.*

## INTRODUCTION

Within the behavior genetic literature, one frequently reads statements relating genetic correlations or genetic factors to sets of loci. For example, it is not uncommon to hear or read that a large genetic correlation implies that two variables have a high proportion of loci in common, that low genetic correlations imply that different sets of loci underlie the two varibles, or that a general genetic factor with specific genetic variances sug-

[1] Institute for Behavioral Genetics and Department of Psychology, University of Colorado, Boulder, Colorado 80309-0447.

gests a common set of genes for the latent factor and specific sets of genes for each variable. In models of development, specific genetic variance at, say, age 6 years is sometimes interpreted as uncovering a set of genes whose effects are first manifest at age 6 years and not before. These inferences *may* be correct. However, genetic correlations and genetic factors are not sufficient evidence to prove the inferences about "sets of loci." Many behavior geneticists realize this point, but it is not clear that all do. This paper is offered in order to clarify the meaning of genetic correlations under different assumptions about additive, polygenic gene action.

## A MODEL FOR GENETIC CORRELATIONS

The model used to demonstrate the relationship between genetic correlations and sets of loci is the simple additive model most often used in introductory quantitative genetics (e.g., Falconer, 1981; Mather and Jinks, 1982). It expresses the genetic correlation in terms of genic effects.

The model assumes two alleles per locus, Hardy–Weinberg–Castle equilibrium, linkage equilibrium, and only additive genic effects. Let $p$ denote the frequency of one allele, and $q = 1 - p$, the frequency of the other allele. Let $h$ denote the genic value of the heterozygote or, in different words, the average phenotypic value for all heterozygotes in the population. Let $a$ denote the additive deviation from the heterozygote. The genic value for one homozygote becomes $h - a$ and the genic value for the other is $h + a$. If this locus is pleiotropic for two traits, the additive genic effect for the second trait can be parameterized as $ba$. Here, $b$ need not have the same value for all loci nor for the same locus across a number of traits. Table I presents the notation used here for the genic effects of a single pleiotropic locus on two variables. The two variables may represent two different traits or the same trait measured at different times.

**Table I.** Model for the Genetic Effect of a Single Pleiotropic Locus in Hardy–Weinberg–Castle Equilibrium on Two Variables

| | Genotype | | |
|---|---|---|---|
| | aa | Aa | AA |
| Genotypic frequency | $q^2$ | $2pq$ | $p^2$ |
| Average genic value, variable 1 | $h_1 - a$ | $h_1$ | $h_1 + a$ |
| Deviation from mean, variable 1 | $-2pa$ | $(q - p)a$ | $2qa$ |
| Average genic value, variable 2 | $h_2 - ba$ | $h_2$ | $h_2 + ba$ |
| Deviation from mean, variable 2 | $-2pba$ | $(q - p)ba$ | $2qba$ |

Some genes may contribute to variability in variable 1 but may not contribute to individual differences in variable 2, and some loci that affect variable 2 may not influence variability in trait 1. Instead of writing $a$ for these loci, let their additive effects be denoted by, respectively, $u_1$ and $u_2$. And let $f = 2pq$ to simplify notation.

The contribution of a pleiotropic locus to the total genotypic variance of variable 1 is

$$p^2(2qa)^2 + 2pq[(q - p)a]^2 + q^2(-2pa)^2 = fa^2.$$

The contribution of a unique locus is $fu_1^2$. The total genotypic variance is the sum of the contribution of all the pleiotropic loci plus the sum of all unique loci or, say, $\sum fa^2 + \sum fu_1^2$. (Subscripts that might denote loci are ignored in the summation.) The total genetic variance of the second trait becomes $\sum f(ba)^2 + \sum fu_2^2$.

The contribution of a pleiotropic locus to the genotypic covariance is

$$p^2(2qa)(2qba) + 2pq[(q - p)a][(q - p)ba] + q^2(-2pa)(-2pba) = fba^2,$$

so the total genetic covariance is $\sum fba^2$. The genetic correlation may then be written as

$$\sum fba^2/\sqrt{(\sum fa^2 + \sum fu_1^2)(\sum f(ba)^2 + \sum fu_2^2)}. \qquad (1)$$

The genetic correlation may now be viewed in the light of two different sets of assumptions or models. First, let each locus contribute equally to the total genetic variance variable and let allelic frequencies be identical across all loci. Loci are effectively equivalent to each other in this model. In Eq. (1), $a$ becomes constant with $u_1 = a$, $b$ becomes constant with $u_2 = ba$, $f$ becomes constant, and Eq. (1) reduces to a direct function of the number of common and unique loci:

$$n_p/\sqrt{(n_p + n_1)(n_p + n_2)}, \qquad (2)$$

where $n_p$ is the number of pleiotropic loci, and $n_1$ and $n_2$ are the number of loci unique to, respectively, variables 1 and 2. Under these conditions, the genetic correlation is directly interpretable as an index of the proportion of loci that two variables have in common.

Inference about sets of loci change when these assumptions are relaxed. Let allelic frequencies vary and let genic effects differ from one locus to another. The same genetic correlation can now arise in different ways. For example, consider the case of a large number of pleiotropic loci, each of relatively small effect, and a small number of unique loci of large effect. This genetic system could give a genetic correlation similar

to a different genetic system in which there were only a few pleiotropic loci with large genic effects and a large number of unique loci of small effect.

Another implication of unequal allelic effects and frequencies is that the genetic correlation need not be unity even though the same set of loci contributes to individual differences in both traits. To see how this occurs, let all $u$'s $= 0$ so that all loci are pleiotropic to traits 1 and 2, and let the $a$'s be scaled so that the genetic variance for variable 1 $= 1.0$. Equation (1) reduces to

$$\sum fba^2/\sqrt{\sum f(ba)^2} \tag{3}$$

and will equal unity only under special conditions (e.g., when $b$ is a scalar constant). Usually, the genetic correlation will be less than unity *even when the same set of loci underlies both variables.* When $b$ is totally random with respect to $a$, it is even mathematically possible to have a genetic correlation equal to zero even though there are no loci unique to either variable.

## SIMULATIONS OF GENETIC CORRELATIONS

Two sets of simulations were conducted to test the influence of unequal allelic effects and frequencies on genetic correlations when there are in fact no loci unique to either trait. In both sets, the effects of 25 loci on two variables were examined. Allelic frequencies were generated from a uniform distribution ranging from 0.01 to 0.99, the typical limits for a common polymorphism. Ten thousand sets of allelic frequencies were generated.

For each set of allelic frequencies, two types of genic values were generated. Each type of genic values represented a model in which genotypes such as AA that contribute to high scores on variable 1 would always contribute to high scores on variable 2. Only the relative magnitude of the genic contributions differed between the two models. If the two variables were quantitative and verbal ability, both simulated models predict that all alleles that increase quantitative ability will also increase verbal ability. Model 1 simulated a perfect rank-order correlation in genic effects on quantitative and verbal ability. That is, the allele that creates the greatest increase in quantitative ability is also that allele giving the largest increase in verbal ability. In model 2, the rank-order correlation is sampled from a distribution with mean $= 0.0$. That is, the allele that creates the greatest increase in quantitative ability will always increase verbal ability, but one cannot predict if it would be the *largest* increase in verbal ability. From analytical considerations of Eq. (1), model 1 should

always generate high genetic correlations. Model 2 should always give positive genetic correlations, but it is unclear exactly how high these correlations might be.[2]

To understand the simulation of genic values, assume that the alleles were rank ordered from 1 to 25 by the magnitude of their effect upon a variable. The genic value for the $j$th locus was taken from an exponential density function, $a_j = \lambda \exp(-\lambda j)$, where $\lambda$ is an arbitrary constant. Two values of $\lambda$ were generated, one for variable 1 and the other for variable 2. To maintain the perfect rank-order correlation of model 1, the same ordering of the loci was maintained for variables 1 and 2. The exponential distribution introduced only nonlinearity of effects; that is, the rank-order correlation between $a_j$ and $b_j$ was always unity but the Pearson correlation need not be unity. To simulate model 2, the ordering of the loci's effects was random.

The value of $\lambda$ was generated from a random uniform distribution ranging from 0.02 to 0.2225. With $\lambda = 0.02$, allelic effects are almost identical from one locus to another; the genic effects decreased almost linearly with $j$, locus 1 contributed 5% to the total genetic variance, and locus 25 contributed 3%. With $\lambda$ at its maximum of 0.2225, the polygenic system is similar to the mixed model (e.g., Elston and Stewart, 1971; Morton and MacClean, 1974) of a major locus with a substantial contribution from polygenic background. Here, genic effects decreased rapidly with $j$, the first locus contributed 20% to the total genetic variance, the total polygenic background contributed 80%, and the last locus contributed 0.1%.

To present data from the simulations, results were collapsed. Simulated polygenic systems were classified into three types by the value of $\lambda$. A $\lambda$ from 0.02 to 0.0875 generated polygenic systems that are arbitrarily termed "equal" for the simulations. Systems generated by $\lambda$ between 0.0875 and 0.155 are termed "weighted," and those with $\lambda$ between 0.155 and 0.2225 are called "mixed." For model 2, the genetic correlations were further subdivided into five groups on the basis of the Spearman rank-order correlation of allelic effects. The five groups were not rectangular because the actual Spearman correlations appeared normally distributed around a mean of 0; post hoc cutoffs were used in order to obtain at least

---

[2] The simulations allow allelic effects to vary only in magnitude, not in sign. It is obvious that genetic correlations will be less than unity under differences in sign. In speaking of this issue, Falconer (1981, p. 281) states that when some alleles increase one character but decrease the other, then "pleiotropy does not necessarily imply a detectable correlation." Also, it is not clear whether allelic effects that differ in sign are of important biological relevance (excluding the direction of scale for a variable). Thus, the interesting issue for research arises when allelic effects do not differ in sign but differ in magnitude across traits.

Table II. Mean Genetic Correlations as a Function of Type of Polygenic System and Rank-Order Correlation of Genic Effects[a]

| Polygenic models for two traits | Model 1 | Model 2: Range of Spearman rank-order correlation | | | | |
|---|---|---|---|---|---|---|
| | | −0.75 −0.45 | −0.44 −0.15 | −0.14 0.14 | 0.15 0.44 | 0.45 0.75 |
| Equal, equal | 99 | 80 | 84 | 88 | 90 | 93 |
| Equal, weighted | 95 | 62 | 68 | 73 | 78 | 82 |
| Equal, mixed | 86 | 50 | 55 | 61 | 67 | 73 |
| Weighted, weighted | 99 | 44 | 52 | 61 | 70 | 79 |
| Weighted, mixed | 98 | 32 | 41 | 51 | 62 | 74 |
| Mixed, mixed | 99 | 24 | 32 | 43 | 56 | 68 |
| Total | | | | | | |
|   Mean | 98 | 52 | 58 | 65 | 71 | 80 |
|   Minimum | 65 | 17 | 10 | 19 | 27 | 55 |
|   Maximum | 99 | 94 | 97 | 97 | 98 | 97 |

[a] Results from 10000 Monte Carlo simulations of a 25-locus polygenic system. Decimal points omitted for the genetic correlations.

100 simulations for the most extreme values of the rank order correlation. All correlations were first zeta transformed before deriving the mean. The mean zeta transforms were then transferred back to correlations for presentation.

Table II gives the results of the simulations. Several points are obvious from consideration of the analytical derivations given above. With a perfect rank-order correlation for allelic effects (model 1), the genetic correlations are uniformly high when two genetic systems of the same type are compared. This occurs because $b$ approaches a constant. Lower values are found with the most discrepant systems, a pairing of equal with mixed in this case. The distribution of correlations for all polygenic systems under model 1 was strongly and negatively skewed. Thus, although means are high, correlations of lower magnitude occasionally occur. The results from model 1 strongly suggest that nonlinearity in allelic effects generally has only a trivial influence upon polygenic systems. Also, because allelic frequencies can differ from one locus to another, the effect of differing allelic frequencies is unimportant when the rank order of allelic effects is unity.

When the rank-order correlations are not unity (model 2), genetic correlations are a function of the two polygenic systems, the rank-order correlation, and interactions of system with rank-order effects. The genetic correlation always increases with an increasing rank-order correlation, reflecting again the fact the $b$ becomes more correlated with $a$ with

increasing rank-order correlation. When allelic effects are approximately equal for both traits, the genetic correlation generally remains high; the actual correlations for two "equal" polygenic systems range between 0.59 and 0.98, with an overall mean of 0.88. However, at the opposite extreme, two mixed polygenic systems can generate relatively low genetic correlations; the correlations here average 0.44 and range from 0.10 to 0.89. Even at the highest value of Spearman's rho in model 2, the mixed–mixed system generated a genetic correlation of 0.56. This correlation explains 31% of the total genetic variance even though the same loci contribute to all the genetic variance in both traits.

The similarity in the type of polygenic system no longer is associated with high genetic correlations in model 2. The two most discrepant polygenic systems (the equal with mixed) give higher correlations than the mixed with mixed. Thus, as the rank-order correlation departs from unity, the veridicality of assumptions about equal effects of alleles in polygenic systems becomes more critical for avoiding errors of inference about genetic correlations.

## DISCUSSION

The analytical development and the simulations given above suggest that a low genetic correlation can arise even when the same genes are involved in two traits and when allelic effects do not simply vary in sign. But do they represent biologically relevant circumstances?

It is obvious that when (1) genes lay down a biological structure and (2) anatomical and/or physiological structure is primarily responsible for individual differences in two or more measurable variables, then the genetic correlation will be high. This case was simulated by model 1 when the two polygenic systems are the same. And here it makes no difference whether allelic effects are the same or differ across loci.

But what happens if the same set of loci affects two related biological structures that influence two different traits, say the concentration of dopamine receptors (DAR) in two different areas of the brain such as the basal ganglia and olfactory lobes? Here, a prediction about the genetic correlation is less clear. On the one hand, loci that have a large effect on DAR in one area may have an equally large effect in the second area. On the other hand, consider what might happen if natural selection begins to operate on both traits. Suppose that the optimum DAR concentrations are not positively and linearly related in both regions (e.g., one region is under directional selection, while the other is under stabilizing selection). Will evolution take the genes that are already there and alter regulatory

mechanisms that change allelic values to obtain the most adaptive DAR concentrations in both brain regions?

Another circumstance where the simulations may be appropriate is during an active phase of development. Some loci may be very important during initial stages of development but become increasingly less important as the organism approaches maturity. Other loci that exert a major effect toward the end of development may have only a minor role during the initial phases. In this case, the same loci could operate during the whole developmental period, but the rank-order correlation for allelic effects would be negative. The simulations suggest marked caution in interpreting genetic correlations in this case. With a negative rank-order correlation, the genetic correlation depends greatly upon the type of underlying polygenic system—equal versus fixed—and could actually approach zero. Low genetic correlation across time periods could lead to the erroneous conclusion that specific loci are coming into play at different points in development, when exactly the opposite is occurring.

Finally, the model and simulations have ignored the influences of nonadditive allelic effects and linkage disequilibrium. Dominance and epistatic effects that differ across traits may also reduce a genetic correlation. Linkage disequilibrium can induce a genetic correlation even when there are *no* loci in common. If the simple additive model can produce such results, nonadditivity and linkage disequilibrium mean that the conclusions must hold a fortiori.

Given that there is some biological relevance to the development herein, how can we reconcile these conclusions with statements such as those made by Falconer (1981, p 281): "the degree of correlation arising from pleiotropy expresses the extent to which two characters are influenced by the same genes"? It appears necessary to distinguish *biological pleiotropism* (in which the same genes physically underlie different traits) from *statistical pleiotropism* (in which allelic effects on one trait predict allelic effects on other characters). Falconer's statement applies to statistical pleiotropism, not to biological pleiotropism. The distinction is crucial, especially for multivariate analysis.

In the statistical sense, the genetic "factors" and specific genetic variances derived from multivariate analysis are hypothetical predictive constructs used to reproduce parsimoniously a genetic covariance matrix. Interpreting genetic factors in terms of biological pleiotropism may be problematic for several reasons. First, there is a logical problem. With more than two traits, the number of potential sets of pleiotropic genes exceeds the number of traits. For example, with four variables, there are 11 possible sets—1 for all four variables, 4 for the variables taken three at a time, and 6 for all pairwise combinations. If a genetic factor defines

a "set of genes," then there will be more possible factors than variables! (An interesting question here is whether it is possible to characterize possible gene sets that satisfy, say, a single factor model.)

Second, a broad general factor can emerge even when there is minimal or no overlap in the genetical systems underlying each trait. For example, consider four traits with the following genetical system: traits 1 and 2 have a certain proportion of loci in common, traits 1 and 3 have another set in common, traits 1 and 4 have yet another, etc. All pairwise genic sets could be taken even if there were no single locus pleiotropic to all four traits. This system could produce a matrix of genetic correlations in which the correlations are approximately equal. The first eigenvalue of this matrix would be very large, and the first eigenvector would have approximately equal loadings from each trait. This broad general factor exists, yet there is not a single locus common to the three traits.[3]

Third and probably most important, genetic factors that are weakly correlated could reflect similar biological systems with different allelic effects. The simulations (and some reflection on the examples given above) suggest that inference about sets of genes may be more valid when genetic correlations are high, the assumption of additivity is robust, and there is good reason to reject a major gene model. However, low genetic correlations do not preclude a common set of loci.

To understand the relationship between biological and statistical pleiotropism, we must trace pathways from gene to character, a topic discussed by Wright (1967, Chap. 5), Mather and Jinks (1982, Chap. 2), and others. Until the time when such data are available, *some* genetic correlations should be interpreted in a statistical sense and terms such as "shared genetic variance" and "specific genetic variance" should be used to refer to statistical pleiotropism. Terms such as "sets of loci in common" and "unique genes" best refer to *biological* pleiotropism.

## ACKNOWLEDGMENTS

## REFERENCES

Elston, R. C., and Stewart, J. (1971). A general model for the genetic analysis of pedigree data. *Hum. Hered.* 21:523–542.

[3] D. I. Boomsma and P. Molenaar (personal communication) point out that the problem extends beyond genes and a general factor. Thompson (1916, 1920) demonstrated that a correlation matrix could give evidence for a general factor when in fact there was not one.

Falconer, D. S. (1981). *Introduction to Quantitative Genetics*, 2nd ed., Longman Press, London.

Mather, K., and Jinks, J. L. (1982). *Biometrical Genetics: The Study of Continuous Variation*, 3rd ed., Chapman and Hall, London.

Morton, N. E., and MacLean, C. J. (1974). Analysis of family resemblances. III. Complex segregation analysis of quantitative traits. *Am. J. Hum. Genet.* **26**:489–504.

Thompson, G. H. (1916). A hierarchy without a general factor. *Br. J. Psychol.* **8**:271–281.

Thompson, G. H. (1920). The general factor fallacy in psychology. *Br. J. Psychol.* **10**:319–326.

Wright, S. (1967). *Evolution and the Genetics of Populations, Vol. 1*, University of Chicago Press, Chicago.

Edited by N. G. Martin