

Joint Linkage and Linkage Disequilibrium Mapping in Natural Populations

Rongling Wu* and Zhao-Bang Zeng†

*Department of Statistics, University of Florida, Gainesville, Florida 32611 and †Program in Statistical Genetics, Department of Statistics, North Carolina State University, Raleigh, North Carolina 27695

Manuscript received April 19, 2000

Accepted for publication October 30, 2000

ABSTRACT

A new strategy for studying the genome structure and organization of natural populations is proposed on the basis of a combined analysis of linkage and linkage disequilibrium using known polymorphic markers. This strategy exploits a random sample drawn from a panmictic natural population and the open-pollinated progeny of the sample. It is established on the principle of gene transmission from the parental to progeny generation during which the linkage between different markers is broken down due to meiotic recombination. The strategy has power to simultaneously capture the information about the linkage of the markers (as measured by recombination fraction) and the degree of their linkage disequilibrium created at a historic time. Simulation studies indicate that the statistical method implemented by the Fisher-scoring algorithm can provide accurate and precise estimates for the allele frequencies, recombination fractions, and linkage disequilibria between different markers. The strategy has great implications for constructing a dense linkage disequilibrium map that can facilitate the identification and positional cloning of the genes underlying both simple and complex traits.

WITH improved techniques for high-throughput identification and genotyping of polymorphisms, it has been possible to genotype molecular markers throughout the genome and construct a dense linkage map covering the entire genome (LANDEGREN *et al.* 1998; WANG *et al.* 1998). For species with homozygous inbred lines available, the linkage analysis of markers is based on the recombinations of a particular chromosome region that are created by hybridization between two genetically divergent inbred lines (MATHER and JINKS 1982). Such a strategy can directly provide an estimate for the linkage relationship of markers as measured by recombination fraction, because there is clear information about parental linkage phase between alleles of different markers. However, linkage mapping has two major limitations. First, for closely linked markers, there will be few recombinations in a segregating generation and, hence, a dense linkage map will provide little extra information about the localization of target genes, unless the number of individuals of the generation is very large (DARVASI *et al.* 1993). For example, using linkage analysis, LONG *et al.* (1995) could only map quantitative trait loci (QTL) affecting bristle numbers in *Drosophila* to regions of ~5–10 cM. Second, homozygous inbred lines used to generate the F₁ parents of *a priori* known linkage phases for the traditional linkage analysis (MATHER and JINKS 1982) are virtually unavailable for natural populations. For many study materials

sampled randomly from a natural population, therefore, uncertainty of linkage phase between markers prevents a direct estimate of their recombination fraction.

For natural populations, the degree of nonrandom genetic association or linkage disequilibrium, produced at a historic time by various evolutionary forces such as mutation, drift, selection, and admixture, is estimated to indirectly infer how strongly these markers are linked on the same chromosome. If the linkage disequilibrium of the markers occurred a long time ago, a strong linkage disequilibrium detected may suggest close physical linkage between the markers because linkage disequilibrium decays with time (KAPLAN *et al.* 1995). This principle has tremendous potential for constructing fine-scale linkage disequilibrium maps for cloning the genes that cause complex qualitative or quantitative traits (reviewed in TEMPLETON 1999). At present, a number of theories or techniques have been well established for linkage disequilibrium-based mapping of target genes in natural populations (HÄSTBACKA *et al.* 1992, 1994; RISCH and MERIKANGAS 1996; XIONG and GUO 1997; TERWILLIGER and WEISS 1998; KRUGLYAK 1999; MEUWISSEN and GODDARD 2000).

The success of linkage disequilibrium mapping is the presence of linkage disequilibrium between different loci arising from the covariance of the population frequencies of nonalleles in the same gamete (LYNCH and WALSH 1998). The degree and extent of linkage disequilibrium reflect the evolutionary history of a population and its interactions with different evolutionary forces (HILL and ROBERTSON 1968; NEI and LI 1973; OHTA 1982a,b; EPPERSON and ALLARD 1987; PETERSON *et al.* 1999; FARNIR *et al.* 2000). A number of earlier studies

Corresponding author: Rongling Wu, Department of Statistics, 533 McCarty Hall C, University of Florida, Gainesville, FL 32611.
E-mail: rwu@stat.ufl.edu

have focused on statistical methods for detecting the existence of linkage disequilibrium. A likelihood-based procedure was developed by HILL (1974) to estimate the coefficient of linkage disequilibrium between two loci in a finite random mating population. BROWN (1975) established a theoretical framework for the sample sizes required for detecting the disequilibrium by the use of data on gametic and zygotic frequencies. WEIR and COCKERHAM (1978) suggested a statistical procedure for calculating the power of testing linkage disequilibrium between different loci of multiple alleles. More recently, LUO (1998) and LUO and SUHAI (1999) proposed statistical approaches for testing and estimating linkage disequilibrium between a polymorphic marker and a putative QTL. All of these analyses have laid a necessary foundation for linkage disequilibrium mapping of disease genes in human populations (HÄSTBACKA *et al.* 1992, 1994; COLLINS and MORTON 1998; ESCAMILLA *et al.* 1999; SERVICE *et al.* 1999).

A major problem with current strategies for linkage disequilibrium mapping is that they provide little insight into the mechanistic basis of linkage disequilibrium detected in a natural population. Without such knowledge, however, the genomic localization and cloning of genes based on linkage disequilibrium may not be successful, because a strong linkage disequilibrium detected between two genetic loci may be due to the recent occurrence of disequilibrium rather than a close physical map distance of the two loci. In human genetics, the cause of linkage disequilibrium can be revealed through a combined linkage and linkage disequilibrium analysis, as shown by a transmission/disequilibrium testing (TDT) approach (ALLISON 1997; RABINOWITZ 1997; CAMP 1998). However, TDT is critically relied upon for nuclear family data with complete records for multiple successive generations. This approach therefore cannot be used for genome mapping in many other situations where no nuclear family records are available or for other undomesticated species, such as wildlife and forest trees. For these situations or species, it is essential for developing a powerful approach that needs no nuclear family but can still provide a simultaneous estimate for linkage and linkage disequilibrium between genetic loci of interest.

In this article, we propose a new strategy for detecting linkage and linkage disequilibrium between polymorphic markers in natural populations. The new strategy is expected to provide a new avenue for studying the evolutionary dynamics of population variation and differentiation. Furthermore, as compared to a pure linkage analysis or linkage disequilibrium analysis, the combined use of linkage and linkage disequilibrium analysis methods can greatly enhance the feasibility of high-resolution mapping of genes of interest and their subsequent genetic manipulation. The strategy is presented in two parts, one on dioecious species and the other on monoecious species. Dioecious species including animals, humans, and many forest trees, such as Ginkgo,

poplar, and willow, display a single sex for an individual and, therefore, are predominantly outcrossing. Monoecious species comprising most crop and horticultural plants and forest trees such as pine, fir, and spruce carry both sexes on every individual and could be both self-compatible and outcrossing. We first deal with a simpler dioecious model. A more complicated statistical model for analyzing natural populations of monoecious species will be reported in a forthcoming companion article.

TWO-LOCUS MODEL

Population structure theory: Consider a panmictic natural population of a dioecious species in Hardy-Weinberg equilibrium. In the population, η neutral co-dominant markers $\mathbf{M}^1, \dots, \mathbf{M}^\eta$ are assumed to be segregating. Let an allele at marker \mathbf{M}^i ($i = 1, \dots, \eta$), designated by M_r^i ($r = 1, \dots, n_i$), have population frequency P_r^i , $\sum_{r=1}^{n_i} P_r^i = 1$, with the number of alleles n_i at the marker being arbitrary.

Assume that a second marker \mathbf{M}^j is located on the same chromosome as \mathbf{M}^i , both markers having a recombination fraction θ^{ij} . These two linked markers are genetically associated in the population with the coefficient of gametic linkage disequilibrium between a pair of nonalleles from the two markers denoted by D_{rs}^{ij} ($s = 1, \dots, n_j$). The population frequency of the gamete (haplotype) at the two markers $M_r^i M_s^j$ can be expressed as

$$P_{rs}^{ij} = P_r^i P_s^j + D_{rs}^{ij}, \tag{1}$$

with the constraints of

$$\begin{aligned} -p^i p^j &\leq D_{rs}^{ij} \leq p^i(1 - p^j) && \text{(LEWONTIN 1964)} \\ \sum_{r=1}^{n_i} D_{rs}^{ij} &= \sum_{s=1}^{n_j} D_{rs}^{ij} = 0 && \text{(WEIR and COCKERHAM 1978)}. \end{aligned}$$

The value of D_{rs}^{ij} may be positive or negative depending on whether nonalleles M_r^i and M_s^j are in coupling ($M_r^i M_s^j$ gametes are overrepresented) or repulsion ($M_r^i M_s^j$ gametes are underrepresented) disequilibrium (LYNCH and WALSH 1998). Because random mating is assumed in the population, $n_i n_j$ two-locus gametes unite to form $\frac{1}{4} n_i n_j (n_i + 1)(n_j + 1)$ unique zygotic genotypes, designated by $M_{r_1}^i M_{r_2}^i M_{s_1}^j M_{s_2}^j$, where r_1, r_2 ($r_1 \leq r_2 = 1, \dots, n_i$) and s_1, s_2 ($s_1 \leq s_2 = 1, \dots, n_j$) are the two alleles of zygotic genotype at markers \mathbf{M}^i and \mathbf{M}^j , respectively. The genotype frequency of $M_{r_1}^i M_{r_2}^i M_{s_1}^j M_{s_2}^j$ in the current population is expressed as

$$P_{r_1 r_2 s_1 s_2}^{ij} = w(P_{r_1 s_1}^{ij} P_{r_2 s_2}^{ij} + P_{r_1 s_2}^{ij} P_{r_2 s_1}^{ij}), \tag{2}$$

where w is the indicator variable relating the marker genotypes to their frequencies,

$$w = \begin{cases} 1/2 & \text{if } r_1 = r_2 \text{ and } s_1 = s_2 \\ 1 & \text{if } r_1 = r_2 \text{ and } s_1 \neq s_2 \text{ or } r_1 \neq r_2 \text{ and } s_1 = s_2 \\ 2 & \text{if } r_1 \neq r_2 \text{ and } s_1 \neq s_2. \end{cases}$$

TABLE 1

The frequencies of different genotypes at two biallelic markers M^i and M^j in a natural population, the numbers (H) of genotypes drawn randomly from the population, and the conditional probabilities of the gamete genotypes (haplotypes) given each sampled plant of a particular marker genotype

Genotype	Sampled plant		Gamete (haplotype) produced by each sampled plant			
	Frequency	Number	$M_1^i M_1^j$	$M_1^i M_2^j$	$M_2^i M_1^j$	$M_2^i M_2^j$
$M_1^i M_1^i M_1^j M_1^j$	$P_{1111}^{ij} = (P_{11}^{ij})^2$	H_{1111}^{ij}	1	0	0	0
$M_1^i M_1^i M_1^j M_2^j$	$P_{1112}^{ij} = 2P_{11}^{ij} P_{12}^{ij}$	H_{1112}^{ij}	1/2	1/2	0	0
$M_1^i M_1^i M_2^j M_2^j$	$P_{1122}^{ij} = (P_{12}^{ij})^2$	H_{1122}^{ij}	0	1	0	0
$M_1^i M_2^i M_1^j M_1^j$	$P_{1211}^{ij} = 2P_{11}^{ij} P_{21}^{ij}$	H_{1211}^{ij}	1/2	0	1/2	0
$M_1^i M_2^i M_1^j M_2^j$	$P_{1212}^{ij} = 2(P_{11}^{ij} P_{22}^{ij} + P_{12}^{ij} P_{21}^{ij})$	H_{1212}^{ij}	$\frac{P_{11}^{ij} P_{22}^{ij} - \theta^j D_{rs}^{ij}}{P_{1212}^{ij}}$	$\frac{P_{12}^{ij} P_{21}^{ij} + \theta^j D_{rs}^{ij}}{P_{1212}^{ij}}$	$\frac{P_{12}^{ij} P_{21}^{ij} + \theta^j D_{rs}^{ij}}{P_{1212}^{ij}}$	$\frac{P_{11}^{ij} P_{22}^{ij} - \theta^j D_{rs}^{ij}}{P_{1212}^{ij}}$
$M_1^i M_2^i M_2^j M_2^j$	$P_{1222}^{ij} = 2P_{12}^{ij} P_{22}^{ij}$	H_{1222}^{ij}	0	1/2	0	1/2
$M_2^i M_2^i M_1^j M_1^j$	$P_{2211}^{ij} = (P_{21}^{ij})^2$	H_{2211}^{ij}	0	0	1	0
$M_2^i M_2^i M_1^j M_2^j$	$P_{2212}^{ij} = 2P_{21}^{ij} P_{22}^{ij}$	H_{2212}^{ij}	0	0	1/2	1/2
$M_2^i M_2^i M_2^j M_2^j$	$P_{2222}^{ij} = (P_{22}^{ij})^2$	H_{2222}^{ij}	0	0	0	1

$M_1^i M_2^i M_1^j M_2^j$, two-locus zygotic genotype; $M_1^i M_2^i$, two-locus gametic genotype; P_{1111}^{ij} , the frequency of two-locus zygotic genotype in the current generation; P_{11}^{ij} , the frequency of two-locus gamete producing the current generation; θ^j , the recombination fraction of the two markers; D_{rs}^{ij} , the linkage disequilibrium of the two markers.

If all zygotes can produce gametes for the next generation, there will be a total of $n_i n_j$ gametes for markers M^i and M^j at the entire population level. But different zygotic genotypes produce different types of gametes; only the genotypes heterozygous at both markers generate all types of gametes whose relative frequencies are affected by recombination fraction and linkage disequilibrium. Table 1 gives nine zygotic genotypes, their population frequencies, and the frequencies of gametes they produce for the next generation under a simpler biallelic model (see APPENDIX A for derivations). According to the population genetics theory (NAGYLAKI 1991), the amount of linkage disequilibrium between any two loci is reduced at the rate of recombination frequency after the population mates at random for one generation. The coefficient of linkage disequilibrium in the new generation is changed to be $(1 - \theta^j) D_{rs}^{ij}$. Thus, the gamete frequencies for haplotypes $M_1^i M_1^j$ in the new generation at the entire population level are

$$Q_{rs}^{ij} = P_r^i P_s^j + (1 - \theta^j) D_{rs}^{ij} \quad (3)$$

Further, these gametes are randomly combined to generate the progeny $M_{1_1}^i M_{1_2}^i M_{1_1}^j M_{1_2}^j$, which are contained in seeds for plants. If there is no overlapping in reproduction between the parental and progeny generations, the frequencies of the genotypes at the two markers are the products of the frequencies of the corresponding gametes.

Sampling theory: A sample of H female plants is ran-

domly selected from the population. The seeds of these sampled plants are collected and germinated into seedlings. In traditional quantitative genetics, these seedlings grown in a regular experimental design initiate a progeny test, which serves as the selection of best parents for the next generations (MCKEAND and BRIDGWATER 1998). Both the sampled plants and their progeny are genotyped for molecular markers M^1, \dots, M^n , each of an arbitrary number of alleles. According to the sampling theory, every randomly selected plant from the original population should be one of the $\frac{1}{4} n_i n_j (n_i + 1) (n_j + 1)$ distinguishable genotypes for the two markers M^i and M^j , each genotype with a frequency of $P_{1_1 2_1 1_2 2_2}^{ij}$ (see Equation 2) and a sample size of $H_{1_1 2_1 1_2 2_2}^{ij}$ (Table 1). For dioecious species, offspring genotypes (contained in seeds) of a sampled plant are formed through combing its maternal gametes with paternal gametes from the pollen pool. These offspring virtually represent a half-sib relationship with the common mother and different (unknown) fathers. The relative fractions of different maternal gametes generated by a sampled plant of a particular marker genotype are given for two biallelic markers in Table 1. The frequencies of paternal gametes in the pollen pool at the entire population level are described by Equation 3. Because of different compositions of maternal gametes (Table 1), different marker genotypes of the sampled plants generate different compositions of progeny genotypes. The conditional probabilities ($Q_{1_1 2_1 1_2 2_2}^{R_1 R_2 S_1 S_2}$) of the progeny genotypes at markers

\mathbf{M}^i and \mathbf{M}^j , given a mother plant, can be derived from Bayes' theorem, where the subscripts index the mother plant's genotype and the superscripts index the genotype of a progeny. As a simpler example, these conditional probabilities are given for two biallelic markers in Table 2. Similarly, we use $N_{r_1 r_2 s_1 s_2}^{R_1 R_2 S_1 S_2}$ to denote the number of progeny with a particular genotype collected from a sampled plant.

Estimation theory: The allele frequencies, linkage, and linkage disequilibrium for the markers \mathbf{M}^i and \mathbf{M}^j in the original population can be estimated using the random sample. To estimate these unknown genetic parameters associated with the two markers $\boldsymbol{\pi}^{ij} = (P^i, P^j, \theta^{ij}, D^{ij})^T$, a two-stage hierarchical likelihood function of the marker data (\mathbf{M}) is formulated from the sampled plants and their half-sib families,

$$L(\mathbf{M}|\boldsymbol{\pi}^{ij}) = \prod_{r_1=1}^{n_i} \prod_{r_2=1}^{n_i} \prod_{s_1=1}^{n_j} \prod_{s_2=1}^{n_j} \prod_{\xi=1}^{H_{r_1 r_2 s_1 s_2}^{ij}} P_{r_1 r_2 s_1 s_2}^{ij}(\xi) \times \prod_{R_1=1}^{n_i} \prod_{R_2=1}^{n_i} \prod_{S_1=1}^{n_j} \prod_{S_2=1}^{n_j} \prod_{\zeta=1}^{N_{R_1 R_2 S_1 S_2}^{ij}} Q_{R_1 R_2 S_1 S_2}^{ij}(\zeta), \quad (4)$$

with the restrictions of $r_1 \leq r_2, s_1 \leq s_2, R_1 \leq R_2, S_1 \leq S_2$, where ξ and ζ are the ξ th sampled plant and the ζ th progeny of the sampled plant, respectively, and the other symbols have been defined as above and are given in Tables 1 and 2 when a biallelic model is assumed.

There have been a number of computational algorithms available to obtain maximum-likelihood estimates (MLEs) of the four unknowns. In this article, the Fisher-scoring algorithm based on iterations is employed (EDWARDS 1984) because it is easy to derive and also very fast. In terms of this algorithm, the estimates at the $(\tau + 1)$ th iteration can be expressed by the score function vector $\mathbf{S}(\boldsymbol{\pi}^{ij})$ and Fisher information matrix $\mathbf{I}(\boldsymbol{\pi}^{ij})$ (APPENDIX B). The values at the τ th iteration are modified by adding to them the scores divided by the information, both evaluated at the τ th iteration. This iteration continues until successive iterates differ by less than some specified amount. It is apparent that the appropriateness of the Fisher-scoring algorithm relies upon the condition that the information is not zero or the information matrix is nonsingular. In practice, it is always desirable to try several different starting values and to compare the likelihoods found after convergence. After obtaining the MLEs of the unknown parameters, the inverse of $\mathbf{I}(\hat{\boldsymbol{\pi}}^{ij})$ is calculated to estimate the sampling variances of $\boldsymbol{\pi}^{ij}$.

For the purpose of linkage mapping, the degree of linkage between the two markers under consideration is important and should be tested statistically. The hypotheses for testing for linkage are H_0 (free recombination), $\theta^{ij} = 0.5$ vs. H_1 (linkage), $\theta^{ij} \neq 0.5$. The likelihood-ratio (LR) test statistic has the form of

$$LR_0 = -2 \log \left[\frac{L(\mathbf{M}|\hat{P}^i, \hat{P}^j, \theta^{ij} = 0.5, \hat{D}^{ij})}{L(\mathbf{M}|\hat{\boldsymbol{\pi}}^{ij})} \right], \quad (5a)$$

TABLE 2

The conditional probabilities of the joint genotypes at two biallelic markers \mathbf{M}^i and \mathbf{M}^j in the progeny population produced by the sampled plants of particular marker genotypes drawn randomly from a natural population

Sampled plants	Progeny										
	$M_1^i M_1^j M_1^i M_1^j$	$M_1^i M_1^j M_1^i M_2^j$	$M_1^i M_1^j M_2^i M_1^j$	$M_1^i M_2^j M_1^i M_2^j$	$M_1^i M_2^j M_2^i M_1^j$	$M_1^i M_2^j M_2^i M_2^j$	$M_2^i M_1^j M_1^i M_1^j$	$M_2^i M_1^j M_1^i M_2^j$	$M_2^i M_1^j M_2^i M_1^j$	$M_2^i M_1^j M_2^i M_2^j$	$M_2^i M_2^j M_2^i M_2^j$
$M_1^i M_1^j M_1^i M_1^j$	Q_{11}^{ij}	Q_{12}^{ij}	Q_{21}^{ij}	Q_{22}^{ij}	Q_{32}^{ij}	0	0	0	0	0	0
$M_1^i M_1^j M_1^i M_2^j$	$\frac{1}{2} Q_{11}^{ij}$	$\frac{1}{2} Q_{12}^{ij}$	$\frac{1}{2} Q_{21}^{ij}$	$\frac{1}{2} Q_{22}^{ij}$	$\frac{1}{2} P_2^i$	$\frac{1}{2} Q_{32}^{ij}$	$\frac{1}{2} Q_{32}^{ij}$	0	0	0	0
$M_1^i M_1^j M_2^i M_1^j$	0	Q_{12}^{ij}	0	0	Q_{21}^{ij}	0	0	0	0	0	0
$M_1^i M_1^j M_2^i M_2^j$	$\frac{1}{2} Q_{12}^{ij}$	Q_{12}^{ij}	$\frac{1}{2} P_1^j$	$\frac{1}{2} P_2^j$	$\frac{1}{2} P_1^j$	$\frac{1}{2} P_2^j$	$\frac{1}{2} Q_{32}^{ij}$	$\frac{1}{2} Q_{32}^{ij}$	0	0	0
$M_1^i M_2^j M_1^i M_1^j$	Q_{11}^{ij}	0	Q_{21}^{ij}	0	Q_{31}^{ij}	$(P_1^i P_2^j + \theta^{ij} D^{ij}) Q_{11}^{ij} + (P_1^i P_2^j - \theta^{ij} D^{ij}) Q_{21}^{ij} + (P_1^i P_2^j + \theta^{ij} D^{ij}) Q_{31}^{ij} + (P_1^i P_2^j - \theta^{ij} D^{ij}) Q_{41}^{ij}$	$(P_1^i P_2^j + \theta^{ij} D^{ij}) Q_{11}^{ij} + (P_1^i P_2^j - \theta^{ij} D^{ij}) Q_{21}^{ij}$	$(P_1^i P_2^j + \theta^{ij} D^{ij}) Q_{31}^{ij} + (P_1^i P_2^j - \theta^{ij} D^{ij}) Q_{41}^{ij}$	$(P_1^i P_2^j + \theta^{ij} D^{ij}) Q_{11}^{ij} + (P_1^i P_2^j - \theta^{ij} D^{ij}) Q_{21}^{ij}$	$(P_1^i P_2^j + \theta^{ij} D^{ij}) Q_{31}^{ij} + (P_1^i P_2^j - \theta^{ij} D^{ij}) Q_{41}^{ij}$	$(P_1^i P_2^j + \theta^{ij} D^{ij}) Q_{11}^{ij} + (P_1^i P_2^j - \theta^{ij} D^{ij}) Q_{21}^{ij}$
$M_1^i M_2^j M_1^i M_2^j$	$\frac{1}{2} Q_{11}^{ij}$	$\frac{1}{2} Q_{12}^{ij}$	0	0	$\frac{1}{2} P_1^j$	$\frac{1}{2} P_2^j$	$\frac{1}{2} Q_{32}^{ij}$	$\frac{1}{2} Q_{32}^{ij}$	0	0	0
$M_1^i M_2^j M_2^i M_1^j$	0	0	Q_{21}^{ij}	0	0	0	0	0	0	0	0
$M_1^i M_2^j M_2^i M_2^j$	0	0	$\frac{1}{2} Q_{21}^{ij}$	0	0	0	0	0	0	0	0
$M_2^i M_1^j M_1^i M_1^j$	0	$\frac{1}{2} Q_{12}^{ij}$	0	0	$\frac{1}{2} P_1^i$	$\frac{1}{2} P_2^i$	$\frac{1}{2} Q_{32}^{ij}$	$\frac{1}{2} Q_{32}^{ij}$	0	0	0
$M_2^i M_1^j M_1^i M_2^j$	0	0	0	0	0	0	0	0	0	0	0
$M_2^i M_1^j M_2^i M_1^j$	0	0	0	0	0	0	0	0	0	0	0
$M_2^i M_1^j M_2^i M_2^j$	0	0	0	0	0	0	0	0	0	0	0

$M_1^i, M_2^i, M_1^j, M_2^j$, two-locus zygotic genotype; P_1^i, P_2^i , the frequency of marker \mathbf{M}^i in the current generation; $Q_{11}^{ij}, Q_{12}^{ij}, Q_{21}^{ij}, Q_{22}^{ij}, Q_{31}^{ij}, Q_{32}^{ij}, Q_{41}^{ij}, Q_{42}^{ij}$, the allele frequency of marker \mathbf{M}^i .

where \wedge and \sim denote the MLEs of the unknowns under H_1 and H_0 , respectively. Linkage disequilibrium is an important population genetic parameter and its existence and degree reflect the dynamics of population evolution. The hypotheses for overall linkage disequilibrium between markers \mathbf{M}^i and \mathbf{M}^j can be formulated as H_0 , all $D_{rs}^{ij} = 0$ vs. H_1 , at least $D_{rs}^{ij} \neq 0$, whose LR test statistic is

$$\text{LR}_D = -2 \log \left[\frac{L(\mathbf{M}|\hat{P}_r^i, \hat{P}_s^j, \hat{\theta}^{ij}, \hat{D}_{rs}^{ij} = 0, r = 1, \dots, n_i, s = 1, \dots, n_j)}{L(\mathbf{M}|\hat{\pi}^{ij})} \right]. \quad (5b)$$

These test statistics (5a and 5b) are $\sim\chi^2$ -distributed with 1 d.f. Alternatively, the hypothesis test about linkage disequilibrium can be based on the collapse of marker data into a few alleles. But such a treatment may change the power of the tests for linkage disequilibrium, as demonstrated in WEIR and COCKERHAM (1978).

If the null hypothesis of (5a) is accepted, then a significant linkage disequilibrium detected by (5b) indicates that linkage disequilibrium between a pair of markers is not due to their strong linkage. In this case, results from pure linkage disequilibrium mapping (LUO 1998; LUO and SUHAI 1999; MEUWISSEN and GODDARD 2000) are ineffective for gene mapping. If nonsignificant linkage disequilibrium is detected for two linked markers, although this may be rare, the two markers are still useful for potential localization of a target gene. Thus, by testing simultaneously for the significance of linkage and linkage disequilibrium, our analytical approach increases both the effectiveness and efficiency of gene mapping in a natural population.

MARKER ORDERING

The principle for a joint linkage and linkage disequilibrium analysis of two markers can be extended to include more than two markers. This extension is based on two assumptions: (1) recombination between any two markers is independent of recombination between any other nonoverlapping two, *i.e.*, no crossover interference; and (2) linkage disequilibrium between one pair of markers is independent of disequilibrium between other pairs. When there are more than two markers, the most likely linkage order should give the highest likelihood value for a particular dataset. With the two assumptions described above, we propose a hidden Markov model to determine an optimal order for different markers (see also LANDER and GREEN 1987).

Assume that all η codominant markers are derived from the same chromosome in a randomly mating population. We use $M_{r_i}^i$ ($r_i = 1, \dots, n_i$) and $M_{r_{i1}}^i M_{r_{i2}}^i$ ($r_{i1} \leq r_{i2} = 1, \dots, n_i$) to denote an allele (gamete) and genotype (zygote) from a marker \mathbf{M}^i , with the population frequencies $P_{r_i}^i$ and $P_{r_{i1}r_{i2}}^i$, respectively. For a particular order $\mathbf{M}^1, \dots, \mathbf{M}^i, \dots, \mathbf{M}^n$, $\theta^{i(i+1)}$ is used to denote a recombination fraction between two adjacent markers.

The coefficient of linkage disequilibrium between a pair of nonalleles r_i and r_{i+1} from two adjacent markers is denoted by $D_{r_i r_{i+1}}^{i(i+1)}$. For a vector of unknowns $\boldsymbol{\pi} = \{P_{r_i}^i, \theta^{i(i+1)}, D_{r_i r_{i+1}}^{i(i+1)}\}^T$, a two-stage hierarchical likelihood function is formulated as

$$\begin{aligned} L(\mathbf{M}|\boldsymbol{\pi}) = & \prod_{i=1}^{\eta} \prod_{r_{i1}=1}^{n_i} \prod_{r_{i2}=1}^{n_i} \prod_{r_{(i+1)1}=1}^{n_{i+1}} \prod_{r_{(i+1)2}=1}^{n_{i+1}} \\ & \times \prod_{\xi=1}^{H_{r_{i1}r_{i2}^{(i+1)}r_{(i+1)1}r_{(i+1)2}}} P_{r_{i1}r_{i2}^{(i+1)}r_{(i+1)1}r_{(i+1)2}}^{i(i+1)}(\xi) \prod_{R_{i1}=1}^{n_i} \prod_{R_{i2}=1}^{n_i} \prod_{R_{(i+1)1}=1}^{n_{i+1}} \\ & \times \prod_{R_{(i+1)2}=1}^{n_{i+1}} \prod_{\zeta=1}^{N_{r_{i1}R_{i2}R_{(i+1)1}R_{(i+1)2}}} Q_{r_{i1}R_{i2}R_{(i+1)1}R_{(i+1)2}}^{i(i+1)}(\zeta), \quad (6) \end{aligned}$$

where there are the restrictions $r_{i1} \leq r_{i2}$, $r_{(i+1)1} \leq r_{(i+1)2}$, $R_{i1} \leq R_{i2}$, and $R_{(i+1)1} \leq R_{(i+1)2}$, and $P_{r_{i1}r_{i2}^{(i+1)}r_{(i+1)1}r_{(i+1)2}}^{i(i+1)}$ and $Q_{r_{i1}R_{i2}R_{(i+1)1}R_{(i+1)2}}^{i(i+1)}$ are accordingly defined by Equations 1 and 3.

Similarly, the MLEs of the unknown vector $\boldsymbol{\pi}$ can be obtained by the Fisher-scoring algorithm based on iterations (APPENDIX B). The hypotheses for linkage and linkage disequilibrium for every two adjacent markers can be tested accordingly. Using the Markov chain model (6), we can only estimate the linkage disequilibria between two adjacent markers and ignore the estimates of disequilibria between distant markers. Such a result may be limited from a population genetic perspective, because one cannot detect all possible linkage disequilibria generated by evolutionary forces. However, this result can definitely facilitate genomic localization and cloning of genes because our objective is to use a nearest marker to manipulate a target gene of interest.

RESULTS

To demonstrate the statistical properties of the method proposed in this article, we analyze examples on the basis of simulations. In these examples, plants for seed collection are supposed to be randomly sampled from a natural population in Hardy-Weinberg equilibrium. The effects of different sampling schemes and parameter values on the estimates for unknowns are examined, respectively.

Effects of sampling schemes: Assume that the total number (1000) of the open-pollinated progeny collected from all sampled plants is fixed. Five different sampling schemes are generated by changing the number of the sampled plants (H), each of which corresponds to a half-sib family, and the size of progeny (N) generated by each sampled plant (Table 3). These five schemes represent few large families, many small families, and moderately sized families of a moderate number. Among all the strategies, the value for each of the genetic parameters $P_{r_i}^i$, $P_{r_{i1}r_{i2}}^i$, θ^{ij} , and D_{rs}^{ij} for two hypothesized biallelic markers \mathbf{M}^i and \mathbf{M}^j is set to be equal (Table 3). The generation of the marker data for the H half-

TABLE 3

MLEs of the genetic parameters for two biallelic markers and standard errors of the MLEs (in parentheses) calculated from 100 simulation runs for different sampling schemes

Sampling scheme	M	N	$P_r^i = 0.5$	$P_s^i = 0.3$	$\theta^j = 0.1$		$D_{rs}^j = 0.12$	
					MLE	Power	MLE	Power
1	10	100	0.4991 (0.0012) (0.0012)	0.3007 (0.0008) (0.0009)	0.1008 (0.0007) (0.0009)	0.95	0.1203 (0.0004) (0.0004)	0.96
2	20	50	0.5004 (0.0013) (0.0014)	0.2993 (0.0008) (0.0008)	0.1005 (0.0007) (0.0007)	0.93	0.1202 (0.0005) (0.0005)	0.96
3	32	32	0.5016 (0.0012) (0.0012)	0.3005 (0.0008) (0.0008)	0.0999 (0.0007) (0.0008)	0.92	0.1201 (0.0006) (0.0007)	0.94
4	50	20	0.4994 (0.0016) (0.0017)	0.3010 (0.0007) (0.0007)	0.1006 (0.0008) (0.0010)	0.91	0.1204 (0.0006) (0.0008)	0.92
5	100	10	0.5007 (0.0011) (0.0011)	0.3009 (0.0007) (0.0007)	0.1007 (0.0007) (0.0008)	0.89	0.1201 (0.0005) (0.0006)	0.90

Standard errors presented in upper parentheses are averaged from 100 runs and those in lower parentheses are derived from the Fisher information index on the basis of a single run. P_r^i and P_s^i are the allele frequencies of markers \mathbf{M}^i and \mathbf{M}^j , and θ^j and D_{rs}^j are the recombination fraction and linkage disequilibrium between the two markers.

sib families of equal size N includes the following two steps:

1. Randomly assign the nine joint genotypes at the two markers \mathbf{M}^i and \mathbf{M}^j to the H sampled plants according to multinomial distribution with the probabilities as given in Table 1.
2. Randomly assign the two-marker genotypes to the open-pollinated progeny generated by each sampled plant of a particular marker genotype according to the probabilities of the marker genotypes of the progeny given in Table 2.

Because the estimates for the four unknowns are based on known marker genotypes, a likelihood-based approach has many desirable properties in the rate of convergence to achieve stable MLEs and the accuracy and power to obtain these estimates (HILL 1974). In this simulation, we compare the predicted variances of the estimates for these parameters from the asymptotic variance-covariance matrix of the MLEs to the empirical estimates calculated from the repeated simulations. The estimates of the parameters based on 100 runs are averaged and their standard errors are calculated. The sampling errors of the estimates based on a single run are calculated using the Fisher information index as described in APPENDIX B.

Table 3 illustrates the MLEs for each of the four unknown parameters and two types of standard errors under different sampling schemes. In all situations, regardless of the combinations of family number (H) and size (N), the MLEs of the allele frequencies for two

hypothesized markers using the estimation procedures developed in this article are adequately consistent with their actual values. The same is also true for the MLEs of recombination fraction and linkage disequilibrium between the two markers. Results from statistical tests based on Equations 5a and 5b indicate that the alleles of these two different markers are physically significantly linked and genetically significantly associated in the population.

In this example, the predicted values for standard errors estimated from the inverse of the information matrix are reasonably approximate to their empirical values from multiple simulation runs (Table 3). This may be partly because our parameter estimates are based on complete marker information without missing data (see also LUO and SUHAI 1999). As assessed by these two types of estimates for standard errors, the method proposed here has good precision for estimating the population genetic parameters of molecular markers. The power to detect significant linkage or linkage disequilibrium between the two simulated markers is higher for few families of large sizes than for many families of small sizes. But in all sampling schemes, the power is 0.90 or higher.

Effects of linkage and linkage disequilibrium: In this simulation, we assume five biallelic markers with a known order on the same chromosome. These markers are jointly sampled from a natural population in which allele frequency is set to be $P_{r_i}^i = 0.40$ for each marker. The sampling strategy used is 10 half-sib families and 100 progeny in each family. Different recombination

fractions and linkage disequilibria of two adjacent markers are hypothesized as given in Table 4 and lead to four combination patterns: (1) tight linkage and weak disequilibrium, (2) tight linkage and strong disequilibrium, (3) loose linkage and weak disequilibrium, and (4) loose linkage and strong disequilibrium. We first use separate analyses for every two adjacent markers, which are then followed by a joint analysis combining all the five markers through a Markov chain model. The MLEs for unknown parameters are obtained from a single run and their sampling errors for the estimates are assessed by the inverse of the information matrix.

Generally, the estimates of allele frequency are not much affected by the degrees of linkage and linkage disequilibrium of markers (Table 4), with consistent results from separate and joint analyses. The estimation precision of recombination fraction and linkage disequilibrium can be much increased when two markers are tightly linked or display low nonrandom association between the allelic frequencies of the markers (Table 4). Both accuracy and precision of parameter estimates from a separate analysis are largely reduced when two markers have loose linkage and strong disequilibrium. However, these can be much improved by using a joint analysis of all the five markers based on a Markov model.

DISCUSSION

The originality of the statistical method proposed in this study is a combined use of the current linkage analysis and linkage disequilibrium-based mapping theory to simultaneously estimate genetic map distances and population genetic associations of markers using random samples drawn from a natural population. Linkage analysis looks for coinheritance of different markers or QTL within a chromosomal region, while linkage disequilibrium looks for differences in the frequency of marker alleles between genotypes of a different marker or different categories of a phenotype. The combined analysis not only can overcome the limitations of linkage analysis, as noted in the Introduction, but also can increase the effectiveness and efficiency of linkage disequilibrium mapping aimed at precise estimation of gene location. The new analytical method can be seen as an extension of linkage disequilibrium mapping for human pedigrees with complete family records toward any types of natural populations.

In this article, a mapping model is developed for dioecious plant species. The progeny of random samples collected from a dioecious population form a series of open-pollinated (or half-sib) families each with a common female parent and different male parents. The experimental strategy for including both the sampled plants and their progeny for genome mapping offers a unique opportunity to study the transmission of genes from the parental to progeny generation, which causes the breakdown of linkage through meiotic recombina-

TABLE 4
MLEs (\pm SE) of the genetic parameters for five biallelic markers calculated from a single run

Marker	Combination pattern	Hypothesized value			MLE by separate analysis			MLE by joint analysis		
		P_r^i	$\theta^{i(i+1)}$	D_k^j	\hat{P}_r^i	$\hat{\theta}^{i(i+1)}$	\hat{D}_k^j	\hat{P}_r^i	$\hat{\theta}^{i(i+1)}$	\hat{D}_k^j
M ¹	1	0.40	0.02	0.02	0.4051 \pm 0.0008	0.0201 \pm 0.0004	0.0205 \pm 0.0005	0.4001 \pm 0.0004	0.0204 \pm 0.0004	0.0205 \pm 0.0006
M ²	2	0.40	0.02	0.12	0.4030 \pm 0.0010	0.0206 \pm 0.0012	0.1285 \pm 0.0074	0.4008 \pm 0.0008	0.0193 \pm 0.0009	0.1205 \pm 0.0034
M ³	3	0.40	0.30	0.02	0.4021 \pm 0.0015	0.3102 \pm 0.0084	0.0242 \pm 0.0056	0.4021 \pm 0.0009	0.3012 \pm 0.0014	0.0215 \pm 0.0024
M ⁴	4	0.40	0.30	0.12	0.3974 \pm 0.0019	0.3194 \pm 0.0109	0.1105 \pm 0.0108	0.4006 \pm 0.0007	0.3046 \pm 0.0054	0.1213 \pm 0.0077
M ⁵		0.40			0.4098 \pm 0.0032			0.4035 \pm 0.0011		

Standard errors (SE) are derived from the Fisher information index. Combination patterns include (1) tight linkage and weak disequilibrium, (2) tight linkage and strong disequilibrium, (3) loose linkage and weak disequilibrium, and (4) loose linkage and strong disequilibrium.

tion and, thus, the dissipation of linkage disequilibrium between two markers. Unlike previous strategies for a linkage disequilibrium analysis (HILL 1974; WEIR and COCKERHAM 1978), the new strategy, therefore, can capture the intrinsic relationship between linkage and linkage disequilibrium. In human genetic mapping, simultaneous estimation of linkage and linkage disequilibrium is based on nuclear family data (ALLISON 1997). Such a strategy cannot take advantage of analyzing random samples from a natural population in that one can collect sample sizes as large as those considered in the simulation study. In practice, the new strategy can make an immediate application to many plant species in which progeny tests have been established in the field for a number of years (McKEAND and BRIDGWATER 1998).

Given a fixed sample size, our simulation study has focused on the influence of different allocations of the samples between and within families on parameter estimation. When a sample size is adequately large, for instance, as is that used in our example, the precise estimation of genetic parameters, allele frequencies, linkage, and linkage disequilibrium for markers can be obtained, irrespective of few large families or many small families. Such an advantage for the strategy proposed in this article results from two reasons. First, our mapping analysis is established on the foundation of both parental generation and open-pollinated progeny generation. As a random sample, the parental generation contains as much full information about marker allele frequencies, linkage, and linkage disequilibrium as the original population. Unlike full-sib families, open-pollinated families used in our strategy contain full information not only about marker linkage but also about marker population genetic properties due to the contribution of the paternal gametes (pollen) from the population. Second, our linkage analysis of known marker genotypes includes no missing information, a situation not analogous to QTL mapping in which the genotypes at QTL are unknown.

In this study, we implement the Fisher-scoring algorithm to obtain the MLEs of unknown parameters defining the likelihood function of a marker dataset. The Fisher-scoring algorithm is computationally faster and can be more easily derived (EDWARDS 1984), as compared to the expectation-maximization (EM) algorithm (DEMPSTER *et al.* 1977) or Markov chain Monte Carlo method (MCMC; HOESCHELE *et al.* 1997). Its implementation permits a multiple sampling technique to be used more conveniently. Also, this algorithm can provide the estimates for the asymptotic variances of the parameter estimates. For parameter estimates of known markers with joint genotypes in a multinomial distribution, the asymptotic variances estimated from the Fisher information matrix can adequately describe their sampling errors, especially for large samples, although this may not be an actual case for QTL mapping as seen in KAO and ZENG (1997). In some situations, the calculation of the

sampling errors from the Fisher information matrix may encounter negative definitive or singular problems of the matrix. Although these may not be serious for the linkage mapping of known markers, it is advisable to try several different starting values for the parameters to be estimated.

In our experience, a simple Fisher-scoring algorithm is sufficient for analyzing informative markers of known genotypes. However, for a real dataset, there may be many marker types of different segregating patterns. Some markers may be dominant and others may be incomplete or misscored. Many questions for treating these noninformative markers are still open. For example, can we extract useful information from these markers to globally enhance our joint linkage and linkage disequilibrium analysis throughout an entire genome? If yes, how do we make this more efficient? Because of the involvement of the markers of missing information, the Fisher-scoring algorithm may be insufficient for parameter estimation. The EM algorithm or MCMC methodology should be developed to effectively handle these missing data. In addition, when the idea for a combined linkage and linkage disequilibrium analysis is extended to map QTL of unknown genotypes, which is viewed as a missing data problem, the Fisher-scoring algorithm may be very limited. For QTL mapping, more advanced approaches, such as EM algorithm or MCMC, should be developed. Although these approaches are computationally demanding, they can take account of the distribution of multilocus marker-QTL genotypes and permit investigators to fit different models of variation at the QTL.

One of the major contributions of this study is to derive general formulas for estimating allelic frequencies, recombination fractions, and linkage disequilibria for multiallelic markers in natural populations. A number of molecular experiments have demonstrated that multiple alleles per genetic locus are very common in undomesticated populations, such as forest trees (DEGEN *et al.* 1999). Also, the capacity to detect multiallelic markers is largely enhanced by the development of new biotechnologies such as microsatellites. Analyses of multiallelic markers can be simplified by collapsing them into a few alleles at each locus. But, as found by WEIR and COCKERHAM (1978), this simplification may change the power of detecting linkage disequilibrium and lose some important information about disequilibrium inferences. Thus, it is especially not advisable to use a collapsed set of data when the aim of a study is precise localization of QTL and its subsequent positional cloning.

Our mapping approach here is based on a two-point analysis. We further extend the simple two-point analysis to include all markers from the same chromosome through a Markov model. Such a joint two-point analysis can increase both accuracy and precision of parameter estimation, as demonstrated by a simulation study. For

two markers that are not strongly linked but strongly associated between their allelic frequencies, the two-point analysis excluding other marker information likely has low precision (Table 4). But when other markers are included, the precision of the analysis of these two markers is much increased. Apart from the improvement of the precision of parameter estimation, a joint two-point analysis based on a Markov model can facilitate the ordering of molecular markers on a chromosome (see also LANDER and GREEN 1987). The other means of ordering markers is to develop a multipoint analysis. Although this is mathematically very complicated, the extension of our method to a multipoint analysis is straightforward. Assume that there are three different markers with an order \mathbf{M}^i , \mathbf{M}^j , and \mathbf{M}^k on the same chromosomes. A three-locus gamete $M_r^i M_s^j M_t^k$ ($r = 1, \dots, n_i$, $s = 1, \dots, n_j$, $t = 1, \dots, n_k$) has the frequency in the original population,

$$P_{rst}^{ijk} = P_r^i P_s^j P_t^k + P_r^i D_{st}^{jk} + P_s^j D_{rt}^{ik} + P_t^k D_{rs}^{ij} + D_{rst}^{ijk}$$

where a three-locus linkage disequilibrium D_{rst}^{ijk} is assumed to exist (WEIR 1996). In a three-point analysis, the relationship between linkage disequilibrium and recombination fraction becomes nonlinear when genes are transmitted from parental to progeny generation (R. L. WU, unpublished results). The gamete frequencies of three-locus gametes in the current and next generations can be expanded to formulate a likelihood function, given the data as in (4).

Many of our species are still in wild states and are of great importance in terms of their economical significance and theoretical values of biological research. For example, as evidenced in TANKSLEY and MCCOUCH (1997), a number of favorable disease-resistant genes in crop plants are currently warehoused in natural populations of their wild relative species and can be made useful to humans if the number, effects, and locations of these genes are understood. From an evolutionary perspective, knowledge about the organization and structure of wild populations helps us to understand the genetic mechanisms of population evolution and make reasonable predictions about the dynamic changes of the populations (BARTON 2000). Unfortunately, the gene-level studies of genetic architecture and inheritance mode for a complex trait in natural populations are surprisingly rare as a result of the paucity of powerful tools to effectively analyze these populations. Although the method reported here is developed for dioecious species, its extension to monoecious species is possible, but requires an additional mathematical manipulation on outcrossing rate, a population genetic parameter that describes the relative importance of outcrossing pollination to selfing pollination within the same plant. With such powerful mapping approaches available to different kinds of species, we are close to addressing many theoretical or practical genetic questions in depth for natural populations.

We are grateful to Dr. George Casella, Dr. Bruce Weir, and Dr. Mark Yang for stimulating discussions about this work; Dr. James Hobert and Dr. Kenneth Portier for helpful readings of this manuscript; and Dr. Gary Churchill and two anonymous referees for thoughtful comments on this manuscript. This research is partially supported by a grant (GM 45344) from the National Institutes of Health.

LITERATURE CITED

- ALLISON, D. B., 1997 Transmission disequilibrium tests for quantitative traits. *Am. J. Hum. Genet.* **60**: 676–690.
- BARTON, N. H., 2000 Estimating multilocus linkage disequilibria. *Heredity* **84**: 373–389.
- BROWN, A. D. H., 1975 Sample sizes required to detect linkage disequilibrium between two or three loci. *Theor. Popul. Biol.* **8**: 184–201.
- CAMP, N. J., 1998 Genomewide transmission/disequilibrium testing—consideration of the genotype relative risks at disease loci. *Am. J. Hum. Genet.* **61**: 1424–1430.
- COLLINS, A., and N. E. MORTON, 1998 Mapping a disease locus by allelic association. *Proc. Natl. Acad. Sci. USA* **95**: 1741–1745.
- DARVASI, A., A. WEINREB, V. MINKE, J. I. WELLER and M. SOLLER, 1993 Detecting marker-QTL linkage and estimating QTL gene effect and map location using a saturated genetic map. *Genetics* **134**: 943–951.
- DEGEN, B., R. STREIFF and B. ZIEGENHAGEN, 1999 Comparative study of genetic variation and differentiation of two pedunculate oak (*Quercus robur*) stands using microsatellite and allozyme loci. *Heredity* **83**: 597–603.
- DEMPSTER, A. P., N. M. LAIRD and D. B. RUBIN, 1977 Maximum likelihood from incomplete data via EM algorithm. *J. R. Stat. Soc. Ser. B* **39**: 1–38.
- EDWARDS, A. W. F., 1984 *Likelihood*. Cambridge University Press, Cambridge, UK.
- EPPELSON, B. K., and R. W. ALLARD, 1987 Linkage disequilibrium between allozymes in natural populations of lodgepole pine. *Genetics* **115**: 341–352.
- ESCAMILLA, M. A., L. A. MCINNES, M. SPESNY *et al.*, 1999 Assessing the feasibility of linkage disequilibrium methods for mapping complex traits: an initial screen for bipolar disorder loci on chromosome 18. *Am. J. Hum. Genet.* **64**: 1670–1678.
- FARNIR, F., W. COPPIETERS, J.-J. ARRANZ, P. BERZI *et al.*, 2000 Extensive genome-wide linkage disequilibrium in cattle. *Genome Res.* **10**: 220–227.
- HÄSTBACKA, J., A. DE LA CHAPELLE, I. KAITILA, P. SISTONEN, A. WEAVER *et al.*, 1992 Linkage disequilibrium mapping in isolated founder populations: diastrophic dysplasia in Finland. *Nat. Genet.* **2**: 204–221.
- HÄSTBACKA, J., A. DE LA CHAPELLE, M. MAHTANI, G. CLINES, M. P. REEVE-DALY *et al.*, 1994 The diastrophic dysplasia gene encodes a novel sulfate transporter: positional cloning by fine-structure linkage disequilibrium mapping. *Cell* **78**: 1073–1087.
- HILL, W. G., 1974 Estimation of linkage disequilibrium in randomly mating populations. *Theor. Appl. Genet.* **33**: 54–78.
- HILL, W. G., and A. ROBERTSON, 1968 Linkage disequilibrium in finite populations. *Theor. Appl. Genet.* **33**: 54–78.
- HOESCHELE, I., P. UIMARI, F. E. GRIGNOLA, Q. ZHANG and K. M. GAGE, 1997 Advances in statistical methods to map quantitative trait loci in outbred populations. *Genetics* **147**: 1445–1457.
- KAO, C.-H., and Z.-B. ZENG, 1997 General formulas for obtaining the MLEs and the asymptotic variance-covariance matrix in mapping quantitative trait loci when using the EM algorithm. *Biometrics* **53**: 653–665.
- KAPLAN, N. L., W. G. HILL and B. S. WEIR, 1995 Likelihood methods for locating disease genes in nonequilibrium populations. *Am. J. Hum. Genet.* **56**: 18–32.
- KRUGLYAK, L., 1999 Prospects for whole-genome linkage disequilibrium mapping of common disease genes. *Nat. Genet.* **22**: 139–144.
- LANDEGREN, U., M. NILSSON and P. Y. KWOK, 1998 Reading bits of genetic information: methods for single-nucleotide polymorphism analysis. *Genome Res.* **8**: 769–776.
- LANDER, E. S., and P. GREEN, 1987 Construction of multilocus ge-

- netic linkage maps in human. *Proc. Natl. Acad. Sci.* **84**: 2363–2367.
- LEWONTIN, R. C., 1964 The interaction of selection and linkage. I. General considerations for heterotic models. *Genetics* **49**: 49–67.
- LONG, A. D., S. L. MULLANEY, L. A. REID, J. D. FRY, C. H. LANGLEY *et al.*, 1995 High resolution mapping of genetic factors affecting abdominal bristle number in *Drosophila melanogaster*. *Genetics* **139**: 1273–1291.
- LUO, Z. W., 1998 Detecting linkage disequilibrium between a polymorphic marker locus and a trait locus in natural populations. *Heredity* **80**: 198–208.
- LUO, Z. W., and S. SUHAI, 1999 Estimating linkage disequilibrium between a polymorphic marker locus and a trait locus in natural populations. *Genetics* **151**: 359–371.
- LYNCH, M., and B. WALSH, 1998 *Genetics and Analysis of Quantitative Traits*. Sinauer Associates, Sunderland, MA.
- MATHER, K., and J. L. JINKS, 1982 *Biometrical Genetics*, Ed. 3. Chapman & Hall, New York.
- McKEAND, S. E., and F. E. BRIDGWATER, 1998 A strategy for the third breeding cycle of loblolly pine in the Southeastern US. *Silvae Genet.* **47**: 223–234.
- MEUWISSEN, T. H. E., and M. E. GODDARD, 2000 Fine mapping of quantitative trait loci using linkage disequilibria with closely linked marker loci. *Genetics* **155**: 421–430.
- NAGYLAKI, T., 1991 *Introduction to Theoretical Population Genetics*. Springer-Verlag, Berlin.
- NEI, M., and W.-H. LI, 1973 Linkage disequilibrium in subdivided populations. *Genetics* **75**: 213–219.
- OHTA, T., 1982a Linkage disequilibrium due to random genetic drift in finite subdivided populations. *Proc. Natl. Acad. Sci.* **79**: 1940–1944.
- OHTA, T., 1982b Linkage disequilibrium with the island model. *Genetics* **101**: 139–155.
- PETERSON, R. J., D. GOLDMAN and J. C. LONG, 1999 Effects of worldwide population subdivision on ALDH2 linkage disequilibrium. *Genome Res.* **9**: 844–852.
- RABINOWITZ, D., 1997 A transmission disequilibrium test for quantitative trait loci. *Hum. Hered.* **47**: 342–350.
- RISCH, N., and K. MERIKANGAS, 1996 The future of genetic studies of complex human diseases. *Science* **273**: 1516–1517.
- SERVICE, S. K., D. W. T. LANG, N. B. FREIMER and L. A. SANDKUJIL, 1999 Linkage-disequilibrium mapping of disease genes by reconstruction of ancestral haplotypes in founder populations. *Am. J. Hum. Genet.* **64**: 1728–1738.
- TANKSLEY, S. D., and S. R. MCCOUCH, 1997 Seed banks and molecular maps: unlocking genetic potential from the wild. *Science* **277**: 1063–1066.
- TEMPLETON, A. R., 1999 Uses of evolutionary theory in the human genome project. *Annu. Rev. Ecol. Syst.* **30**: 23–49.
- TERWILLIGER, J. D., and K. M. WEISS, 1998 Linkage disequilibrium mapping of complex disease: fantasy or reality? *Curr. Opin. Biotechnol.* **9**: 578–594.
- WANG, D. G., J. B. FAN, C. J. SIAO *et al.*, 1998 Large-scale identification, mapping, and genotyping of single-nucleotide polymorphisms in the human genome. *Science* **280**: 1077–1082.
- WEIR, B. S., 1996 *Genetic Data Analysis II*. Sinauer Associates, Sunderland, MA.
- WEIR, B. S., and C. C. COCKERHAM, 1978 Testing hypotheses about linkage disequilibrium with multiple alleles. *Genetics* **88**: 633–642.
- XIONG, M. M., and S. W. GUO, 1997 Fine-linkage genetic mapping based on linkage disequilibrium: theory and applications. *Am. J. Hum. Genet.* **60**: 1513–1531.

Communicating editor: G. A. CHURCHILL

APPENDIX A: DERIVATION OF GAMETE FREQUENCIES

When a zygotic genotype is homozygous for both markers M^i and M^j , only one type of gamete is produced and, thus, recombinant and nonrecombinant gametes are mixed. When a zygotic genotype is homozygous for a marker but heterozygous for the other marker, two

types of gametes are produced. In this case, one still cannot distinguish between recombinant and nonrecombinant gametes, because each gamete type is mixed. However, for a genotype that is heterozygous at both markers, four types of gametes can be produced. The frequency for each of the four gamete types produced by a two-marker heterozygous genotype is dependent on the recombination fraction of the markers and the frequency with which the heterozygous genotype was yielded through gamete combination in the previous generation. There are two ways for yielding the heterozygous genotype $M^i_1 M^i_2 M^j_1 M^j_2$ ($r_1 \neq r_2$, $s_1 \neq s_2$): (1) via the combination of two gametes $M^i_1 M^j_1$ and $M^i_2 M^j_2$ and (2) via the combination of two gametes $M^i_1 M^j_2$ and $M^i_2 M^j_1$. These two different ways therefore produce two different diplotypes. The probability of the diplotype produced the first way is $2P^i_{r_1 s_1} P^j_{r_2 s_2}$, whereas the probability of the diplotype produced the second way is $2P^i_{r_1 s_2} P^j_{r_2 s_1}$ (see Equation 2). The frequencies of the four gamete types produced by the heterozygous genotype are $\frac{1}{2}(1 - \theta^i)$, $\frac{1}{2}\theta^i$, $\frac{1}{2}\theta^j$, and $\frac{1}{2}(1 - \theta^j)$ in the first way and $\frac{1}{2}\theta^j$, $\frac{1}{2}(1 - \theta^j)$, $\frac{1}{2}(1 - \theta^i)$, and $\frac{1}{2}\theta^i$ in the second way for $M^i_1 M^j_1$, $M^i_1 M^j_2$, $M^i_2 M^j_1$, and $M^i_2 M^j_2$, respectively. Thus, it is not difficult to derive the frequencies of the four gametes $M^i_1 M^j_1$, $M^i_1 M^j_2$, $M^i_2 M^j_1$, and $M^i_2 M^j_2$ produced by a heterozygous genotype in the entire population as $P^i_{r_1 s_1} P^j_{r_2 s_2} - \theta^j D^j$, $P^i_{r_1 s_2} P^j_{r_2 s_1} + \theta^j D^j$, $P^i_{r_1 s_2} P^j_{r_2 s_1} + \theta^j D^j$, and $P^i_{r_1 s_1} P^j_{r_2 s_2} - \theta^j D^j$ (see Table 1), respectively, where $D^j = P^i_{r_1 s_1} P^j_{r_2 s_2} - P^i_{r_1 s_2} P^j_{r_2 s_1}$. The conditional probabilities of these gametes are derived according to Bayes' theorem (see Table 1).

APPENDIX B: FISHER-SCORING ALGORITHMS FOR OBTAINING MLES OF π

For the Fisher-scoring algorithm based on iterative steps, the estimates at the $(\tau + 1)$ th iteration can be expressed by

$$\pi^{i,j(\tau+1)} = \pi^{i,j(\tau)} + \mathbf{I}^{-1}(\pi^{i,j(\tau)})\mathbf{S}(\pi^{i,j(\tau)}),$$

where

$$\begin{aligned} \mathbf{S}(\pi^{i,j(\tau)}) &= \frac{\partial}{\partial \pi^j} \ln L(\mathbf{M}|\pi^j) = (S_{P_r^i}, S_{P_s^i}, S_{\theta^j}, S_{D_{rs}^j})^T \\ &= \left(\frac{\partial \ln L(\mathbf{M}|\pi^j)}{\partial P_r^i}, \frac{\partial \ln L(\mathbf{M}|\pi^j)}{\partial P_s^i}, \frac{\partial \ln L(\mathbf{M}|\pi^j)}{\partial \theta^j}, \frac{\partial \ln L(\mathbf{M}|\pi^j)}{\partial D_{rs}^j} \right)^T \end{aligned}$$

is the score function, and

$$\mathbf{I}(\pi^j) = -E[\mathbf{S}(\pi^j)\mathbf{S}(\pi^j)]$$

$$= - \begin{bmatrix} E\left(\frac{\partial^2 \ln L(\mathbf{M}|\pi^j)}{\partial (P_r^i)^2}\right) & E\left(\frac{\partial^2 \ln L(\mathbf{M}|\pi^j)}{\partial P_r^i \partial P_s^i}\right) \\ E\left(\frac{\partial^2 \ln L(\mathbf{M}|\pi^j)}{\partial P_s^i \partial P_r^i}\right) & E\left(\frac{\partial^2 \ln L(\mathbf{M}|\pi^j)}{\partial (P_s^i)^2}\right) \\ E\left(\frac{\partial^2 \ln L(\mathbf{M}|\pi^j)}{\partial \theta^j \partial P_r^i}\right) & E\left(\frac{\partial^2 \ln L(\mathbf{M}|\pi^j)}{\partial \theta^j \partial P_s^i}\right) \\ E\left(\frac{\partial^2 \ln L(\mathbf{M}|\pi^j)}{\partial D_{rs}^j \partial P_r^i}\right) & E\left(\frac{\partial^2 \ln L(\mathbf{M}|\pi^j)}{\partial D_{rs}^j \partial P_s^i}\right) \end{bmatrix}$$

$$\begin{bmatrix} E\left(\frac{\partial^2 \ln L(\mathbf{M}|\boldsymbol{\pi}^j)}{\partial P_i^j \partial \theta^j}\right) & E\left(\frac{\partial^2 \ln L(\mathbf{M}|\boldsymbol{\pi}^j)}{\partial P_i^j \partial D_{rs}^j}\right) \\ E\left(\frac{\partial^2 \ln L(\mathbf{M}|\boldsymbol{\pi}^j)}{\partial P_s^j \partial \theta^j}\right) & E\left(\frac{\partial^2 \ln L(\mathbf{M}|\boldsymbol{\pi}^j)}{\partial P_s^j \partial D_{rs}^j}\right) \\ E\left(\frac{\partial^2 \ln L(\mathbf{M}|\boldsymbol{\pi}^j)}{\partial (\theta^j)^2}\right) & E\left(\frac{\partial^2 \ln L(\mathbf{M}|\boldsymbol{\pi}^j)}{\partial \theta^j \partial D_{rs}^j}\right) \\ E\left(\frac{\partial^2 \ln L(\mathbf{M}|\boldsymbol{\pi}^j)}{\partial D_{rs}^j \partial \theta^j}\right) & E\left(\frac{\partial^2 \ln L(\mathbf{M}|\boldsymbol{\pi}^j)}{\partial (D_{rs}^j)^2}\right) \end{bmatrix}$$

is the Fisher information matrix. More specifically, the score function and the Fisher information index can be derived using

$$\mathbf{S}(\boldsymbol{\pi}^{j|\pi}) = \frac{\partial}{\partial \boldsymbol{\pi}^j} \ln L(\mathbf{M}|\boldsymbol{\pi}^j)$$

$$\begin{aligned} &= \prod_{r_1=1}^{n_i} \prod_{r_2=1}^{n_i} \prod_{s_1=1}^{n_j} \prod_{s_2=1}^{n_j} H_{r_1 r_2 s_1 s_2}^{ij} \frac{\partial P_{r_1 r_2 s_1 s_2}^{ij}}{\partial \boldsymbol{\pi}^j} \\ &\times \prod_{R_1=1}^{n_i} \prod_{R_2=1}^{n_i} \prod_{S_1=1}^{n_j} \prod_{S_2=1}^{n_j} H_{R_1 R_2 S_1 S_2}^{ij} \frac{\partial P_{R_1 R_2 S_1 S_2}^{ij}}{\partial \boldsymbol{\pi}^j}, \end{aligned}$$

$$\mathbf{I}(\boldsymbol{\pi}^j) = -E\left[\frac{\partial^2}{\partial (\boldsymbol{\pi}^j)^2} \ln L(\mathbf{M}|\boldsymbol{\pi}^j)\right]$$

$$\begin{aligned} &= -\prod_{r_1=1}^{n_i} \prod_{r_2=1}^{n_i} \prod_{s_1=1}^{n_j} \prod_{s_2=1}^{n_j} H_{r_1 r_2 s_1 s_2}^{ij} \frac{1}{P_{r_1 r_2 s_1 s_2}^{ij}} \left[\frac{\partial P_{r_1 r_2 s_1 s_2}^{ij}}{\partial \boldsymbol{\pi}^j}\right]^2 \\ &\times \prod_{R_1=1}^{n_i} \prod_{R_2=1}^{n_i} \prod_{S_1=1}^{n_j} \prod_{S_2=1}^{n_j} H_{R_1 R_2 S_1 S_2}^{ij} \frac{1}{Q_{R_1 R_2 S_1 S_2}^{ij}} \left[\frac{\partial Q_{R_1 R_2 S_1 S_2}^{ij}}{\partial \boldsymbol{\pi}^j}\right]^2. \end{aligned}$$