

Detecting Marker-Disease Association by Testing for Hardy-Weinberg Disequilibrium at a Marker Locus

Dahlia M. Nielsen,^{1,2} M. G. Ehm,¹ and B. S. Weir²

¹Bioinformatics Department, Glaxo Wellcome, Inc., Research Triangle Park, NC; and ²Program in Statistical Genetics, Department of Statistics, North Carolina State University, Raleigh

Summary

We review and extend a recent suggestion that fine-scale localization of a disease-susceptibility locus for a complex disease be done on the basis of deviations from Hardy-Weinberg equilibrium among affected individuals. This deviation is driven by linkage disequilibrium between disease and marker loci in the whole population and requires a heterogeneous genetic basis for the disease. A finding of marker-locus Hardy-Weinberg disequilibrium therefore implies disease heterogeneity and marker-disease linkage disequilibrium. Although a lack of departure of Hardy-Weinberg disequilibrium at marker loci implies that disease susceptibility-weighted linkage disequilibria are zero, given disease heterogeneity, it does not follow that the usual measures of linkage disequilibrium are zero. For disease-susceptibility loci with more than two alleles, therefore, care is needed in the drawing of inferences from marker Hardy-Weinberg disequilibria.

Introduction

We will refer to fine mapping as attempting to narrow what may be a 10-cM region indicated by linkage analysis to an $\sim \leq 1$ -cM region containing the disease-susceptibility locus. Fine-mapping methods for qualitative and quantitative phenotypic traits have been under constant development in recent years. Simple Mendelian traits with high penetrance are often fine mapped by recombinant mapping: typing markers every 1–2 cM, determining haplotypes by use of extended family information, and identifying recombination events on either side of the supposed disease-susceptibility locus

(Boehnke 1994). Glaser et al. (1995) illustrate this approach in a search for the gene responsible for familiar hyperinsulinism. In the absence of high penetrance or sufficient numbers of patients, linkage-disequilibrium methods in isolated populations have been used. Hästbacka et al. (1992) used linkage disequilibrium to map diastrophic dysplasia (DTD) in Finland and indicated that the DTD gene should lie <0.06 cM from the CSF1R gene, which was later confirmed (Hästbacka et al. 1994). Although the identification of disease-susceptibility loci for complex traits has been slow, several fine-mapping methods have been employed. Identification of IDDM2 (the insulin gene) was accomplished by use of linkage disequilibrium (Bennett et al. 1995). Association studies played a key role in implicating the apolipoprotein E gene in late-onset Alzheimer disease and heart disease (Corder et al. 1993).

Although methods for mapping simple Mendelian diseases use extended families collected for the genomic scan to refine disease-gene locations, fine-mapping techniques for complex diseases use samples with varying characteristics. Several methods have been proposed. Association methods use unrelated cases and controls. Traditional transmission/disequilibrium tests require affected children and their parents (Spielman et al. 1993; Kaplan et al. 1997). Several tests using data sets that are simpler to collect and are more similar to those used in a genomic scan have been proposed. Trégouët et al. (1997) have proposed making use of estimating equations to estimate association parameters in samples of nuclear families of varying sizes and in mixtures of related and unrelated individuals. Martin et al. (1997) have proposed two test statistics for association that use data from all affected children (and their parents) in a nuclear family. Spielman and Ewens (1998) proposed a test statistic that tests for linkage disequilibrium by use of affected and unaffected siblings. Feder et al. (1996) have suggested that fine localization of a disease-susceptibility locus could be accomplished by use of deviation from Hardy-Weinberg (HW) equilibrium among affected individuals.

Feder et al. (1996) studied hereditary hemochromatosis (HH), a common autosomal recessive disorder of

Received June 26, 1998; accepted for publication September 9, 1998; electronically published October 26, 1998.

Address for correspondence and reprints: Dr. M. G. Ehm, Glaxo Wellcome, Inc., Bioinformatics Department, 5 Moore Drive, Research Triangle Park, NC 27709. E-mail: mge37216@glaxowellcome.com

© 1998. The American Society of Human Genetics. All rights reserved.
0002-9297/98/6305/0031\$02.00

iron metabolism. As described in their paper, previous localization of the HH gene placed it near the major histocompatibility complex on chromosome 6p and <1-2 cM from the HLA-A gene, although many reports have been contradictory. Linkage-disequilibrium studies confirmed the existence of a founder effect. To proceed with the localization of a gene involved in this disease, Feder et al. (1996) developed 45 short tandem-repeat polymorphism and single nucleotide polymorphism (SNP) markers lying within an 8-cM region suspected to contain the gene. All 45 markers were typed in 101 HH patients and 64 controls.

To estimate the position of the gene relative to these closely spaced markers, Feder et al. (1996) used the measure p_{excess} (Bengtsson and Thomson 1981; Lehesjoki et al. 1993). This is a measure of linkage disequilibrium that, in the presence of linkage, is expected to be maximized at the marker nearest the gene. Feder et al. (1996) plotted p_{excess} for each marker along the marker map. This plot had a peak representing the maximum p_{excess} in the region; however, the peak was not very sharp, causing concern about the accuracy of these results.

In examining the data used in this study, Feder et al. (1996) noted that, among the affected individuals, there appeared to be an excess of homozygosity at the marker loci. They considered several explanations for this, the most likely of which proved to provide the basis for a new measure for linkage disequilibrium. They noted that, for heterogeneous recessive traits such as those which they were studying, not only will an excess of homozygosity exist among affected individuals, but also this excess homozygosity should decrease with decreased linkage disequilibrium between the marker loci and the disease-susceptibility locus. A disequilibrium measure based on this observation has the advantage that only affected individuals need to be collected and genotyped, as opposed to the case-control type studies necessary for most measures of association, such as p_{excess} . In their paper, Feder et al. (1996) plotted a measure of HW disequilibrium within their HH-affected individuals, for each marker along the marker map. This plot had a maximum at approximately the same point in the map as did p_{excess} , but the peak was much sharper. From this, Feder et al. (1996) concluded that their initial results from the p_{excess} measure had been confirmed and that the region in which the gene lies had been more accurately defined.

We were impressed with these results and were interested in exploring the properties of this measure. We have extended the model and have examined some general results for this and other measures. We also compare a test for HW disequilibrium, to a direct test for linkage disequilibrium.

Methods

Recessive Disease Model

Feder et al. (1996) examined a heterogeneous recessive model in which a subset of the disease cases are due to a mutation in the region of interest and in which other disease cases are due to either unrelated genetic loci or nongenetic factors. If “A” is used to denote the disease allele and “ \bar{A} ” is used to denote all other alleles at the disease-susceptibility locus, this model can be summarized as $\Pr(\text{Affected}|AA) = 1$, $\Pr(\text{Affected}|A\bar{A}) = \psi$, and $\Pr(\text{Affected}|\bar{A}\bar{A}) = \psi$, where ψ is the probability that an individual will exhibit the disease because of causes other than this locus. With the assumption of random mating in the population, genotype and allele probabilities at the disease-susceptibility locus among affected individuals can be calculated and include

$$\Pr(AA|\text{Affected}) = P_{AA|\text{Affected}} = p_A^2/\phi,$$

$$\Pr(A|\text{Affected}) = p_{A|\text{Affected}} = p_A(p_A + \psi p_{\bar{A}})/\phi,$$

where ϕ is the prevalence of the disease in the population. We have used “ $p_{A|\text{Affected}}$ ” and “ $P_{AA|\text{Affected}}$ ” to differentiate between frequencies among affected individuals and the whole-population frequencies, denoted as “ p_A ” and “ P_{AA} ,” respectively. For this model, $\phi = p_A^2 + \psi(1 - p_A^2)$.

Departure from HW equilibrium at the disease-susceptibility locus can be measured by the disequilibrium coefficient $\mathcal{D}_{AA} = P_{AA} - p_A^2$ (Weir 1996). Among affected individuals, this coefficient becomes

$$\begin{aligned} \mathcal{D}_{AA|\text{Affected}} &= P_{AA|\text{Affected}} - p_{A|\text{Affected}}^2 \\ &= \psi(1 - \psi)p_A^2(1 - p_{\bar{A}})^2/\phi^2. \end{aligned}$$

Feder et al. (1996) quantified departure from HW equilibrium at the disease-susceptibility locus by use of the measure F_A , defined as $(H_o - H_e)/(1 - H_e)$, where H_o and H_e are the observed and the expected homozygosities, respectively. Although they did not give an explicit expression for this quantity, it appears to us that they used the formulation

$$\begin{aligned} F_A &= \frac{P_{AA|\text{Affected}} + P_{\bar{A}\bar{A}|\text{Affected}} - p_{A|\text{Affected}}^2 - p_{\bar{A}|\text{Affected}}^2}{1 - p_A^2 - p_{\bar{A}}^2} \\ &= 2\mathcal{D}_{AA|\text{Affected}}/(2p_A p_{\bar{A}}) \\ &= \psi(1 - \psi)p_A p_{\bar{A}}/\phi^2. \end{aligned}$$

Association between the disease-susceptibility allele A and a marker allele M can be expressed by use of the

linkage-disequilibrium measure $D_{AM} = P_{AM} - p_A q_M$, where q_M is the frequency of marker allele M . This quantity compares the frequency (P_{AM}) of haplotypes carrying both alleles A and M with the product of the separate frequencies of the two alleles. D_{AM} is positive when marker allele M is more likely to be associated with disease-susceptibility allele A than would be expected by chance.

Feder et al. (1996) also discussed HW disequilibrium at a biallelic marker locus. With the assumption of random mating in the whole population, probabilities for the marker alleles and marker genotypes conditioned on having the disease include

$$P_{MM|Affected} = [(1 - \psi)(p_A q_M + D_{AM})^2 + \psi q_M^2] / \phi ,$$

$$q_{M|Affected} = [\psi q_M + (1 - \psi)p_A(p_A q_M + D_{AM})] / \phi .$$

The HW disequilibrium coefficient at the marker locus among affected individuals is

$$D_{MM|Affected} = \psi(1 - \psi)D_{AM}^2 / \phi^2 . \quad (1)$$

This is non-0 only if ψ is neither 1 nor 0, implying that the disease must be heterogeneous, and that, if there is linkage disequilibrium, $D_{AM} \neq 0$. The HW-departure measure of Feder et al. (1996) for the marker locus is

$$F_M = \frac{P_{MM|Affected} + P_{\bar{M}\bar{M}|Affected} - q_M^2 - q_{\bar{M}}^2}{1 - q_M^2 - q_{\bar{M}}^2}$$

$$= D_{MM|Affected} / (q_M q_{\bar{M}})$$

$$= \psi(1 - \psi)D_{AM}^2 / (\phi^2 q_M q_{\bar{M}}) .$$

As stated by Feder et al. (1996),

$$F_M = \Delta_{AM}^2 F_A , \quad (2)$$

where $\Delta_{AM}^2 = D_{AM}^2 / p_A p_A q_M q_{\bar{M}}$.

Equations (1) and (2) capture the essential point that HW disequilibrium at a marker locus among affected individuals depends on the whole-population linkage disequilibrium between the marker locus and the disease locus. Although it is the latter quantity that is of interest, it is easier to test for the former. A test for HW disequilibrium at the marker locus can serve as a test for linkage disequilibrium. It should be noted, however, that the measure F_M proposed by Feder et al. (1996) depends on the values q_M and $q_{\bar{M}}$, which are whole-population parameters and cannot be estimated by use of affected individuals alone.

A common direct measure of linkage disequilibrium is the quantity p_{excess} (Bengtsson and Thomson 1981; Lehesjoki et al. 1993). This measure compares the frequency of a marker allele M among affected individuals

($q_{M|Affected}$) versus the frequency among unaffected individuals ($q_{M|Unaffected}$). It is defined as

$$p_{\text{excess}} = \frac{q_{M|Affected} - q_{M|Unaffected}}{1 - q_{M|Unaffected}} .$$

For the model of Feder et al. (1996),

$$q_{M|Affected} = [\psi q_M + (1 - \psi)p_A(p_A q_M + D_{AM})] / \phi ,$$

$$q_{M|Unaffected} = (1 - \psi)[q_M - p_A(p_A q_M + D_{AM})] / (1 - \phi) ,$$

so that

$$p_{\text{excess}} = \frac{(1 - \psi)p_A D_{AM}}{\phi(1 - \phi)[q_M + (1 - \psi)p_A D_{AM} / (1 - \phi)]} . \quad (3)$$

Therefore p_{excess} is proportional to D_{AM} , and it reaches its maximum at the marker with the greatest disequilibrium with the disease. Note that ψ must be < 1 . HW disequilibrium is proportional to the square of disequilibrium, so that F_M is expected to be a more sensitive indicator of linkage in the presence of linkage disequilibrium (eq. [2]). This appears to have been the case in the analyses reported by Feder et al. (1996).

General Disease Model

We wished to know whether equation (1) might be generalized to other disease models, and so we considered a more general model with disease susceptibility affected by a locus with an arbitrary number of alleles, denoted by " A_r ." Under this model, the conditional probability that an individual has the disease, given that the individual has genotype $A_r A_s$ at the disease-susceptibility locus, is ϕ_{rs} . We will refer to these values as "penetrances," although we recognize that, for some of the $A_r A_s$ genotypes, the ϕ_{rs} values should properly be called "phenocopy rates." These values could equivalently be called "prevalences": they represent the prevalence of the disease within a genotypic class. This relates the notation " ϕ_{rs} " to the use of " ϕ ," the whole-population prevalence (the unconditional probability that an individual has the disease). This value is $\phi = \sum_i \sum_s \phi_{rs} p_r p_s$, where p_r is the population frequency of allele A_r at the disease-susceptibility locus and where HW equilibrium is assumed.

In addition to the genotypic penetrances ϕ_{rs} , we find it convenient to define an allelic penetrance, ϕ_r : $\phi_r = \sum_s p_s \phi_{rs}$, which is the conditional probability that an individual will have the disease, given that the individual has allele A_r (the other allele being a random allele from the population). Note that $\phi = \sum_r p_r \phi_r$.

We consider a marker locus with alleles M_i occurring at frequencies q_i . For such a marker, we are likely to concentrate on those alleles that show positive associ-

ations with the disease, meaning that they have higher frequencies among affected than among unaffected individuals. If P_{ri} is the population frequency of haplotypes carrying disease-susceptibility allele A_r and marker allele M_i , then the population linkage disequilibrium D_{ri} between these alleles is defined as $D_{ri} = P_{ri} - p_r q_i$. These coefficients sum to 0 over all the alleles at either locus, so that $\sum_r D_{ri} = \sum_i D_{ri} = 0$. We also wish to describe the linkage disequilibrium between marker allele M_i and the disease locus as a whole, and we do so by weighting the D_{ri} terms by the allelic penetrances. This measure is written as $\delta_i = \sum_r \phi_r D_{ri} = \sum_r \sum_s p_s \phi_{rs} D_{ri}$ and sums to 0 over i . The quantity δ_i is 0 if all disease susceptibility-locus alleles have the same penetrances.

For a disease-susceptibility locus with two alleles, A_1 and A_2 , δ_i is a multiple of D_{1i} and so is proportional to the usual linkage-disequilibrium coefficient and will maximize at the same point as does linkage disequilibrium. In this two-allele case, if the penetrances are not the same, a 0 value of δ_i implies that there is no linkage disequilibrium between disease susceptibility and marker loci. For a disease-susceptibility locus with more than two alleles, however, it is possible for δ_i to be near or equal to 0 even when there is linkage disequilibrium, since the values of D_{ri} do not all have the same sign and may have a (penetrance-weighted) sum close to 0.

The penetrance-weighted linkage-disequilibrium coefficient allows simple expressions for marker-allele frequencies among affecteds: $q_{i|Affected} = q_i + (\delta/\phi)$, as is shown in Appendix A. This equation shows that marker-allele frequencies among affected individuals deviate from the overall population frequencies by an amount that depends on the strength of association between the marker allele and the disease-susceptibility alleles, weighted by the penetrances of those alleles. A similar expression holds for the marker-allele frequency among unaffected individuals, $q_{i|Unaffected} = q_i - [\delta_i/(1 - \phi)]$, so that, in the whole population, $q_i = \phi q_{i|Affected} + (1 - \phi)q_{i|Unaffected}$.

As a generalization of equation (3), the quantity p_{excess} for marker allele M_i becomes

$$p_{\text{excess}_i} = \frac{\delta_i}{\phi(1 - \phi)[(1 - q_i) + \delta_i/(1 - \phi)]}$$

If M_i is a marker allele showing a positive association with the disease, then $p_{\text{excess}_i} \geq 0$, so that $\delta_i \geq 0$, and these two quantities are maximized together. However, it is not necessary that each individual linkage-disequilibrium coefficient D_{ri} be positive.

For the general disease model, discussion of marker-locus HW disequilibrium requires an additional summary measure of linkage disequilibrium. This quantity, δ_{ij} , is defined for pairs of marker alleles, M_i and M_j , instead of for single marker alleles: $\delta_{ij} = \sum_r \sum_s \phi_{rs} D_{ri} D_{sj}$.

We term it “genotypic disequilibrium,” as opposed to the “allelic disequilibrium” δ_i . Note that $\sum_i \delta_{ij} = \sum_j \delta_{ij} = 0$. Among affected individuals, the marker-locus-homozygote HW disequilibrium coefficients can now be written as

$$D_{ii|Affected} = P_{ii|Affected} - q_{i|Affected}^2 = \frac{\phi \delta_{ii} - \delta_i^2}{\phi^2},$$

and the heterozygote disequilibria (Weir 1996) are

$$D_{ij|Affected} = P_{ij|Affected} - 2q_{i|Affected}q_{j|Affected} = \frac{2(\phi \delta_{ij} - \delta_i \delta_j)}{\phi^2}.$$

For this more general model, it is not clear that the HW disequilibria, $D_{ij|Affected}$ and $D_{ii|Affected}$, are maximized when the linkage disequilibria, δ_i , are maximized. It is clear, however, that some patterns of non-0 linkage disequilibrium will result in 0 departure from HW equilibrium at a marker locus. Conversely, a departure from HW equilibrium at a marker locus provides evidence both for linkage disequilibrium between marker and disease-susceptibility loci and for heterogeneity of disease susceptibility.

Test Statistics

We have discussed two measures that can be used to characterize marker/disease associations. One is p_{excess} , which is directly proportional to linkage disequilibrium measured in unrelated affected and unaffected individuals, and the other is the HW-disequilibrium coefficient measured among affected individuals. To compare these two approaches, we consider the statistical power of corresponding test statistics.

A widely used statistical test for association based on unrelated affected and unaffected individuals—that is, a case-control design—is the $(m - 1)$ -df χ^2 test based on the statistic χ_{CC}^2 when the marker locus has m alleles. When the marker alleles have sample frequencies $\tilde{p}_{i|Affected}$ and $\tilde{p}_{i|Unaffected}$ among n affecteds and n unaffecteds,

$$\chi_{CC}^2 = 4n \sum_i \frac{(\tilde{p}_{i|Affected} - \tilde{p}_{i|Unaffected})^2}{\tilde{p}_{i|Affected} + \tilde{p}_{i|Unaffected}}. \tag{4}$$

When alternatives to the null hypothesis of no disequilibrium are of the Pitman type (i.e., departures tend toward 0 with sample size), the noncentrality parameter of this statistic is (Meng and Chapman 1966)

$$\lambda_{CC} = 4n \sum_i \frac{(q_{i|Affected} - q_{i|Unaffected})^2}{q_{i|Affected} + q_{i|Unaffected}}$$

$$= 4n \sum_i \frac{\delta_i^2}{\phi^2(1 - \phi)^2 \{2q_i + (1 - 2\phi)\delta_i / (\phi(1 - \phi))\}}$$

Elsewhere (Kaplan et al. 1997), we have written the sum in this expression as I^* .

To test for HW disequilibrium at the marker locus among the same total number of individuals, $2n$ affecteds, the test statistic χ^2_{HW} is (Weir 1996)

$$\chi^2_{HW} = n \sum_i \frac{(\tilde{p}_{i|Affected} - \tilde{q}_{i|Affected}^2)^2}{\tilde{q}_{i|Affected}^2}$$

$$+ 2n \sum_{i < j} \frac{(\tilde{p}_{ij|Affected} - 2\tilde{q}_{i|Affected}\tilde{q}_{j|Affected})^2}{2\tilde{q}_{i|Affected}\tilde{q}_{j|Affected}}$$
(5)

This has $m(m - 1)/2$ df and a noncentrality parameter of

$$\lambda_{HW} = 2n \sum_i \sum_j \frac{(\phi\delta_{ij} - \delta_i\delta_j)^2}{\phi^2(\phi q_i + \delta_i)(\phi q_j + \delta_j)}$$

When there are just two marker alleles, $m = 2$, the power of the two χ^2 tests can be compared directly by comparing λ_{CC} with λ_{HW} . For this case, HW disequilibrium decays at a rate proportional to the square of linkage disequilibrium (Appendix B). This indicates that the measure of HW disequilibrium should be a more sensitive indicator of position, decaying more quickly than the measure of linkage disequilibrium as the distance between the marker locus and the disease-susceptibility locus increases.

Simulations

To illustrate our theoretical results, we performed simulations of evolving populations segregating for a biallelic disease-susceptibility locus and several biallelic markers. We performed these simulations under four different disease models, representing special cases of the general model. Analytical results for these special cases can be found in Appendix B. For the four models, we performed χ^2 tests for linkage disequilibrium and for HW disequilibrium and compared the estimated power of the results.

For all four models, we considered a marker allele M , at frequency $q_M = .20$, that had a positive association with the disease allele. Our first simulated model was the heterogeneous recessive model of Feder et al. (1996),

with $p_A = .10$ and $\psi = .05$. For these parameters, the maximum linkage disequilibrium expected is .08. The second model was also of the Feder et al. (1996) type, but with different parameter values. For this model, we chose the parameters $p_A = .05$ and $\psi = .05$. Since p_A is smaller in the second model, less linkage disequilibrium is possible, reaching a maximum expected value of .04, one-half of what was expected in the first model.

The third model was an additive model for penetrance. We set the effect of the disease-causing allele (A) at .50 and set the effect of the nondisease allele (\bar{A}) at 0. This yields $\phi_{AA} = 1.0$, $\phi_{A\bar{A}} = .5$, and $\phi_{\bar{A}\bar{A}} = .0$. The frequency of the disease allele in the population, p_A , was .10. For the additive model, HW disequilibrium is expected to be negative (Appendix B) and will increase in absolute value with increasing linkage disequilibrium.

A multiplicative model for penetrance was assumed for the fourth set of simulations. We set the effect of the disease-causing allele (A) at .9 and set the effect of the nondisease allele (\bar{A}) at .05. This leads to $\phi_{AA} = .8100$, $\phi_{A\bar{A}} = .0450$, and $\phi_{\bar{A}\bar{A}} = .0025$. The frequency of the disease allele in the population, p_A , was .10. We did not expect to see any HW disequilibrium among the affected individuals (Appendix B). A summary of the parameter values used in these four models can be found in table 1.

For our simulated populations, we considered marker loci positioned at distances of 0–2 cM from the disease-susceptibility locus, considering one marker every 0.25 cM. The populations started at generation G_0 with complete association between the disease allele and one allele at each marker locus, then evolved for 50 generations of random mating. For each model, we retained the first 100 populations, which, after 50 generations, had not experienced substantial genetic drift at the disease locus. For a population to be accepted, the frequency of the disease allele at the end of the evolution could not deviate from the original frequency by $>.05$. We made no adjustments for genetic drift at the marker locus.

Results

Power

To determine the power to detect HW and linkage disequilibria, we performed the χ^2 tests χ^2_{CC} and χ^2_{HW} (eqs.

Table 1

Parameters of the Simulated Disease Models

Model (Type)	ϕ_{AA}	$\phi_{A\bar{A}}$	$\phi_{\bar{A}\bar{A}}$	p_A	p_B	D_{max}^a
1 (heterogeneous recessive)	1.00	.05	.05	.10	.20	.08
2 (heterogeneous recessive)	1.00	.10	.10	.05	.20	.04
3 (additive)	1.00	.50	0	.10	.20	.08
4 (multiplicative)	.81	.045	.0025	.10	.20	.08

^a Values are maximum expected disequilibrium.

[4] and [5], respectively) on samples taken from each population. For the case-control test, we sampled 50 affected and 50 unaffected individuals from each population. For the test for HW disequilibrium, we sampled 100 affected individuals. We repeated both tests 5,000 times for each population, recording the percentage of times that we rejected the hypothesis of no disequilibrium. This rejection percentage gave us an estimate of the power of the respective tests. The comparisons of these results can be seen in figure 1. The symbols in this figure are box plots of the results; the bottom and top edges of the box are located at the sample 25th and 75th percentiles, the point joined by the connecting line is the median, and the whiskers extend the range of the results. This figure shows that, for the Feder et al.-type models (figs. 1A and 1B), the power to detect HW disequilibrium is greater in general than the power to detect linkage disequilibrium. This is particularly noteworthy in the case of the second Feder et al.-type model (fig. 1B), in which the power to detect linkage disequilibrium by use of χ^2_{CC} is not very different from the $\alpha = .05$ nominal level. The additive model (fig. 1C) shows high power for both tests. For the multiplicative model (fig. 1D), we expected to find no HW disequilibrium among affected individuals. These experiments showed that, although the power to detect HW disequilibrium was very low for this model, it was very frequently above the $\alpha = .05$ nominal level. This appears to be caused by increased variance of HW-disequilibrium values, which is created by sampling of affected individuals.

Figure 1 reveals the variability of the power of the two χ^2 tests. For several of these experiments, the power of χ^2_{HW} varied from the nominal .05 level to values close to 1.0. χ^2_{CC} was less variable in its power.

In these experiments, we generated linkage disequilibrium in the presence of physical linkage. Thus, both tests showed reduced power at greater distances between the loci. As expected, in the three models in which we expected to find HW disequilibrium, the power to detect HW disequilibrium decayed more quickly than did the power to detect linkage disequilibrium.

Size

To determine the size of our tests, we simulated a second set of populations under the same four disease models but segregating for a biallelic marker located at 50% recombination from the disease locus. We performed the same sampling and testing experiments as described above. The results from these experiments are displayed in figure 2. These results showed that the case-control test, χ^2_{CC} , was conservative; in most cases examined, it rejected the true null hypothesis at a rate less than the $\alpha = .05$ nominal level. The rejection rate of the test for HW disequilibrium, χ^2_{HW} , appeared to be cen-

tered around the nominal $\alpha = .05$ level for the first three models but was higher for model 4, the multiplicative model. For model 4, the variance of the rejection rate appeared to be quite large. It was, however, very similar to that seen in the populations generated with linked markers, as shown in figure 1.

Discussion

We have examined departures from HW equilibrium that are created by sampling of individuals on the basis of the presence of a disease phenotype. These departures from equilibrium are created because the selection criterion is based on disease-susceptibility genotypes, rather than on independently selected alleles. Alleles within genotypes that confer greater susceptibilities are represented in the sample at disproportionately high rates. Disequilibrium is expected to be greatest at the disease-susceptibility locus itself, since this is the factor that determines the selection criterion. Loci that are phenotypically neutral but are somehow associated with the disease-susceptibility locus, such as genetic markers in linkage disequilibrium with the disease-susceptibility locus, also experience disproportionate genotype selection. As the degree of association between disease susceptibility and marker loci decreases, HW disequilibrium at the marker locus is also expected to decrease. We have examined measures that capture this relationship, potentially offering fine-mapping techniques that can be performed on samples of affected individuals when an appropriate control sample is not available. For a general disease model that considers an arbitrary number of alleles at the disease-susceptibility locus, we have proposed the measures δ_i and δ_{ij} . These are summary measures, useful in quantifying the linkage disequilibria between marker allele M_i and the disease-susceptibility alleles. Since these measures allow for a simple expression of the marker-allele frequencies within affected and unaffected individuals, they can be readily incorporated into many established measures of association. This allows for simpler interpretation of these measures.

Under certain disease models, and when physical linkage and linkage disequilibrium exist between the markers and a disease-susceptibility locus, conventional tests for HW disequilibrium at marker loci can be used to fine map disease-susceptibility loci. For biallelic locus models (in which both the disease-susceptibility locus and the marker loci have only two alleles per locus), HW disequilibrium is proportional to the square of linkage disequilibrium (Appendix B). This indicates that measures of HW disequilibrium are expected to decay more rapidly than direct measures of linkage disequilibrium as linkage disequilibrium diminishes. The results of Feder et al. (1996) illustrate this: their curve plotting the marker map versus HW equilibrium is sharper than

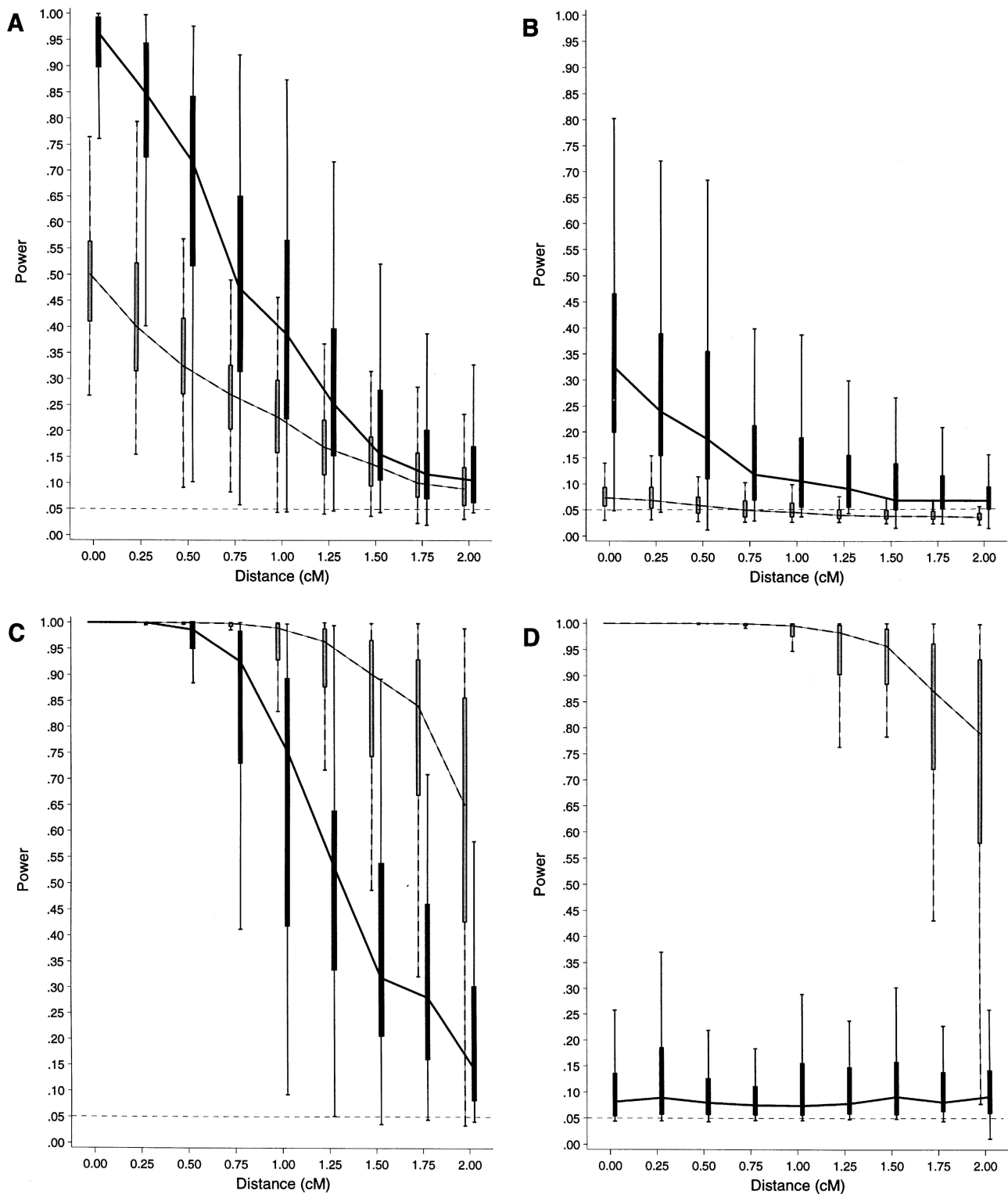


Figure 1 Power results for the χ^2 tests for linkage disequilibrium (gray-shaded boxes) and HW disequilibrium (blackened boxes). A and B, first and second heterogeneous recessive models of Feder et al. (1996); C, additive model; and D, multiplicative model. The symbols represent the range of the proportions of times the hypothesis of no disequilibrium was rejected for the 100 populations. The bottom and top edges of the box represent the sample 25th and 75th percentiles, the point joined by the connecting line is the median, and the whiskers extend the range of the results.

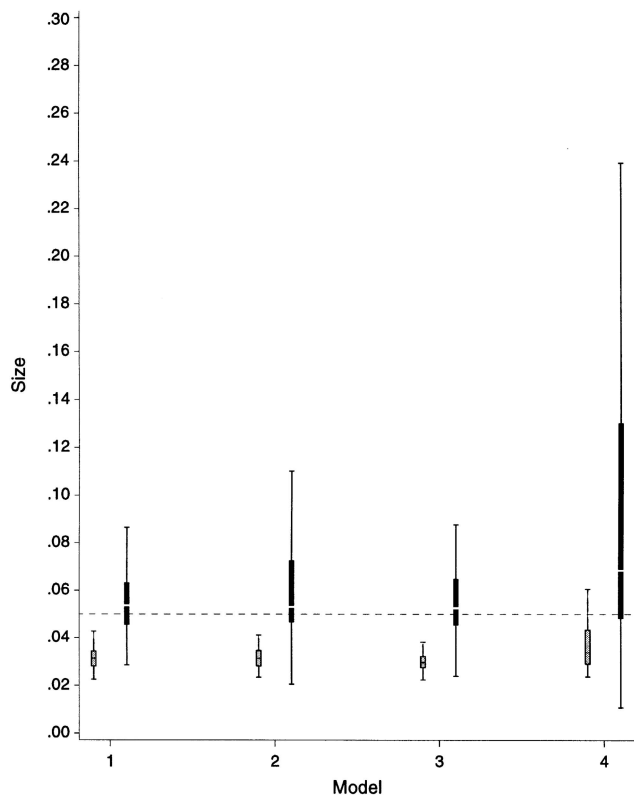


Figure 2 Size of the χ^2 tests for linkage disequilibrium (gray-shaded boxes) and HW disequilibrium (blackened boxes). Models 1 and 2 were the first and second heterogeneous recessive models of Feder et al. (1996), model 3 was the additive model, and the multiplicative model was model 4. These symbols represent the proportion of times a true null hypothesis was rejected.

their curve plotting the marker map versus p_{excess} , a measure of linkage disequilibrium. For a general disease model, allowing for two or more alleles at the marker and disease-susceptibility loci, the relationship between linkage disequilibrium and HW disequilibrium becomes less clear. However, departure from HW equilibrium at a marker locus provides evidence both for linkage disequilibrium between marker and susceptibility loci and for heterogeneity of disease susceptibility, so tests for HW disequilibrium could still be useful.

There are some caveats that should be considered when a general disease model is examined. For simple biallelic-locus models, the interpretation of disequilibrium measures is straightforward. However, when more than two alleles exist at the marker and/or disease-susceptibility locus, complications arise. For these models, summary measures may be used to quantify association between loci; however, although disequilibria between specific alleles may exist, these disequilibria may cancel out when combined into the summary measure. This poses a challenge in the mapping of loci involved in

complex traits, since it is doubtful that many of the traits of interest are biallelic. With the use of SNPs as genetic markers, some of the problems with multiple alleles disappear. In this case, $\delta_2 = -\delta_1$. However, if there are more than two alleles at the disease-susceptibility locus, then the problem of the disequilibria between the disease-susceptibility alleles and the marker allele canceling to within δ_1 is still a concern.

Tests for HW disequilibrium will be the most powerful when large amounts of disequilibrium within a sample of affected individuals are expected. The amount of HW disequilibrium expected depends on both the degree to which the disease-susceptibility locus affects disease status and the manner in which the alleles within a genotype interact. In the sampling of affected individuals, genotypes will be sampled proportionally to the rate of disease susceptibility that they confer. By definition, HW disequilibrium is the difference between genotype proportions and the product of the proportions of the composite alleles. If the alleles within a genotype act in a multiplicative manner to cause increased levels of disease susceptibility, the genotypes are expected to be selected proportionally to the product of the allele frequencies. Thus, with disease models in which the alleles act in a multiplicative manner, HW disequilibrium is not expected to be created in the sample. The more the effects of alleles deviate from multiplicative interactions, the greater the amount of HW disequilibrium that is expected. This is seen in the theoretical results of Appendix B and is illustrated in the results of the simulations that we have performed. We note that, if the penetrances ϕ_{rs} are regarded as genotypic values, much of the theory in this paper can be applied to the study of quantitative traits.

Acknowledgments

This work was supported in part by National Institutes of Health grant GM45344 to North Carolina State University. Particular thanks are due to Dr. Michael Wagner, Glaxo Wellcome, Inc., for suggesting that we investigate more carefully the methods discussed.

Appendix A

For disease susceptibility-locus homozygotes, the two-locus genotypes and their frequencies are

$$A_r A_r M_i M_i \quad (p_r q_i + D_{ri})^2,$$

$$A_r A_r M_i M_j \quad 2(p_r q_i + D_{ri})(p_r q_j + D_{rj}), \quad i \neq j;$$

those for disease susceptibility-locus heterozygotes are

$$A_r A_s M_i M_i \quad 2(p_r q_i + D_{ri})(p_s q_i + D_{si}), \quad r \neq s,$$

$$A_r A_s M_i M_j \quad 2(p_r q_i + D_{ri})(p_s q_j + D_{sj})$$

$$+ 2(p_r q_j + D_{rj})(p_s q_i + D_{si}), \quad r \neq s, i \neq j.$$

Among affected people, therefore, the marker genotype frequencies are

$$P_{ii|Affected} = \frac{1}{\phi} \sum_r \sum_s \phi_{rs} (p_r q_i + D_{ri})(p_s q_i + D_{si})$$

$$= q_i^2 + \frac{2q_i \delta_i}{\phi} + \frac{\delta_{ii}}{\phi},$$

$$P_{ij|Affected} = \frac{1}{\phi} \sum_r \sum_s \phi_{rs} [(p_r q_i + D_{ri})(p_s q_j + D_{sj})$$

$$+ (p_r q_j + D_{rj})(p_s q_i + D_{si})]$$

$$= 2q_i q_j + \frac{2(q_i \delta_j + q_j \delta_i)}{\phi} + \frac{2\delta_{ij}}{\phi}, \quad i \neq j,$$

where $\delta_i = \sum_r \sum_s p_s \phi_{rs} D_{ri} = \sum_i \phi_r D_{ri}$ and $\delta_{ij} = \sum_r \sum_s \phi_{rs} D_{ri} D_{sj}$. Adding over genotypes provides the marker allele frequencies:

$$q_i|Affected = P_{ii|Affected} + \frac{1}{2} \sum_{j \neq i} P_{ij|Affected}$$

$$= \sum_j \left[q_i q_j + \frac{(q_i \delta_j + q_j \delta_i)}{\phi} + \frac{\delta_{ij}}{\phi} \right]$$

$$= q_i + \frac{\delta_i}{\phi}.$$

Appendix B

Heterogeneous Recessive Model

For disease susceptibility-locus alleles A and \bar{A} , marker alleles M and \bar{M} , $\phi_{AA} = 1$, and $\phi_{A\bar{A}} = \phi_{\bar{A}\bar{A}} = \psi$,

$$\phi = p_A^2 + \psi(1 - p_A^2),$$

$$\delta_M = (1 - \psi)p_A D_{AM},$$

$$\delta_{MM} = (1 - \psi)D_{AM}^2,$$

$$\mathcal{D}_{MM|Affected} = \frac{\psi(1 - \psi)D_{AM}^2}{\phi^2} \geq 0.$$

General Biallelic Model

For disease susceptibility-locus alleles A and \bar{A} and marker alleles M and \bar{M} ,

$$\delta_M = [p_A(\phi_{AA} - \phi_{A\bar{A}})$$

$$+ (1 - p_A)(\phi_{A\bar{A}} - \phi_{\bar{A}\bar{A}})]D_{AM},$$

$$\delta_{MM} = (\phi_{AA} - 2\phi_{A\bar{A}} + \phi_{\bar{A}\bar{A}})D_{AM}^2,$$

$$\mathcal{D}_{MM|Affected} = \frac{(\phi_{AA}\phi_{\bar{A}\bar{A}} - \phi_{A\bar{A}}^2)D_{AM}^2}{\phi^2}.$$

Additive Susceptibilities

If $\phi_{rs} = \alpha_r + \alpha_s$, then

$$\phi = 2 \sum_r \alpha_r p_r,$$

$$\delta_i = \sum_r \alpha_r D_{ri},$$

$$\delta_{ij} = 0,$$

$$\mathcal{D}_{ij|Affected} = - \left(\frac{\sum_r \alpha_r D_{ri}}{2 \sum_r \alpha_r p_r} \right) \left(\frac{\sum_s \alpha_s D_{sj}}{2 \sum_s \alpha_s p_s} \right) \leq 0, \text{ if } r = s.$$

Multiplicative Susceptibilities

If $\phi_{rs} = \alpha_r \alpha_s$, then

$$\phi = \left(\sum_r \alpha_r p_r \right)^2,$$

$$\delta_i = \left(\sum_r \alpha_r p_r \right) \left(\sum_r \alpha_r D_{ri} \right),$$

$$\delta_{ij} = \left(\sum_r \alpha_r D_{ri} \right) \left(\sum_s \alpha_s D_{sj} \right),$$

$$\mathcal{D}_{ij|Affected} = 0.$$

References

- Bengtsson BO, Thomson G (1981) Measuring the strength of associations between HLA antigens and diseases. *Tissue Antigens* 18:356-363
- Bennett ST, Lucassen AM, Gough SCL, Powell EE, Undlien DE, Pritchard LE, Merriman ME, et al (1995) Susceptibility to human type 1 diabetes at IDDM2 is determined by tan-

- dem repeat variation at the insulin gene mini satellite locus. *Nat Genet* 9:284-292
- Boehnke, M (1994) Limits of resolution of genetic linkage studies: implications for the positional cloning of human disease genes. *Am J Hum Genet* 55:379-390
- Corder EH, Saunders AM, Strittmatter, Schmechel DE, Gaskell PC, Small GW, Roses AD, et al (1993) Gene dose of apolipoprotein E type 4 allele and the risk of Alzheimer's disease in late onset families. *Science* 261:921-923
- Feder JN, Gnirke A, Thomas W, Tsuchihashi Z (1996) A novel MHC class I-like gene is mutated in patients with hereditary haemochromatosis. *Nat Genet* 13:399-408
- Glaser B, Chiu KC, Liu L, Anker A, Nestorowicz A, Cox NJ, Landau N, et al (1995) Recombinant mapping of the familial hyperinsulinism gene to an 0.8 cM region on chromosome 11p15.1 and demonstration of a founder effect in Ashkenazi Jews. *Hum Mol Genet* 4:879-886
- Hästbacka J, de la Chapelle A, Kaitila I, Sistonen P, Weaver A, Lander ES (1992) Linkage disequilibrium mapping in isolated founder populations: Diastrophic dysplasia in Finland. *Nat Genet* 2:204-211
- Hästbacka J, de la Chapelle A, Mahtani MM, Clines G, Reeve-Daly MP, Daly M, Hamilton BA, et al (1994) The diastrophic dysplasia gene encodes a novel sulfate transporter: Position cloning by fine-structure linkage disequilibrium mapping. *Cell* 78:1073-1087
- Kaplan NL, Martin ER, Weir BS (1997) Power studies for the transmission/disequilibrium tests with multiple alleles. *Am J Hum Genet* 60:691-702
- Lehesjoki A-E, Koskineemi M, Norio R, Tirrito S, Sistonen P, Lander E, de la Chapelle A (1993) Localization of the EPM1 gene for progressive myoclonus epilepsy on chromosome 21: Linkage disequilibrium allows high resolution mapping. *Hum Mol Genet* 2:1229-1234
- Martin ER, Kaplan NL, Weir BS (1997) Tests for linkage and association in nuclear families. *Am J Hum Genet* 61:439-448
- Meng RC, Chapman DG (1966) The power of chi-square tests for contingency tables. *J Am Stat Assoc* 61:965-975
- Spielman RS, Ewens WJ (1998) A sibship test for linkage in the presence of association: the S-TDT. *Am J Hum Genet* 62:450-458
- Spielman RS, McGinnis RE, Ewens WJ (1993) Transmission test for linkage disequilibrium: the insulin gene region and insulin-dependent diabetes mellitus (IDDM). *Am J Hum Genet* 52:506-516
- Trégouët D-A, Ducimetire P, Tiret L (1997) Testing association between candidate-gene markers and phenotype in related individuals, by use of estimating equations. *Am J Hum Genet* 61:189-199
- Weir BS (1996) *Genetic data analysis II*. Sinauer, Sunderland, MA