

A Quantitative Trait Locus Mixture Model That Avoids Spurious LOD Score Peaks

Bjarke Feenstra¹ and Ib M. Skovgaard

Department of Natural Sciences, Royal Veterinary and Agricultural University, DK-1871 Frederiksberg C, Denmark

Manuscript received December 5, 2003

Accepted for publication February 13, 2004

ABSTRACT

In standard interval mapping of quantitative trait loci (QTL), the QTL effect is described by a normal mixture model. At any given location in the genome, the evidence of a putative QTL is measured by the likelihood ratio of the mixture model compared to a single normal distribution (the LOD score). This approach can occasionally produce spurious LOD score peaks in regions of low genotype information (*e.g.*, widely spaced markers), especially if the phenotype distribution deviates markedly from a normal distribution. Such peaks are not indicative of a QTL effect; rather, they are caused by the fact that a mixture of normals always produces a better fit than a single normal distribution. In this study, a mixture model for QTL mapping that avoids the problems of such spurious LOD score peaks is presented.

FOR more than a decade, interval mapping (LANDER and BOTSTEIN 1989) has been the most commonly used method for quantitative trait locus (QTL) mapping in experimental crosses. Often, interval mapping is used to identify regions of interest in the genome, which are then analyzed with more refined methods such as composite interval mapping (ZENG 1993, 1994) or multiple interval mapping (KAO *et al.* 1999). In cases where interval mapping suggests the existence of a QTL in a region that is sparsely covered with markers, it may be decided to develop more markers in this region to map the putative QTL more accurately. There may, however, be situations where interval mapping produces strong evidence for a QTL, when in fact there is none. If, for instance, the residual environmental variation does not follow a normal distribution, interval mapping can result in spurious LOD score peaks in regions of low genotype information (*e.g.*, widely spaced markers or much missing marker data; BROMAN 2003).

In standard interval mapping the distribution of the phenotype is modeled as a mixture of two (or more) components corresponding to the two (or more) different genotypes at the putative QTL (LANDER and BOTSTEIN 1989). When a specific basic distribution like the normal is used for each component this approach has the side effect that even without genetic (marker) information the distribution is a mixture of two or more normals when a QTL is included in the model, while under the null hypothesis of no QTL there is only a single component. If the basic distribution is not normal the model including a QTL may fit to data much better

than the single-component model, even in a model without any genetic (marker) information and even if there is no real QTL.

As an example, consider the following preliminary data set from an ongoing study of yellow rust (*Puccinia striiformis*) resistance in wheat (*Triticum aestivum*): 55 doubled haploid lines (DHLs; see, for example, LYNCH and WALSH 1998) were scored for rust resistance using a 0–9 scale in which 0 is no rust and 9 is total infection. The phenotypes were taken to be the scores divided by 10 and arc sine square root transformed, a transformation often used for observations on a finite interval. The DHLs were genotyped for a suite of microsatellite markers and interval mapping was performed (Figure 1).

As can be seen from Figure 1, there were three large LOD score peaks that all occurred in regions where the markers were very far apart (80–100 cM). Also, it was noted that when all genotype information was disregarded and a mixture of two normal distributions was fitted to the phenotypic data, this resulted in a LOD score of 9.83 compared to a single normal distribution. The fact that the three LOD score peaks were of the same order of magnitude as the LOD score based on no genotype information and that the peaks occurred in regions of little genotype information strongly suggests that these peaks are artifacts.

In the wheat data set, visual inspection of the phenotype distribution (Figure 2) hints that in areas of little genotype information a mixture of two normal distributions could produce a better fit than a single normal distribution, and hence that spurious LOD score peaks could occur.

In other cases, however, it may be less clear whether a LOD score peak is an artifact or not. We present a new model (Equation 1) that is a mixture of two

¹Corresponding author: Department of Natural Sciences, Royal Veterinary and Agricultural University, Thorvaldsensvej 40, DK-1871 Frederiksberg C, Denmark. E-mail: bjarke@dina.kvl.dk

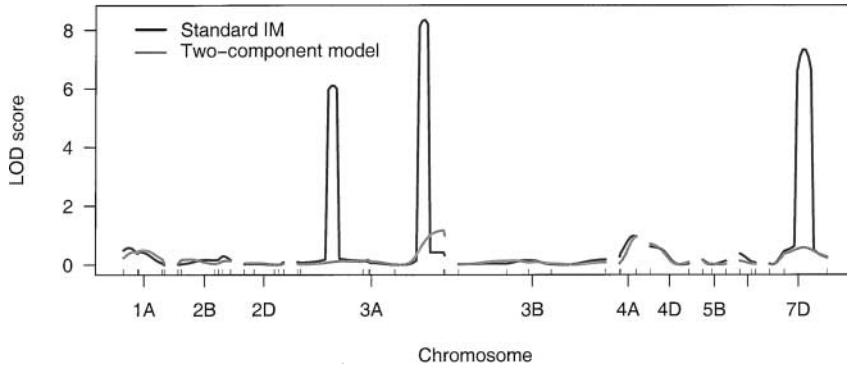


FIGURE 1.—Spurious LOD score peaks produced by standard interval mapping of data from 55 DHLs in wheat. The corresponding LOD score curve from the two-component mixture model (1) is included for comparison. The segments 1A, 2B, . . . , 7D represent different chromosomes.

components whether a QTL is present or not and therefore avoids the problems of such spurious LOD score peaks (see Figure 1).

METHODS

For simplicity, we consider a sample of n individuals from a backcross (BC) population (see, for example, LYNCH and WALSH 1998), but the results extend easily to other kinds of crosses. Let y_i and \mathbf{m}_i denote the quantitative phenotype and the multipoint marker data, respectively, for individual i .

To avoid the problem of spurious LOD score peaks we make sure that the model satisfies the following requirements: the distribution has the same number of components whether a QTL is present or not; without genetic information the model with and without a QTL is the same; and the model contains our original genetic model as a special case.

More concretely, the likelihood function of the parameter vector $\theta = (\mu_1, \mu_2, \sigma, \pi_1, \pi_2)$ is given by

$$L(\theta) = \prod_i \sum_j p_{ij} (\pi_j f(y_i; \mu_1, \sigma) + (1 - \pi_j) f(y_i; \mu_2, \sigma)), \tag{1}$$

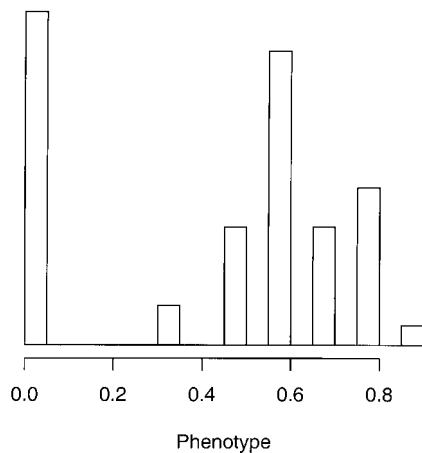


FIGURE 2.—Histogram of transformed disease resistance scores of 55 DHLs in wheat. Approximately 30% of the DHLs showed no sign of rust infection.

where $f(y; \mu, \sigma)$ is the density function for a normal distribution with mean μ and standard deviation σ . The index j may be thought of as the genotype at the putative QTL. The number p_{ij} is the conditional probability, given the marker data and the QTL position, that individual i has genotype j . The distribution of the phenotype of individuals with genotype j is now (unconventionally) modeled as a mixture of the two normal components with weights π_j and $1 - \pi_j$, respectively. Under normal assumptions we would like to see the estimates of these weights at a QTL position close to zero or one, indicating that a given genotype essentially results in a single-component normal distribution.

The likelihood function may be rewritten as

$$L(\theta) = \prod_i (c_i f(y_i; \mu_1, \sigma) + (1 - c_i) f(y_i; \mu_2, \sigma)), \tag{2}$$

where $c_i = \sum_j p_{ij} \pi_j$ is the weight of the first component in the two-component mixture distribution for individual i .

Now, the null hypothesis of no QTL effect is

$$H_0: \pi_j = \pi, \text{ for all } j,$$

implying that the distribution does not depend on the genotype of the putative QTL. The corresponding likelihood function is

$$L(\theta) = \prod_i (\pi f(y_i; \mu_1, \sigma) + (1 - \pi) f(y_i; \mu_2, \sigma)), \tag{3}$$

which, again, is a mixture of two normal distributions as required. In this case, however, the mixture coefficients do not depend on the QTL genotypes. Thus, the likelihood under H_0 is calculated just once.

Under the full model, we obtain maximum-likelihood estimates of the parameters with a form of the expectation-maximization (EM) algorithm (DEMPSTER *et al.* 1977). In the following, let z_i be an unobserved variable indicating whether the observation y_i comes from the first component ($z_i = 1$) or from the second component ($z_i = 2$) of the mixture. Let q_i be another unobserved variable indicating the true genotype at the putative QTL for individual i (*i.e.*, $q_i = 1$ or $q_i = 2$). Assume at iteration $s + 1$ we have estimates of the parameters $\hat{\theta}^{(s)}$.

In the E-step we must find $E(L(\hat{\boldsymbol{\theta}}^{(s)})|y_i)$, the conditional mean of the complete data log-likelihood function given the observed phenotypes. To do so, we calculate three different weights for each individual. First, for each of the two components in the mixture distribution,

$$\begin{aligned} w_{i1}^{(s+1)} &= \Pr(z_i = 1|y_i, \mathbf{m}_i, \hat{\boldsymbol{\theta}}^{(s)}) \\ &= \frac{\hat{c}_i^{(s)} f(y_i; \hat{\boldsymbol{\mu}}_1^{(s)}, \hat{\boldsymbol{\sigma}}^{(s)})}{\hat{c}_i^{(s)} f(y_i; \hat{\boldsymbol{\mu}}_1^{(s)}, \hat{\boldsymbol{\sigma}}^{(s)}) + (1 - \hat{c}_i^{(s)}) f(y_i; \hat{\boldsymbol{\mu}}_2^{(s)}, \hat{\boldsymbol{\sigma}}^{(s)})} \end{aligned}$$

and

$$w_{i2}^{(s+1)} = 1 - w_{i1}^{(s+1)}.$$

Second, for each of the two possible QTL genotypes,

$$\begin{aligned} u_{i1}^{(s+1)} &= \Pr(q_i = 1|y_i, \mathbf{m}_i, \hat{\boldsymbol{\theta}}^{(s)}) \\ &= \frac{p_{11} \hat{\pi}_1^{(s)} f(y_i; \hat{\boldsymbol{\mu}}_1^{(s)}, \hat{\boldsymbol{\sigma}}^{(s)}) + p_{12} (1 - \hat{\pi}_1^{(s)}) f(y_i; \hat{\boldsymbol{\mu}}_2^{(s)}, \hat{\boldsymbol{\sigma}}^{(s)})}{\hat{c}_i^{(s)} f(y_i; \hat{\boldsymbol{\mu}}_1^{(s)}, \hat{\boldsymbol{\sigma}}^{(s)}) + (1 - \hat{c}_i^{(s)}) f(y_i; \hat{\boldsymbol{\mu}}_2^{(s)}, \hat{\boldsymbol{\sigma}}^{(s)})} \end{aligned}$$

and

$$u_{i2}^{(s+1)} = 1 - u_{i1}^{(s+1)}.$$

Third, for the combination of mixture component and QTL genotype,

$$\begin{aligned} v_i^{(s+1)} &= \Pr(z_i q_i = 1|y_i, \mathbf{m}_i, \hat{\boldsymbol{\theta}}^{(s)}) \\ &= \frac{p_{11} \hat{\pi}_1^{(s)} f(y_i; \hat{\boldsymbol{\mu}}_1^{(s)}, \hat{\boldsymbol{\sigma}}^{(s)})}{\hat{c}_i^{(s)} f(y_i; \hat{\boldsymbol{\mu}}_1^{(s)}, \hat{\boldsymbol{\sigma}}^{(s)}) + (1 - \hat{c}_i^{(s)}) f(y_i; \hat{\boldsymbol{\mu}}_2^{(s)}, \hat{\boldsymbol{\sigma}}^{(s)})}. \end{aligned}$$

In the M-step, updated estimates of μ_1 , μ_2 , σ , π_1 , and π_2 are given by

$$\hat{\boldsymbol{\mu}}_l^{(s+1)} = \frac{\sum_i i w_{il}^{(s+1)} y_i}{\sum_i i w_{il}^{(s+1)}} \quad (4)$$

$$\hat{\boldsymbol{\sigma}}^{(s+1)} = \sqrt{\frac{1}{n} \sum_i \sum_l (y_i - \hat{\boldsymbol{\mu}}_l^{(s+1)})^2 w_{il}^{(s+1)}} \quad (5)$$

$$\hat{\pi}_1^{(s+1)} = \frac{\sum_i v_i^{(s+1)}}{\sum_i u_{i1}^{(s+1)}} \quad (6)$$

$$\hat{\pi}_2^{(s+1)} = \frac{\sum_i (w_{i1}^{(s+1)} - v_i^{(s+1)})}{\sum_i u_{i2}^{(s+1)}}, \quad (7)$$

where $l = 1, 2$ is the index of the mixture component. Initial values for the EM algorithm may, for example, be obtained by letting $\pi_j = 0.5$ and taking $w_{il}^{(0)} = \sum_j p_{ij} \pi_j$, and by letting $\hat{\boldsymbol{\mu}}_1^{(0)}$ and $\hat{\boldsymbol{\mu}}_2^{(0)}$ equal the estimates of μ_1 and μ_2 that are obtained under the null hypothesis. We iterate until the estimates converge.

Under the null hypothesis, we also use a form of the EM algorithm to obtain maximum-likelihood estimates of the parameters. As before, let z_i indicate which one of the two mixture components the observation y_i comes from. Under the null hypothesis, there is no QTL effect, so z_i is the only unobserved variable. In the E-step we calculate weights for each individual and for each of the two components in the mixture distribution,

$$\begin{aligned} w_{i1}^{(s+1)} &= \Pr(z_i = 1|y_i, \mathbf{m}_i, \hat{\boldsymbol{\theta}}^{(s)}) \\ &= \frac{\hat{\pi}^{(s)} f(y_i; \hat{\boldsymbol{\mu}}_1^{(s)}, \hat{\boldsymbol{\sigma}}^{(s)})}{\hat{\pi}^{(s)} f(y_i; \hat{\boldsymbol{\mu}}_1^{(s)}, \hat{\boldsymbol{\sigma}}^{(s)}) + (1 - \hat{\pi}^{(s)}) f(y_i; \hat{\boldsymbol{\mu}}_2^{(s)}, \hat{\boldsymbol{\sigma}}^{(s)})} \end{aligned} \quad (8)$$

and

$$w_{i2}^{(s+1)} = 1 - w_{i1}^{(s+1)}.$$

In the M-step, we obtain updated parameter estimates of μ_1 , μ_2 , and σ using Equations 4 and 5 and estimate π by the following equation:

$$\hat{\pi}^{(s+1)} = \frac{\sum_i i w_{i1}^{(s+1)}}{n}. \quad (9)$$

We initiate the EM algorithm by taking $w_{il}^{(0)} = 0.5$, which, however, causes $\hat{\boldsymbol{\mu}}_1^{(0)}$ and $\hat{\boldsymbol{\mu}}_2^{(0)}$ to be equal and $\hat{\pi}^{(0)}$ to be 0.5. In that case, as is seen from Equation 8, the weights and estimates are not changed by the iterations. This is a consequence of the symmetry of the model in the two components; in fact $\mu_1 = \mu_2 = \bar{y}$ is a stationary point on the likelihood surface. Thus, to prevent the algorithm from getting stuck, we offset the initial μ values slightly in opposite directions. We iterate until the estimates converge.

SIMULATIONS

To illustrate the properties of the two-component mixture model and to compare its performance with standard interval mapping, we performed a small simulation study. We assessed the occurrence of spurious LOD score peaks by simulating 80 BC individuals under a null model of no QTL. We simulated 12 chromosomes, each 120 cM long and each with four to nine randomly distributed markers. A random 10% of the marker genotype data was missing. Phenotypes were simulated from a threshold model; first a random number was drawn from a standard normal distribution and then it was rounded upward to the nearest of the following thresholds: 0, 0.5, 1, 1.5, 2, 2.5, 3, 3.5, 4, 4.5, and 5.0. The phenotype was taken to be the threshold value in question. A total of 2000 simulations were done and in each case the data were analyzed both with a standard interval mapping model and the two-component mixture model. For each simulated data set, the maximum LOD score and the length of the interval where it occurred were recorded. Figure 3 shows the maximum LOD score as a function of interval length. *A priori*, when there is no QTL, one would expect no dependence of maximum LOD score on interval length, but the figure suggests otherwise for standard interval mapping; 59 maximum LOD scores exceeded 4 and they almost exclusively occurred in intervals >40 cM. In contrast, the two-component model showed no such trend; only 7 LOD scores were >4 and there was no tendency of increasing LOD scores with increasing interval length (Figure 3). Also, in standard interval mapping the number of maximum

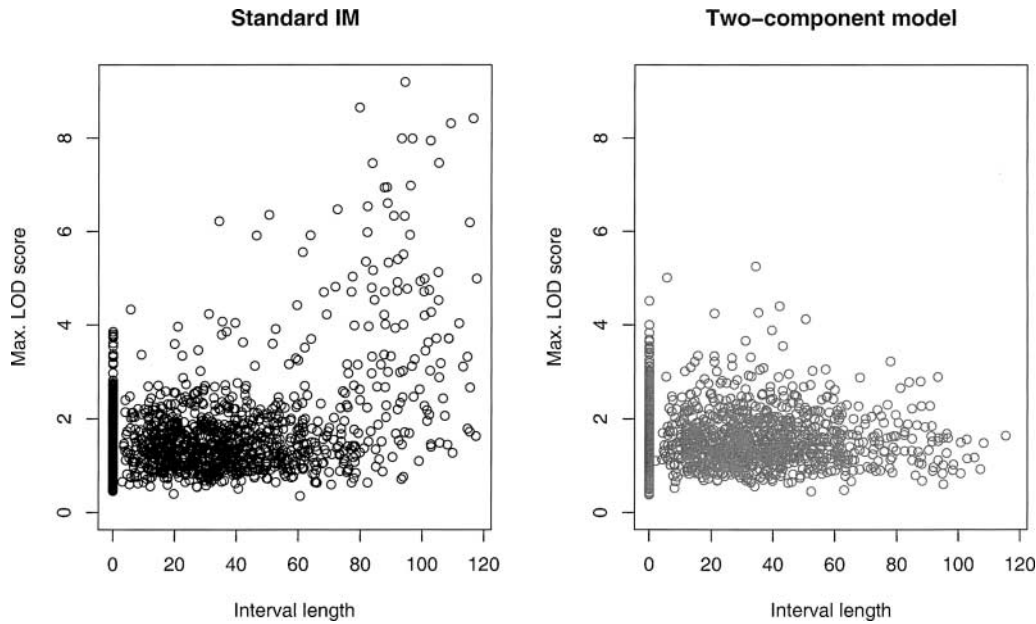


FIGURE 3.—Maximum LOD score as a function of the length of the interval where the maximum occurred for 2000 simulated data sets.

LOD scores in intervals >80 cM was twice that of the two-component model (108 *vs.* 50).

It might be expected that the price of extending the model, as we do in the two-component mixture model, is a loss of power. To compare power and precision of the two-component model with the standard interval mapping model, we simulated 200 BC individuals under a single-QTL model. We simulated five chromosomes, each 100 cM long and each with 11 randomly distributed markers and a QTL at position 60 cM on chromosome 1. We considered six different values of the additive effect of the QTL: 0 (null model), 0.12, 0.20, 0.26, 0.32, and 0.38. The trait value of an individual was determined by a random (environmental) variable drawn from a standard normal distribution plus the QTL effect (QTL genotype 2) *or* minus the QTL effect (QTL genotype 1). We performed 5000 simulations and analyzed all six QTL effects with both standard interval mapping and the two-component mixture model. We obtained genome-wide LOD thresholds from the data with no QTL effect, as the 95th percentiles of the maximum LOD score. The LOD thresholds for standard interval mapping and the two-component mixture model were 2.26 and 2.48, respectively. Figure 4 shows a simulation example. LOD scores were calculated and plotted at every 2 cM. It can be seen from the figure that in a data set not leading to spurious LOD score peaks, the evidence obtained by standard interval mapping and the two-component model may be very similar.

The power of the two methods was estimated as the proportion of the simulation replicates for which the maximum LOD score exceeded the corresponding LOD threshold. As can be seen in Figure 5A, the two-component mixture model had similar although slightly

lower power compared to that of standard interval mapping. To translate the power loss to the relative number of observations we used the approximate relationship

$$\beta(Q) \approx 1 - \Phi(z_{\alpha/2} - Q\sqrt{n}C) + \Phi(-z_{\alpha/2} - Q\sqrt{n}C),$$

where β is the power function, Q is the QTL effect, C is a constant, n is the number of individuals, Φ is the standard normal distribution function, and $z_{\alpha/2}$ denotes the upper $\alpha/2$ quantile of the standard normal distribution. On the basis of this relationship, the power loss of the two-component model corresponded to $\sim 12\%$ fewer observations in the standard interval mapping model. The approximation holds in general for two-sided tests of a parameter in a well-behaved statistical model (see VAN DER VAART 1998, Chap. 14), but in the present setting we use it only empirically without claiming any theoretical justification. We also estimated the precision in locating the QTL by means of the root-mean-square (RMS) error of the estimated QTL position (Figure 5B). The two methods had very similar precision of QTL localization, although interval mapping had a marginally greater precision (smaller RMS error) compared to that of the two-component model.

The additive QTL effect was estimated in somewhat different ways under the two models. Since in the simulations the QTL genotype indexed by $j = 2$ corresponded to a positive additive effect, the QTL effect under standard interval mapping was estimated as $\hat{a}_M = 0.5 \cdot (\hat{\mu}_2 - \hat{\mu}_1)$. In the case of the two-component model, the QTL effect was estimated as $\hat{a}_{2C} = 0.5 \cdot (\hat{\pi}_2\hat{\mu}_1 + (1 - \hat{\pi}_2)\hat{\mu}_2 - \hat{\pi}_1\hat{\mu}_1 - (1 - \hat{\pi}_1)\hat{\mu}_2)$. In each case, the QTL effect was estimated at the position of the maximum LOD score. True and estimated effect sizes are shown in Table 1; both models

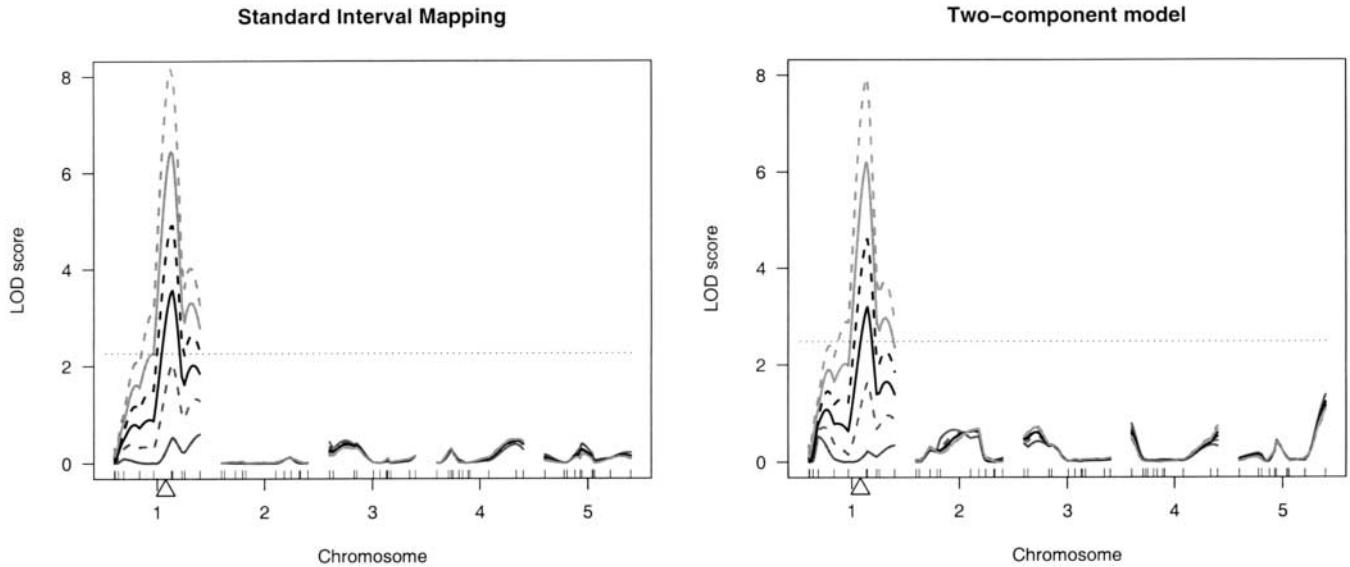


FIGURE 4.—A simulation example of QTL mapping on a population of 200 BC individuals. The QTL position is indicated by the triangle on the x -axis. The additive effect of the QTL increases in six steps from 0 (bottom solid curve) to 0.38 (top dashed curve). Standard interval mapping (left) and the two-component mixture model (right) applied to the same data set are compared. Genome-wide LOD thresholds are indicated by the dotted horizontal lines.

produced estimates slightly lower than the true values, since sometimes \hat{a}_{IM} and \hat{a}_{2C} were negative by chance. Note, however, that the estimates come very close to the true value as the QTL effect increases.

DISCUSSION

We have demonstrated that the commonly used standard interval mapping method may occasionally result in spurious LOD score peaks. In interval mapping the distribution is a mixture of two (or more) components when a QTL is included in the model while under the

null hypothesis of no QTL there is only a single component. Now, if the phenotype distribution is not normal, the two- (or more) component model may fit to data much better than the single-component model, even in a model without any genetic information and even if there is no real QTL. Thus, in cases where the phenotype distribution deviates from a normal distribution, false-positive results may be obtained in regions of low genotype information (*e.g.*, widely spaced markers, low degree of polymorphism, or much missing marker data). The problem was seen in an application (Figure 1). Close inspection of Figure 1 reveals that the LOD

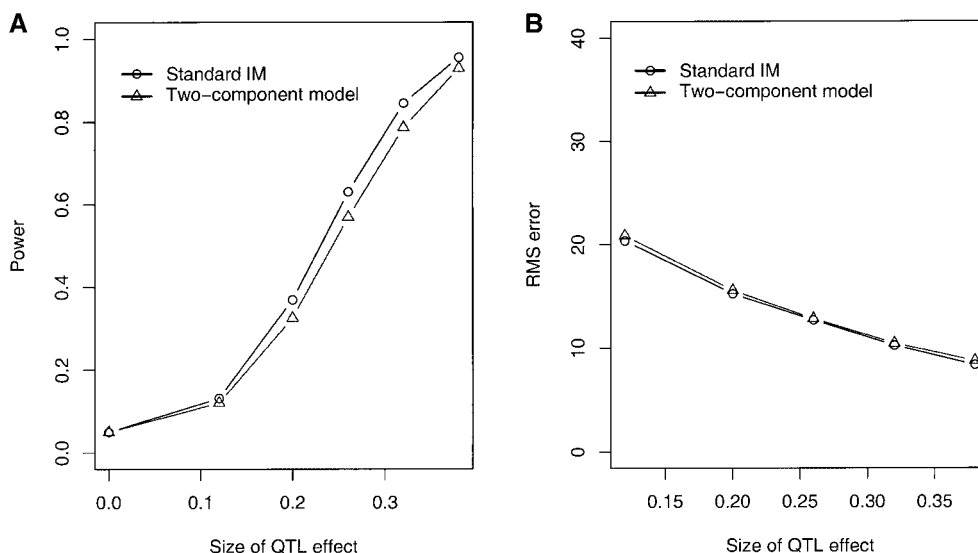


FIGURE 5.—(A) Estimated power to detect a QTL, based on 5000 simulation replicates. The standard error on the estimates ranged from 0.003 to 0.007. (B) Estimated root-mean-square (RMS) error of the estimated QTL location. Results are shown for standard interval mapping and for the two-component mixture model.

TABLE 1
True and estimated QTL effects from 5000 simulation replicates

	True effect					
	0	0.12	0.20	0.26	0.32	0.38
Estimate: interval mapping	-0.004	0.083	0.174	0.249	0.319	0.380
Estimate: two-component model	-0.002	0.076	0.161	0.234	0.306	0.371

Standard error of the means ranged from 0.0012 to 0.0025.

score curve jumps rather abruptly at the peaks. This is due partly to numerical difficulties in finding the global maximum of the likelihood function in the vicinity of the peaks. Thus, improved algorithms would widen and smoothen the peaks, but would not diminish their size.

We have presented a mixture model for QTL mapping that avoids this artifact. Our model is a mixture of two normal distributions (BC or DHL data) whether or not a QTL is included in the model; the QTL affects the mixing probabilities instead of the number of components. Our simulation results indicate that the two-component mixture model has only a minor loss of power and comparable precision to standard interval mapping in locating QTL over a range of QTL effects.

The results of analysis with the two-component mixture model must be interpreted with some care. In the case of a backcross population, we would like the absolute difference between π_1 and π_2 of Equation 1 to be close to 1 at a QTL position. This would indicate that the QTL genotypes from the parental lines each result in a single (different) normal distribution. In our simulations, increasing the additive QTL effect from 0.12 to 0.38 caused the mean of $|\hat{\pi}_1 - \hat{\pi}_2|$ at the true QTL location for data sets with $\text{LOD} > 2.48$ at that position to increase from 0.56 to 0.73 (data not shown). While these numbers are not that close to 1, it should be kept in mind that the residual variance used in the simulations was quite large at 1 compared to the additive QTL effects of 0.12–0.38. Also, it was noted that for a given QTL effect, the estimated difference between π_1 and π_2 increased with increasing LOD score. Still, it appears that the QTL effect needs to be larger compared to the residual variance for the mixing parameters π_1 and π_2 to be better estimated.

Several different numerical optimizations may be considered; the EM algorithm is often found to be somewhat slow but fairly robust and easy to program. As with other methods, there is no guarantee that it will find the global maximum rather than a local maximum, or even get stuck in a local minimum, but in our examples it seemed to work well, as judged from the LOD scores and other results obtained.

Other methods for QTL mapping have been developed for cases where the phenotype distribution is non-normal. If, for example, there is a spike in the phenotype

distribution (a large portion of the individuals share a common phenotype value) this may be modeled by a two-part parametric model (BROMAN 2003). However, the two-part model may also produce spurious LOD score peaks since one of its two parts is a mixture of two (or more) normal distributions when a QTL is included in the model, but only a single normal distribution under the null hypothesis. Thus, while the part corresponding to the common phenotype alleviates the problem, it may still occur if the remaining phenotype values deviate from a single normal distribution. One might also take a nonparametric approach to mapping QTL in the case of nonnormal phenotype distributions (KRUGLYAK and LANDER 1995; BROMAN 2003). Although generally a powerful alternative, nonparametric methods provide only a test for the presence of a QTL, whereas parametric methods also estimate the phenotypic effect of the QTL.

With the advent of extremely dense marker maps in a large number of species, it might be argued that researchers need not be concerned about getting spurious LOD score peaks from interval mapping. However, in many agriculturally important species only few markers have been developed, and even in species with many markers available, initial analyses may be undertaken with few markers to identify important regions of the genome. Moreover, the marker map may be dense and yet the genetic data may have poor information content, if, for example, the markers are dominant or if the proportion of missing data at certain marker loci is high. Also, it should be noted that the type of cross influences the risk of spurious LOD peaks from interval mapping. In the case of F_2 intercross populations (see, for example, LYNCH and WALSH 1998), the phenotype is modeled as a mixture of three components. In regions of low genotype information, the three-component mixture distribution produces a better fit than a two-component mixture distribution. Thus, spurious LOD peaks are expected to be more of a problem in F_2 intercrosses compared to, for instance, backcrosses or DHLs. For F_2 intercrosses, our two-component model may be extended to three components in a straightforward manner. However, problems with false or no convergence generally increase with the number of components in mixture models.

Finally, it should be stressed that spurious LOD peaks may arise for reasons other than the ones discussed here. For example, it is well known that analyzing a chromosome holding two linked QTL with a single-QTL model may result in a so-called “ghost” QTL (KNOTT and HALEY 1992; MARTÍNEZ and CURNOW 1992). While avoiding spurious LOD score peaks from low genotype information, the two-component model is not guaranteed to avoid all spurious LOD score peaks.

In conclusion, the mixture model presented here may be preferred in situations of large intermarker distances, dominant markers, low sample sizes, low degree of polymorphism in markers, or much missing marker information. Spurious LOD score peaks arising from standard interval mapping may be identified if they cannot be reproduced by the mixture model presented here.

The simulations and interval mapping analyses in this article were done with R/qtl (BROMAN *et al.* 2003), an add-on package for the general statistical software, R (IHAKA and GENTLEMAN 1996). The simulation study benefited from parallelization using the R-package SNOW (ROSSINI *et al.* 2003). The two-component model was implemented with a couple of new functions within the framework of R/qtl; the source code can be obtained by contacting the corresponding author.

We are grateful to two anonymous reviewers for their constructive comments, which led to improvements of the presentation and of the numerical optimization method. We thank Merethe J. Christiansen at Sejet Plant Breeding for kindly providing the DHL data.

LITERATURE CITED

- BROMAN, K. W., 2003 Mapping quantitative trait loci in the case of a spike in the phenotype distribution. *Genetics* **163**: 1169–1175.
- BROMAN, K. W., H. WU, S. SEN and G. A. CHURCHILL, 2003 R/qtl: QTL mapping in experimental crosses. *Bioinformatics* **19**: 889–890.
- DEMPSTER, A. P., N. M. LAIRD and D. B. RUBIN, 1977 Maximum likelihood from incomplete data via the EM algorithm. *J. R. Stat. Soc. B* **39**: 1–38.
- IHAKA, R., and R. GENTLEMAN, 1996 R: a language for data analysis and graphics. *J. Comp. Graph. Stat.* **5**: 299–314.
- KAO, C.-H., Z.-B. ZENG and R. D. TEASDALE, 1999 Multiple interval mapping for quantitative trait loci. *Genetics* **152**: 1203–1216.
- KNOTT, S. A., and C. S. HALEY, 1992 Aspects of maximum likelihood methods for the mapping of quantitative trait loci in line crosses. *Genet. Res.* **60**: 139–151.
- KRUGLYAK, L., and E. S. LANDER, 1995 A nonparametric approach for mapping quantitative trait loci. *Genetics* **139**: 1421–1428.
- LANDER, E. S., and D. BOTSTEIN, 1989 Mapping Mendelian factors underlying quantitative traits using RFLP linkage maps. *Genetics* **121**: 185–199.
- LYNCH, M., and B. WALSH, 1998 *Genetics and Analysis of Quantitative Traits*. Sinauer, Sunderland, MA.
- MARTÍNEZ, O., and R. N. CURNOW, 1992 Estimating the locations and the sizes of the effects of quantitative trait loci using flanking markers. *Theor. Appl. Genet.* **85**: 480–488.
- ROSSINI, A., L. TIERNEY and N. LI, 2003 Simple parallel statistical computing in R. University of Washington Biostatistics Working Paper Series 193, University of Washington, Seattle.
- VAN DER VAART, A. W., 1998 *Asymptotic Statistics*. Cambridge University Press, Cambridge, UK.
- ZENG, Z.-B., 1993 Theoretical basis of separation of multiple linked gene effects on mapping quantitative trait loci. *Proc. Natl. Acad. Sci. USA* **90**: 10972–10976.
- ZENG, Z.-B., 1994 Precision mapping of quantitative trait loci. *Genetics* **136**: 1457–1468.

Communicating editor: R. DOERGE

