

Z. Hu · X. Zhang · C. Xie · G. R. McDaniel
D. L. Kuhlers

A correlation method for detecting and estimating linkage between a marker locus and a quantitative trait locus using inbred lines

Received: 19 April 1994 / Accepted: 22 November 1994

Abstract The advent of molecular genetic markers has stimulated interest in detecting linkage between a marker locus and a quantitative trait locus (QTL) because the marker locus, even without direct effect on the quantitative trait, could be useful in increasing the response to selection. A correlation method for detecting and estimating linkage between a marker locus and a QTL is described using selfing and sib-mating populations. Computer simulations were performed to estimate the power of the method, the sample size (N) needed to detect linkage, and the recombination value (r). The power of this method was a function of the expected recombination value $E(r)$, the standardized difference (d) between the QTL genotypic means, and N. The power was highest at complete linkage, decreased with an increase in $E(r)$, and then increased at $E(r)=0.5$. A larger d and N led to a higher power. The sample size needed to detect linkage was dependent upon $E(r)$ and d . The sample size had a minimum value at $E(r)=0$, increased with an increase in $E(r)$ and a decrease in d . In general, the r was overestimated. With an increase in d , the r was closer to its expectation. Detection of linkage by the proposed method under incomplete linkage was more efficient than estimation of recombination values. The correlation method and the method of comparison of marker-genotype means have a similar power when there is linkage, but the former has a slightly higher power than the latter when there is no linkage.

Key words Genetic marker · Inbred line · Linkage · Quantitative trait locus (QTL)

Introduction

Use of molecular genetic markers has encouraged the development of new methods for detecting linkage between markers and quantitative trait loci (QTLs). Several methods have been described to detect and estimate linkage for various population types: including F_2 progeny (Weller 1986; Lander and Botstein 1989; Simpson 1989; Zhang et al. 1992), double-haploid (Snape 1988; Knapp 1991), recombinant-inbred (Knapp 1991), backcross (Knapp 1991) and replicated progeny (Soller and Beckman 1990), as well as single-seed-descent populations (Snape 1988).

Weller (1986) described a maximum-likelihood method to detect linkage. Simpson (1989) proposed a method to detect linkage by comparing marker-genotype means. The power and sample size needed to detect linkage by these two methods are a function of the recombination value (r) between two linked loci and the standardized difference (d) between genotypic means at a QTL. Generally, the maximum-likelihood method has a higher power and requires a smaller sample size than the method of comparison of marker-genotype means (Simpson 1989). In some cases, however, it is difficult to obtain maximum-likelihood estimates even if the function can be constructed (Weller 1987; Zhang et al. 1992). Lander and Botstein (1989) and Knapp et al. (1990) developed a method that uses two flanking markers instead of a conventional individual marker to map QTLs. This approach can map a QTL in one experiment but requires large sample sizes to obtain the double-crossovers needed to map the QTL in flanking-marker methods (Zhang et al. 1992). The efficiency of detecting the impact of a QTL on the trait mean is increased by selectively genotyping the most extreme individuals for the trait (Lander and Botstein 1989). However, the variance of the trait is biased up if only extreme individuals are genotyped for the quantitative trait of interest (Weller and Wyler 1992).

Alabama Agricultural Experiment Station Journal No. 12-944766

Communicated by L. d. Van Vleck

Z. Hu
Department of Biology, Wuhan University, Wuhan, Hubei, PRC

X. Zhang (✉) · G. R. McDaniel
Department of Poultry Science, Auburn University, Auburn,
AL 36849, USA

C. Xie
Department of Agronomy and Soils, Auburn University, Auburn,
AL 36849, USA

D. L. Kuhlers
Department of Animal and Dairy Sciences, Auburn University,
Auburn, AL 36849, USA

In the present paper an alternative method is presented to detect and estimate linkage between a marker locus and a QTL in selfing and sib-mating populations, to estimate the power of this method, and to determine the sample size needed to detect linkage.

Materials and methods

Consider that a population homozygous at a marker locus with a genotype *MM* and a QTL with genotype *QQ* is mated to another population homozygous at the marker locus with the genotype *mm* and a QTL with genotype *qq* to produce a F_1 population with genotype *MmQq*. Selfing in plants, or sib-mating in animals, for multiple generations from the F_1 progeny eventually leads to four types of inbred lines (*MMQQ*, *MMqq*, *mmQQ*, and *mmqq*) homozygous at the marker locus and the QTL. The genotypic frequency (*f*) in the population (across lines) is a function of *r* (Haldane and Waddington 1931). Genotype *MM* is scored as 1 and *mm* as -1. Genotypic values of *QQ* and *qq* are denoted by *a* and -*a*, respectively. When the marker locus and the QTL are linked in the coupling phase, the genetic constitution of the inbred-line populations is as presented in Table 1.

Detection of linkage

The simple correlation coefficient (*c*) between the scores at the marker locus and genotypic values at the QTL (Table 1) has an expected value $E(c)$:

$$E(c) = 1 - 2f = \frac{1 - 2r}{1 + 2r} \quad (1)$$

for selfing populations, and

$$E(c) = 1 - 2f = \frac{1 - 2r}{1 + 6r} \quad (2)$$

for sib-mating populations. The sampling variance (S_c^2) of *c* is

$$S_c^2 = \frac{1 - c^2}{N - 2} \quad (3)$$

$$t(c) = \frac{c}{\sqrt{\frac{1 - c^2}{N - 2}}} = \sqrt{\frac{(1 - 2r)^2}{8r}} (N - 2) \quad (4)$$

where *N* is the sample size. The *t*-statistic can be used to detect linkage between the two loci by testing the significance of the difference between *c* and zero, for selfing populations and

$$t(c) = \frac{c}{\sqrt{\frac{1 - c^2}{N - 2}}} = \sqrt{\frac{(1 - 2r)^2 (N - 2)}{(1 + 6r)^2 - (1 - 2r)^2}} \quad (5)$$

for sib-mating populations.

Table 1 Genetic constitution of inbred-line populations

Genotype	Frequency (<i>f</i>)	Genotypic value of marker locus (<i>x</i>)	Genotype value of QTL (<i>y</i>)
<i>MMQQ</i>	1/2 (1- <i>f</i>) ^a	1	<i>a</i>
<i>MMqq</i>	1/2 <i>f</i>	1	- <i>a</i>
<i>mmQQ</i>	1/2 <i>f</i>	-1	<i>a</i>
<i>mmqq</i>	1/2 (1- <i>f</i>)	-1	- <i>a</i>

^a $f = 2r/(1+2r)$ for selfing populations and $4r/(1+6r)$ for sib-mating populations (Haldane and Waddington 1931), where *r* is recombination value between the two loci

Sample size (*N*) needed to detect linkage

When the marker locus and the QTL are linked, the sample size *N* required to detect linkage at α significance level can be obtained from (4) and (5) as

$$N = \left(\frac{1}{c^2} - 1 \right) t_{\alpha}^2 + 2 = \frac{8r}{(1 - 2r)^2} t_{\alpha}^2 + 2 \quad (6)$$

for selfing populations, and

$$N = \left(\frac{1}{c^2} - 1 \right) t_{\alpha}^2 + 2 = \frac{16r(1 + 2r)}{(1 - 2r)^2} t_{\alpha}^2 + 2 \quad (7)$$

for sib-mating populations, respectively, where t_{α} is a normal approximation to Student's *t*-distribution at α significance level.

Estimation of recombination value

The estimate of *r* can be obtained from (1) and (2) as

$$r = \frac{1 - c}{2(1 + c)} \quad (8)$$

for selfing populations, and

$$r = \frac{1 - c}{2(1 + 3c)} \quad (9)$$

for sib-mating populations, respectively. The approximate variance of *r* is derived from Taylor's expansion (Chiang 1980) (see Appendix):

$$V(r) = \frac{1 - c^2}{(N - 2)(1 + c)^4} = \frac{r(1 + 2r)^2}{2(N - 2)} \quad (10)$$

for selfing populations, and

$$V(r) = \frac{4(1 - c^2)}{(N - 2)(1 + 3c)^4} = \frac{r(1 + 2r)(1 + 6r)^2}{4(N - 2)} \quad (11)$$

for sib-mating populations. Note that if the marker-genotype *MM* is coded as -1 and *mm* as 1, $E(c) = -(1 - 2f)$ and the *c* value in equations (8), (9), (10), and (11) needs be replaced by the absolute value of *c*.

Monte Carlo simulations and statistical analysis

Sets of data consisting of 100, 200, or 400 observations were simulated for various combinations of the expected recombination value $E(r)$ and the standardized difference between the QTL genotypic means (*d*). The *d* is equal to $(m_Q - m_q)/s$, where m_Q and m_q are means of the genotypic values of *QQ* and *qq*, respectively, and *s* is the standard deviation of the population. The marker-genotypes *MM* and *mm* had equal numbers of observations in each set of data and were scored as 1 and -1, respectively. The phenotypic value for the QTL of each individual was generated from each of the two normally-distributed populations using the NORMAL function of SAS (SAS Institute 1989). One-hundred replicates were run for each combination of $E(r)$ and *d* for a given sample size.

The simple correlation coefficient *c* and recombination value *r* between the two linked loci were computed from each combination. The *t*-statistics (4) and (5) were used to test linkage for selfing and sib-mating populations, respectively. The power was estimated from each of the simulated data sets using a normal approximation to a *t*-distribution. The sample sizes needed to detect linkage at the 5% significance level were calculated based on (6) and (7) from each of the simulated data sets. The power and sample size at the 5% significance level calculated by this method were compared to those by the method of comparison of marker-genotype means. Simpson (1989) used at-statistic $t(m)$ and calculated the sample size *N* required to detect linkage as

$$t(m) = \frac{(1 - 2f)ds}{\sqrt{2s^2(1 + d^2f(1 - f))/N}} \quad (12)$$

$$N = \frac{2t_{\alpha}^2(1+d^2f(1-f))}{(1-2f)^2d^2} \quad (13)$$

where f is the same as in Table 1, d is the standardized difference between the marker-genotype means, and s is the standard deviation of the QTL phenotypic values. The t -statistic is used to detect linkage and calculate the power of the test.

Description of an example

The F_1 plants of common vetch (*Vicia sativa* L.) heterozygous at a QTL that affects plant height were allowed to self-pollinate for seven generations to produce two classes of inbred lines differentiating for plant height (55–75 cm for one class and 76–100 cm for the other). A total of 215 inbred-line plants was measured for a random amplified polymorphic DNA (RAPD) marker. The presence of the RAPD marker was coded as 1 and the absence as -1 (see Table 6).

Results and discussion

Power of linkage tests

The statistical test for linkage is based on H_0 : there is no linkage ($r=0.5$), vs H_1 : There is linkage ($r<0.5$). The power of the test is the probability at which H_1 is accepted when H_1 is true, i.e., $P(\text{accept } H_1 | H_1 \text{ is true})$. The power of the method presented here was a function of $E(r)$, d , and N (Table 2). The power was highest at complete linkage [$E(r)=0$], decreased with an increase in $E(r)$ (Table 2), and then increased at $E(r)=0.5$ (Table 3). A larger d and N led to a higher power (Table 2).

The correlation method was compared to the method of comparison of marker-genotype means (Simpson 1989) in terms of the power at various combinations of $E(r)$ and d (Table 3). The two methods did not differ at complete linkage [$E(r)=0$]. Detection of incomplete linkage by testing the significance of the correlation coefficient was nearly as powerful as by testing the significance of the marker-genotype means at $d > 0.5$. The power of the correlation method and the method of comparison of marker-genotype means was similar at $d \leq 0.5$ and $E(r) < 0.5$. However, the correlation method had a higher power than the method of comparison of marker-genotype means at $E(r)=0.5$ and $d < 2$ (Table 3). In other words, the absence of linkage suggested by the correlation method was more convincing than by the method of comparison of marker-genotype means. The proposed method is useful when $E(r) \leq 0.3$ and $d \geq 1$. However, this method might not produce correct results at $0.3 < E(r) < 0.5$ and $d < 1$. Similarly, the maximum-likelihood method does not yield meaningful parameter estimates for a QTL of genotypic effects smaller than a 1.0 phenotypic standard deviation (Weller 1986).

Sample size (N) needed to detect linkage

The sample size needed to detect linkage was dependent upon $E(r)$ and d . The sample size was minimum at $E(r)=0$, increased with an increase in $E(r)$ and a decrease in d (Ta-

Table 2 The power to detect linkage in different sample sizes (100, 200, and 400) using the correlation method at the 5% significance level

Mating type	d^a	Sample size	Expected recombination value			
			0.0	0.1	0.2	0.3
Selfing	0.50	100	0.81	0.47	0.23	0.11
		200	0.96	0.71	0.43	0.21
		400	1.00	0.97	0.61	0.37
	1.00	100	1.00	0.95	0.57	0.21
		200	1.00	1.00	0.87	0.53
		400	1.00	1.00	1.00	0.68
	2.00	100	1.00	1.00	0.97	0.43
		200	1.00	1.00	1.00	0.85
		400	1.00	1.00	1.00	1.00
	3.00	100	1.00	1.00	1.00	0.67
		200	1.00	1.00	1.00	0.99
		400	1.00	1.00	1.00	1.00
Sib-mating	0.50	100	0.81	0.24	0.14	0.09
		200	0.98	0.52	0.30	0.15
		400	1.00	0.78	0.38	0.25
	1.00	100	1.00	0.72	0.24	0.12
		200	1.00	0.98	0.60	0.22
		400	1.00	1.00	0.79	0.31
	2.00	100	1.00	1.00	0.65	0.20
		200	1.00	1.00	0.95	0.38
		400	1.00	1.00	1.00	0.65
	3.00	100	1.00	1.00	0.84	0.24
		200	1.00	1.00	0.99	0.53
		400	1.00	1.00	1.00	0.91

^a Standardized difference between the QTL genotypic means

ble 4). The correlation method needed a slightly larger N at $E(r) \leq 0.2$ and $d \geq 1$ and a smaller N at $E(r) > 0.2$ and $d < 1$ than the method of comparison of marker-genotype means for selfing populations. For sib-mating populations, the correlation method required a slightly larger N at $E(r) \leq 0.1$ and $d \geq 1$ and a smaller N at $E(r) > 0.1$ and $d < 1$ than the method of comparison of marker-genotype means. In other situations, the correlation method and the comparison of marker-genotype means required a similar N to detect linkage (Table 4). The correlation method needs smaller N to detect loose linkage than the method of comparison of marker-genotype means. According to Simpson (1989), genotypic effects of most QTLs on quantitative traits are less than 1 d. Thus, a sample consisting of hundreds of fully selfing or sib-mating individuals would be required to detect linkage. For example, at least 140 selfing individuals are needed by the correlation method for detection of linkage at $E(r)=0.1$ and $d=0.5$ with a power of only 0.70 and at the 0.05 significance level.

Estimation of recombination value

The r values estimated by the proposed method were affected by d except for $E(r)=0.5$. With an increase in d , estimates of r were closer to its expectation. The smaller $E(r)$, the greater the bias in estimates of r from the $E(r)$ (Table

Table 3 The power of the correlation method and the method of comparison of marker-genotype means at the 5% significance level using a sample size of 200 observations

Mating type	d ^a	Method	Expected recombination value					
			0.0	0.1	0.2	0.3	0.4	0.5
Selfing	0.50	t(c) ^b	0.96	0.71	0.43	0.21	0.13	0.92
		t(m) ^c	0.96	0.71	0.43	0.22	0.18	0.83
	1.00	t(c)	1.00	1.00	0.87	0.53	0.20	0.95
		t(m)	1.00	1.00	0.87	0.53	0.20	0.89
	2.00	t(c)	1.00	1.00	1.00	0.85	0.32	0.98
		t(m)	1.00	1.00	1.00	0.85	0.32	0.97
3.00	t(c)	1.00	1.00	1.00	0.99	0.36	1.00	
	t(m)	1.00	1.00	1.00	0.99	0.36	1.00	
Sib-mating	0.50	t(c)	0.98	0.52	0.30	0.15	0.08	0.95
		t(m)	0.98	0.52	0.31	0.16	0.08	0.93
	1.00	t(c)	1.00	0.98	0.60	0.22	0.08	0.97
		t(m)	1.00	0.98	0.60	0.22	0.08	0.95
	2.00	t(c)	1.00	1.00	0.95	0.38	0.07	0.99
		t(m)	1.00	1.00	0.95	0.38	0.07	0.99
3.00	t(c)	1.00	1.00	0.99	0.53	0.04	1.00	
	t(m)	1.00	1.00	0.99	0.53	0.04	1.00	

^a Standardized difference between the QTL genotypic means

^b t-statistic for the correlation method

^c t-statistic for the method of comparison of marker-genotype means (Simpson 1989)

Table 4 Sample size needed to detect linkage at the 5% significance level

Mating type	d ^a	Method	Expected recombination value		
			0.1	0.2	0.3
Selfing	0.50	t(c) ^c	138	338	819
		t(m) ^c	143	352	1041
	1.00	t(c)	44	102	269
		t(m)	39	101	303
	2.00	t(c)	21	50	109
		t(m)	13	38	119
3.00	t(c)	11	29	83	
	t(m)	9	26	85	
Sib-mating	0.50	t(c)	247	682	2115
		t(m)	257	874	3169
	1.00	t(c)	75	220	839
		t(m)	73	254	937
	2.00	t(c)	29	92	363
		t(m)	27	99	373
3.00	t(c)	20	67	268	
	t(m)	18	71	268	

^a Standardized difference between the QTL genotypic means

^b t-statistic for the correlation method

^c t-statistic for the method of comparison of marker-genotype means (Simpson 1989)

Table 5 Recombination values estimated by the correlation method from simulations of 200 sample observations

Mating type	d ^a	Expected recombination value					
		0.0	0.1	0.2	0.3	0.4	0.5
Selfing	0.50	0.30	0.37	0.41	0.44	0.47	0.49
	1.00	0.19	0.28	0.34	0.40	0.45	0.50
	2.00	0.06	0.18	0.26	0.35	0.43	0.50
	3.00	0.05	0.15	0.24	0.33	0.41	0.50
Sib-mating	0.50	0.31	0.39	0.43	0.46	0.48	0.48
	1.00	0.19	0.32	0.39	0.44	0.47	0.48
	2.00	0.09	0.24	0.33	0.41	0.46	0.49
	3.00	0.05	0.21	0.31	0.39	0.45	0.49

^a Standardized difference between the QTL genotypic means

Table 6 The frequency distribution of plant height (cm) at a RAPD marker locus after heterozygous F₁ plants of common vetch were allowed to self-pollinate for seven generations

Genotype	No. of observations	Coding of marker locus	Plant height (cm)
<i>MMQQ</i>	87	1	88
<i>MMqq</i>	15	1	65
<i>mmQQ</i>	20	-1	88
<i>mmqq</i>	93	-1	65

5). In general, r was overestimated, contrary to results from the maximum-likelihood method, which underestimated r in a segregating F₂ population (Weller 1986). The correlation method and the method of comparison of marker-genotype means did not differ in r estimates at any combination of $E(r)$ and d .

Estimation of recombination values by the proposed method under incomplete linkage was less efficient than

detection of linkage. Accuracy of estimates of r decreased as $E(r)$ decreased (Table 5) because the number of recombinants decreased. This characteristic was also illustrated by Knapp et al. (1990) and Weller (1986).

An example in Table 6 shows the frequency distribution of the four possible genotypes. The estimate of c between the marker and QTL was 0.68, significant at 0.01 with 213 degrees of freedom. Therefore, the two loci are

linked with an estimated r of 0.1. Note that the value of d at the QTL was less than 0.2.

The proposed method was based on only one QTL, or a tightly linked cluster of QTLs, controlling a quantitative trait where the effect of the QTLs was considered to be additive. If unlinked QTLs also control the trait of interest, the value of r will be overestimated because the amount of variation due to the QTL under consideration is mixed with variation due to other unlinked QTLs. Although the power to detect linkage also decreases, a balance can be achieved by increasing sample sizes. The maximum-likelihood method produced unreasonable results when a few QTLs controlled a trait (Weller 1987). The comparison of marker-genotype means encountered a similar problem (Simpson 1989).

The proposed method to detect linkage is simple and straightforward. This is an advantage over the maximum-likelihood method. A major advantage of the correlation method over the method of comparison of marker-genotype means is that it has a higher power to detect linkage at $E(r)=0.5$. Serious biases in estimates of r with tight linkage and $d \leq 1$, however, still have not been improved by this method.

In short, detection of linkage by the correlation method under incomplete linkage was more efficient than the estimation of recombination values. The correlation method and the method of comparison of marker-genotype means have similar powers when two loci are tightly linked, but the former has a higher power than the latter when there is no linkage. Therefore, false linkage conclusions would be excluded more often by the correlation method than by the method of comparison of marker-genotype means.

Appendix

Equations (1) and (2) in the text were derived directly from the definition of the simple correlation coefficient, i.e.,

$$E(c) = \frac{\sum f'xy - [(\sum f'x)(\sum f'y)]}{\sqrt{[(\sum f'x^2 - (\sum f'x)^2)] [\sum f'y^2 - (\sum f'y)^2]}} \quad (14)$$

where f' is the frequency of genotypes, x is the score of the marker locus, and y is the genotypic value at the QTL (Table 1).

Equations (10) and (11) were derived from (8) and (9), respectively, according to Taylor's expansion concerning the approximate variance of a function of random variables (Chiang 1980), as follows:

$$V(r) = \left(\frac{\partial r}{\partial c} \right)^2 S_c^2 \quad (15)$$

where S_c^2 is the sampling variance of c , which was given in (3), and

$$\frac{\partial r}{\partial c} = -\frac{1}{(1+c)^2} \quad (16)$$

Because $c = (1-2r)/(1+2r)$ for selfing populations, equation (10) was obtained by applying equations (3) and (16) to (15). Equation (11) was derived similarly.

References

- Chiang CL (1980) An introduction to stochastic processes and their applications. Krieger Publishing Company, Huntington, New York
- Haldane JBS, Waddington CH (1931) Inbreeding and linkage. *Genetics* 16:357-374
- Knapp SJ (1991) Using molecular markers to map multiple quantitative trait loci: models for backcross, recombinant-inbred, and doubled-haploid progeny. *Theor Appl Genet* 81:333-338
- Knapp SJ, Bridges WC, Birkes D (1990) Mapping quantitative trait loci using molecular-marker linkage maps. *Theor Appl Genet* 79:583-592
- Lander ES, Botstein D (1989) Mapping Mendelian factors underlying quantitative traits using RFLP linkage maps. *Genetics* 121:185-199
- SAS Institute (1989) SAS/STAT guide for personal computers. Version 6 edn. SAS Inst., Cary, North Carolina
- Simpson SP (1989) Detection of linkage between quantitative trait loci and restriction fragment length polymorphisms using inbred lines. *Theor Appl Genet* 77:815-819
- Snape JW (1988) The detection and estimation of linkage using doubled-haploid or single-seed-descent populations. *Theor Appl Genet* 76:125-128
- Soller M, Beckmann JS (1990) Marker-based mapping of quantitative trait loci using replicated progenies. *Theor Appl Genet* 80:205-208
- Weller JI (1986) Maximum-likelihood techniques for the mapping and analysis of quantitative trait loci with the aid of genetic markers. *Biometrics* 42:627-640
- Weller JI (1987) Mapping and analysis of quantitative trait loci in *Lycopersicon* (tomato) with the aid of genetic markers using approximate maximum-likelihood methods. *Heredity* 59:413-421
- Weller JI, Wyler A (1992) Power of different sampling strategies to detect quantitative trait loci variance effects. *Theor Appl Genet* 83:582-588
- Zhang XF, Mosjidis JA, Hu ZL (1992) Methods for detection and estimation of linkage between a marker locus and quantitative trait loci. *Plant Breed* 109:35-39