

LOD significance thresholds for QTL analysis in experimental populations of diploid species

JOHAN W. VAN OOIJEN*

Centre for Biometry Wageningen (CBW), DLO-Centre for Plant Breeding and Reproduction Research (CPRO-DLO), PO Box 16, 6700 AA Wageningen, The Netherlands

Linkage analysis with molecular genetic markers is a very powerful tool in the biological research of quantitative traits. The lack of an easy way to know what areas of the genome can be designated as statistically significant for containing a gene affecting the quantitative trait of interest hampers the important prediction of the rate of false positives. In this paper four tables, obtained by large-scale simulations, are presented that can be used with a simple formula to get the false-positives rate for analyses of the standard types of experimental populations with diploid species with any size of genome. A new definition of the term 'suggestive linkage' is proposed that allows a more objective comparison of results across species.

Keywords: linkage, molecular marker, QTL, significance, statistics.

Introduction

Since the introduction of molecular genetic markers, linkage analysis has become one of the most important tools in biological research. A special type of linkage analysis can be carried out for quantitative traits. Quantitative traits, often called complex traits in contrast to simple Mendelian traits, are traits where the relation between genotype and phenotype cannot be observed directly. A gene affecting a quantitative trait is called a quantitative trait locus (QTL). The genetic dissection of a quantitative trait — so-called QTL analysis — is usually carried out using interval mapping (Lander & Botstein, 1989) or a related method (Haley & Knott, 1992; Jansen, 1992, 1993; Martinez & Curnow, 1992; Zeng, 1993, 1994). For this purpose an experimental population segregating for the quantitative trait is created and its linkage map of molecular markers is calculated. The basic procedure of the QTL analysis is such that on many positions on the linkage map a test statistic is calculated. In analogy with the genetic mapping of simple Mendelian traits, this statistic is the LOD score. This is essentially a likelihood ratio statistic. Subsequently, regions on the genome are identified that show significant values of the test statistic; such regions are supposed to contain a QTL. This procedure,

however simple, has a major problem: what value of the test statistic constitutes a significant value? A single LOD score is approximately related to a chi-squared distribution; the distribution of the maximum of a series of LOD scores, however, cannot be determined in a straightforward manner. Because of linkage, tests on neighbouring positions on the genome are not independent — closely linked markers will have equivalent test statistics. Also, the larger the genome, the more tests will be performed, thus increasing the probability that a fixed LOD threshold will be exceeded. Hence, if for the QTL analyses in various species an equal experiment-wise significance level is desired — usually 5% — the appropriate LOD thresholds will depend on the genome size of the species (in terms of recombination). Genome size varies greatly over species. Although most chromosome map lengths lie within the range of 50–250 cM, the numbers of chromosome pairs of diploid species start at two for *Haplopappus gracilis* (a plant), there are several species with just three pairs (e.g. *Crocus balansae* (a plant), *Crepis capillaris* (a plant), *Tipula maxima* (an insect)) and they go to beyond 50 pairs (John & Lewis, 1975; Dyer, 1979). Most agronomically interesting species have fewer than 25 pairs.

Several papers address the problem of statistical significance in QTL analysis and present solutions that are based on complex mathematical formulae, on cumulative distribution functions of the LOD score for specific situations obtained by simulation, or on

*Correspondence. E-mail: j.w.vanooijen@cpro.dlo.nl

permutation tests (Lander & Botstein, 1989; Van Ooijen, 1992; Feingold *et al.*, 1993; Churchill & Doerge, 1994; Rebaï *et al.*, 1994; Doerge & Churchill, 1996; Doerge & Rebaï, 1996; Dupuis & Siegmund, 1999). All solutions, however, have their own specific drawbacks. The mathematical formulae are seldom employed in the literature; apparently they are too complex, despite an available computer program (Rebaï *et al.*, 1994). The cumulative distribution functions of the LOD score are available only for a very limited number of genome sizes and population types. The permutation test requires very long computation times and is difficult to use in multiple-QTL models. The lack of an easy and relatively quick way of obtaining the appropriate significance threshold for the biological species under investigation results in QTL analyses that employ a certain LOD significance threshold, for which the significance level is not known, and which is possibly set at too low a value. As a consequence, the literature describing QTL analyses might contain false-positive QTLs at too high a rate. Lander & Kruglyak (1995) present in their paper on guidelines for QTL analysis in human, mouse and rat genetics, a number of precalculated threshold values that should be used in several types of analyses and experimental situations. This relieves geneticists from difficult or inconvenient computations. Because the LOD significance threshold depends on genome size, which varies greatly over species, it will be clear that the genetical research in all other biological species definitely is in need of an uncomplicated way to obtain the appropriate, correct LOD significance threshold. This paper presents the tails of the cumulative distribution functions of the LOD score in various situations as obtained by extensive simulations. With a simple method — a basic formula and four tables — the LOD threshold at the desired significance level for the standard experimental populations for most diploid species with their own genome sizes can be obtained easily. Besides F_2 and first-generation backcross (and equivalent types), the population types include a recombinant inbred (RI) family and a full-sib (FS) family of a cross between noninbred genotypes of an outbreeding species; for the last type no approximating formula has been published so far.

Methods

The tables present the results of stochastic (i.e. Monte Carlo) simulation of a diploid species with one chromosome on which no QTL was segregating. Various configurations were simulated. Of each configuration 1 000 000 repetitions were simulated. The chromosome map length was 50 up to 250 cM (Haldane mapping function) with a fully informative marker every 1 cM;

this is considered a reasonable approximation of a dense map. The population types were (a) a first-generation backcross (BC_1 , e.g. $F_1 \times P_1$) (b) an F_2 ($F_1 \times F_1$) (c) an RI family of the tenth generation (F_{10} , which is the ninth generation after the initial segregating meioses of the F_1) derived by single-seed descent from an F_2 , or (d) an FS family of a cross between noninbred genotypes of an outbreeding species (four alleles per marker). The population size was always 100. The individuals were assigned a random quantitative trait value according to a normal distribution regardless of the genotype, that is, there was no QTL. In each simulated population the LOD score, as defined by Lander & Botstein (1989), was calculated at all marker positions by fitting the appropriate model with two (BC_1), three (F_2), two (RI family) or four (FS family) QTL phenotypes. For the RI family the occasional heterozygous QTL phenotypes were fitted as strictly intermediate between the two homozygous QTL phenotypes. Subsequently, the maximum LOD on the chromosome was determined and recorded. From these data the cumulative distribution function of the maximum LOD score under the null hypothesis that no segregating QTL is present, was determined.

The LOD significance threshold

In experimentation one always wants to know the probability of arriving at the wrong conclusions. In QTL analysis these are (a) the conclusion that there is a segregating QTL whereas in reality there is not, or (b) not detecting a QTL which actually is present. The first type of error results in a false positive (type I), the second in a false negative (type II). The probability of false positives — the significance level — is controlled by choosing the appropriate significance threshold. The rate of false negatives is determined by the experimental set-up and the sizes of the genetic effects of the QTLs.

In a QTL study an experiment is set up to create a population that segregates for the quantitative trait. The values of the quantitative trait of the individuals in the population are recorded. The genotypes of segregating markers are determined and the linkage map is estimated. In the subsequent QTL analysis, tests for the presence of a segregating QTL are performed at many map positions on the genome — say every 1 cM. For these tests the LOD score is used. The areas on the genome are identified that show high values of the LOD score which are unlikely to occur if no QTL were segregating. It is concluded that statistically significant areas on the genome contain a segregating QTL. To know the significance level one needs to know the distribution of the test statistic under the null hypothesis (H_0) that no segregating QTL is present. Although under H_0 a single LOD score is approximately a chi-square random

variable (multiplied by a constant), the maximum of a series of LOD scores on a chromosome behaves according to a more complex type of distribution. Owing to the very nature of linkage these series of tests on a chromosome are not mutually independent. Because of the independent assortment of chromosomes

in meiosis, tests on different chromosomes are mutually independent.

In biology one usually desires an *experiment-wise* significance level of 5%. In the QTL analysis of a single segregating population this is equivalent to the *genome-wide* significance: the probability of obtaining a LOD

Table 1 Cumulative distribution function† of the maximum LOD on a chromosome for QTL analysis based on two QTL genotypes‡

LOD	Chromosome map length				
	50 cM	100 cM	150 cM	200 cM	250 cM
1.2	0.886342	0.810859	0.740251	0.677634	0.620000
1.3	0.908234	0.845325	0.785499	0.731602	0.680823
1.4	0.925729	0.873797	0.823399	0.777259	0.734005
1.5	0.940141	0.896915	0.855119	0.816276	0.779330
1.6	0.951725	0.915947	0.881135	0.848913	0.817736
1.7	0.961038	0.931704	0.902691	0.875919	0.849759
1.8	0.968463	0.944699	0.920616	0.898229	0.876651
1.9	0.974473	0.955127	0.935340	0.916764	0.898950
2.0	0.979307	0.963660	0.947465	0.931932	0.917565
2.1	0.983245	0.970420	0.957397	0.944578	0.932684
2.2	0.986510	0.976124	0.965499	0.954971	0.945223
2.3	0.989123	0.980653	0.972099	0.963468	0.955612
2.4	0.991252	0.984348	0.977345	0.970414	0.963999
2.5	0.992959	0.987388	0.981737	0.976064	0.970820
2.6	0.994363	0.989770	0.985230	0.980607	0.976327
2.7	0.995484	0.991737	0.988160	0.984404	0.980677
2.8	0.996334	0.993323	0.990414	0.987461	0.984357
2.9	0.997093	0.994571	0.992314	0.989849	0.987350
3.0	0.997717	0.995630	0.993764	0.991763	0.989701
3.1	0.998156	0.996512	0.994944	0.993389	0.991639
3.2	0.998513	0.997198	0.995956	0.994692	0.993218
3.3	0.998792	0.997733	0.996712	0.995749	0.994517
3.4	0.999032	0.998157	0.997387	0.996552	0.995607
3.5	0.999246	0.998525	0.997898	0.997189	0.996493
3.6	0.999393	0.998840	0.998311	0.997749	0.997209
3.7	0.999506	0.999050	0.998643	0.998202	0.997780
3.8	0.999604	0.999249	0.998882	0.998566	0.998248
3.9	0.999679	0.999390	0.999092	0.998852	0.998571
4.0	0.999746	0.999526	0.999258	0.999066	0.998844
4.1	0.999799	0.999618	0.999395	0.999246	0.999077
4.2	0.999852	0.999695	0.999509	0.999392	0.999260
4.3	0.999878	0.999759	0.999610	0.999520	0.999409
4.4	0.999899	0.999810	0.999691	0.999622	0.999514
4.5	0.999921	0.999860	0.999760	0.999703	0.999615
4.6	0.999937	0.999886	0.999804	0.999767	0.999693
4.7	0.999959	0.999905	0.999848	0.999812	0.999755
4.8	0.999965	0.999925	0.999879	0.999855	0.999802
4.9	0.999970	0.999934	0.999905	0.999887	0.999838
5.0	0.999980	0.999950	0.999927	0.999904	0.999866
5.1	0.999982	0.999963	0.999943	0.999924	0.999889
5.2	0.999984	0.999972	0.999955	0.999942	0.999911
5.3	0.999987	0.999977	0.999959	0.999950	0.999929
5.4	0.999989	0.999981	0.999964	0.999957	0.999941
5.5	0.999993	0.999986	0.999971	0.999964	0.999949

Table 1 (Continued)

LOD	Chromosome map length				
	50 cM	100 cM	150 cM	200 cM	250 cM
5.6	0.999994	0.999988	0.999974	0.999975	0.999958
5.7	0.999994	0.999993	0.999977	0.999977	0.999966
5.8	0.999996	0.999993	0.999984	0.999980	0.999969
5.9	0.999996	0.999994	0.999987	0.999985	0.999973
6.0	0.999996	0.999994	0.999988	0.999987	0.999981
6.1	0.999996	0.999996	0.999991	0.999989	0.999983
6.2	0.999997	0.999997	0.999991	0.999992	0.999986
6.3	0.999997	0.999997	0.999992	0.999993	0.999991

†The tail of the distribution with function values approximately from 0.9 to 0.99999.

‡Applicable to a first-generation backcross, a population of haploids or doubled haploids, and an F_2 for which in the analysis the heterozygous QTL genotype is fixed as strictly intermediate; all populations derived from a single heterozygous F_1 genotype as a parent.

above the threshold somewhere on the whole genome just by chance is 5%. A genome-wide threshold will depend on the number and length of the chromosomes, but also on the numbers of markers on the chromosomes. When just a few markers are tested per chromosome — the so-called sparse map case — a lower threshold is needed at the same genome-wide significance level than when many markers are tested per chromosome — the so-called dense map case (Lander & Botstein, 1989). Lander & Kruglyak (1995; and Kruglyak & Lander, 1995) strongly recommend the use of the dense map threshold, regardless of the actual density of the map used. One of the reasons is that geneticists will always deploy many additional markers in regions that show signs of a segregating QTL after an initial sparse map search. With modern marker techniques many markers will be used anyway.

To obtain a genome-wide significance the average map length of the chromosomes of the investigated species is used, because (a) usually the chromosome length does not vary much within a genome, (b) the genome-wide threshold is predominantly determined by the total genome length and nearly independent of the number of chromosome pairs in the genome (Kruglyak & Lander, 1995), and (c) it is difficult to think of any other easy solution. If chromosomes are assumed to be of equal length, then the property of independent chromosome assortment at meiosis can be used to obtain the relationship between the genome-wide and the corresponding chromosome-wide significance. By analogy the latter is the probability of obtaining a LOD above the threshold somewhere on a single chromosome just by chance. Suppose the required genome-wide significance level is α_g , the corresponding chromosome-wide significance level is α_c , the number of chromosome pairs is n and the

average chromosome length is l (in cM). Then the following relationship holds:

$$1 - \alpha_c = \sqrt[n]{(1 - \alpha_g)}.$$

The LOD threshold for the genome-wide significance level α_g can now be obtained from the cumulative distribution function (c.d.f.) of the maximum LOD under H_0 on a single chromosome of length l by looking up the LOD that has a c.d.f. value of $1 - \alpha_c$. Tables 1–4 present for several situations the tail of the c.d.f. of the maximum LOD score under H_0 on a single chromosome. The data were obtained by stochastic simulation. The situations comprise chromosome map lengths of 50 up to 250 cM (at multiples of 50 cM), which should suffice for most biological species. Further, the situations comprise experimental populations segregating for two, three and four QTL genotypes in the first meiotic generation and an RI family in the tenth generation; this should suffice for most experimental population types. The population types are discussed below.

The way to use these tables is as follows. Calculate $1 - \alpha_c$ with the above formula for the required α_g . Look up the LOD score at $1 - \alpha_c$ in the table for the c.d.f. of the maximum LOD for the appropriate population type. Usually the average map length l will not be a multiple of 50 cM. Therefore, look up the LOD under the two map lengths below and above l , and subsequently interpolate to obtain the required LOD threshold. For example, if we have an F_2 of a species with eight chromosome pairs and 120 cM average chromosome length, and we want a genome-wide false-positives rate of 5% ($\alpha_g = 0.05$), we obtain $1 - \alpha_c = 0.9936$. When we look up 0.9936 in Table 2, which applies to an F_2 , under

Table 2 Cumulative distribution function† of the maximum LOD on a chromosome for QTL analysis based on three QTL genotypes‡

LOD	Chromosome map length				
	50 cM	100 cM	150 cM	200 cM	250 cM
1.9	0.882329	0.797265	0.720667	0.650829	0.588102
2.0	0.901938	0.829155	0.762847	0.701098	0.644262
2.1	0.918482	0.856570	0.799526	0.745443	0.695338
2.2	0.932396	0.880204	0.831091	0.784116	0.740122
2.3	0.944162	0.900206	0.858222	0.817939	0.779605
2.4	0.953843	0.917013	0.881358	0.846986	0.813790
2.5	0.961776	0.930960	0.901063	0.872130	0.843510
2.6	0.968472	0.942992	0.917632	0.893086	0.869070
2.7	0.974115	0.952816	0.931646	0.911237	0.890815
2.8	0.978771	0.960964	0.943425	0.926185	0.909291
2.9	0.982569	0.967886	0.953326	0.938683	0.924753
3.0	0.985735	0.973588	0.961555	0.949201	0.937766
3.1	0.988289	0.978305	0.968167	0.958078	0.948615
3.2	0.990403	0.982314	0.973734	0.965476	0.957576
3.3	0.992171	0.985562	0.978308	0.971659	0.964960
3.4	0.993594	0.988167	0.982253	0.976677	0.971149
3.5	0.994737	0.990300	0.985484	0.980760	0.976326
3.6	0.995692	0.992096	0.988098	0.984230	0.980710
3.7	0.996470	0.993547	0.990254	0.987053	0.984155
3.8	0.997121	0.994749	0.992052	0.989435	0.987015
3.9	0.997672	0.995735	0.993472	0.991408	0.989388
4.0	0.998097	0.996539	0.994623	0.992992	0.991328
4.1	0.998451	0.997183	0.995602	0.994313	0.992946
4.2	0.998744	0.997678	0.996424	0.995364	0.994260
4.3	0.998979	0.998151	0.997147	0.996186	0.995328
4.4	0.999172	0.998504	0.997666	0.996930	0.996154
4.5	0.999338	0.998809	0.998068	0.997519	0.996891
4.6	0.999456	0.999025	0.998445	0.997951	0.997501
4.7	0.999558	0.999212	0.998739	0.998331	0.997974
4.8	0.999622	0.999375	0.998976	0.998644	0.998377
4.9	0.999700	0.999505	0.999148	0.998891	0.998689
5.0	0.999760	0.999599	0.999314	0.999108	0.998930
5.1	0.999806	0.999681	0.999434	0.999262	0.999153
5.2	0.999841	0.999742	0.999549	0.999415	0.999300
5.3	0.999869	0.999795	0.999619	0.999497	0.999410
5.4	0.999893	0.999836	0.999700	0.999600	0.999519
5.5	0.999913	0.999861	0.999749	0.999675	0.999622
5.6	0.999928	0.999885	0.999785	0.999731	0.999690
5.7	0.999942	0.999913	0.999828	0.999776	0.999739
5.8	0.999952	0.999924	0.999852	0.999819	0.999788
5.9	0.999959	0.999936	0.999884	0.999854	0.999831
6.0	0.999967	0.999953	0.999902	0.999877	0.999859
6.1	0.999972	0.999961	0.999920	0.999896	0.999889
6.2	0.999977	0.999965	0.999938	0.999914	0.999910
6.3	0.999984	0.999975	0.999949	0.999940	0.999930
6.4	0.999990	0.999977	0.999962	0.999953	0.999941
6.5	0.999992	0.999981	0.999970	0.999964	0.999950
6.6	0.999994	0.999984	0.999978	0.999971	0.999956
6.7	0.999996	0.999988	0.999981	0.999975	0.999963
6.8	0.999996	0.999992	0.999987	0.999981	0.999966
6.9	0.999996	0.999992	0.999989	0.999985	0.999970
7.0	0.999997	0.999994	0.999989	0.999986	0.999977

Table 2 (Continued)

LOD	Chromosome map length				
	50 cM	100 cM	150 cM	200 cM	250 cM
7.1	0.999997	0.999995	0.999991	0.999989	0.999981
7.2	0.999998	0.999995	0.999993	0.999992	0.999984
7.3	0.999998	0.999995	0.999995	0.999994	0.999988
7.4	0.999998	0.999996	0.999997	0.999994	0.999990

†The tail of the distribution with function values approximately from 0.9 to 0.99999.

‡Applicable to an F_2 , where in the analysis there are no restrictions on the heterozygous QTL genotype, i.e. any level of dominance is allowed.

100 cM and 150 cM, we find that the corresponding LODs are 3.7 and 3.9, respectively, so that by interpolating a LOD of 3.8 (rounded upwards in the safe direction) is obtained as the desired 5% genome-wide significance threshold. In (rare) cases where the average chromosome length is larger than 250 cM, use must be made of the already mentioned fact that the genome-wide threshold is predominantly determined by the total genome length and nearly independent of the number of chromosomes in the genome. For instance, the 5% LOD thresholds for a genome with a single 200 cM chromosome, for one with two 100 cM chromosomes, and for one with four 50 cM chromosomes lie within a 0.1 LOD range of each other.

Population types

The population types to which the four tables apply, differ with respect to the number of QTL phenotype classes freely fitted in the analysed model, with respect to there being either one or two initial heterozygous parental genotypes, and with respect to the generation number after the initial segregating meiosis/meioses. All populations of Tables 1–3 are derived from a single (or two identical) heterozygous F_1 genotype(s) as the parent(s) that generate(s) the segregation. The populations of Tables 1, 2 and 4 are the first generation, and that of Table 3 is the ninth generation after the initial segregating meioses. Table 1 is for segregation into two QTL genotypes, that is, just one of the parents causes segregation or there is only one parent. Thus, Table 1 applies to a first-generation backcross (BC_1), a population of haploids such as of some fungi, or a population of doubled haploids. Table 2 is for segregation into three QTL genotypes, which applies to an F_2 (i.e. $F_1 \times F_1$) where no restrictions are imposed on the heterozygous QTL phenotype in the analysis — any level of dominance is allowed. Table 3 is for an RI family in the tenth generation (F_{10}), where the QTL segregates predominantly into two homozygous QTL genotypes. Table 4 is for segregation into four QTL

genotypes, which applies to an FS family of a cross between noninbred genotypes of an outbreeding species. For Table 4, and until recently also for Table 3 (Dupuis & Siegmund, 1999), no approximating formulae have been published, although the corresponding situations are quite frequent in experimental set-ups; for instance RI families are often used in plant science (Burr & Burr, 1991), and using FS families is important in forest and fruit tree genetics (Grattapaglia & Sederoff, 1994; Hemmat *et al.*, 1994; Grattapaglia *et al.*, 1996; Maliepaard *et al.*, 1997, 1998).

For an F_2 a model is often fitted in which the heterozygous QTL phenotype is strictly intermediate (also called an additive model). In such a case in effect just two QTL phenotypes are fitted. The difference from the BC_1 , where also two QTL phenotypes are fitted, is that both parents, instead of one, have a segregating meiosis. This results in a slightly lower correlation between tests on linked markers in the F_2 , so that its c.d.f. of the maximum LOD on a chromosome is slightly different. This was verified by simulation. For a chromosome length of 50 cM and at c.d.f. values of 0.95, 0.99 and 0.999 the c.d.f. values for the BC_1 were approximately 0.003, 0.0003 and 0.00003, respectively, larger. For a chromosome length of 250 cM the differences were approximately 0.003, 0.0006 and 0.0001, respectively. In all these instances the differences were much smaller than the differences with c.d.f. values at 0.1 LOD smaller or larger. This means that for normal practice Table 1 can also be used for an F_2 where the heterozygous QTL phenotype is modelled as strictly intermediate.

RI families are employed in varying generations, but usually not before the F_5 . Because there is recombination from generation to generation, the correlation between tests on linked markers declines in later generations. Therefore, the 50 and 250 cM cases of the F_5 and the F_{20} were also simulated and compared to the F_{10} . For a chromosome length of 50 cM and at c.d.f. values of 0.95, 0.99 and 0.999 the c.d.f. values for the F_5 were approximately 0.004, 0.0006 and 0.0001,

respectively, larger. For a chromosome length of 250 cM the differences were approximately 0.003, 0.0005 and 0.0001, respectively. For the F_{20} the c.d.f. values were hardly different. In all these instances the differences were much smaller than the differences with c.d.f. values

at 0.1 LOD smaller or larger. Therefore, Table 3 can be used reliably for RI families of the usual generation numbers. The reason for these small differences is that most recombination occurs before the F_5 , whereas after the F_5 recombination has little effect because of fixation.

Table 3 Cumulative distribution function† of the maximum LOD on a chromosome for QTL analysis based on two QTL genotypes‡ in an RI family

LOD	Chromosome map length				
	50 cM	100 cM	150 cM	200 cM	250 cM
1.4	0.887975	0.805201	0.729732	0.661914	0.600372
1.5	0.909196	0.839955	0.775564	0.716586	0.662671
1.6	0.926162	0.868769	0.814346	0.763826	0.716884
1.7	0.940172	0.892565	0.847080	0.804054	0.763679
1.8	0.951659	0.912086	0.874381	0.838358	0.803901
1.9	0.960733	0.928499	0.897260	0.867343	0.838175
2.0	0.968371	0.941812	0.915997	0.891175	0.867001
2.1	0.974343	0.952903	0.931496	0.911077	0.890896
2.2	0.979187	0.961742	0.944206	0.927402	0.910979
2.3	0.983147	0.968965	0.954613	0.940935	0.927522
2.4	0.986479	0.974907	0.963054	0.951956	0.940853
2.5	0.989100	0.979658	0.970007	0.960851	0.951730
2.6	0.991162	0.983550	0.975677	0.968283	0.960743
2.7	0.992835	0.986829	0.980378	0.974244	0.968176
2.8	0.994175	0.989424	0.984107	0.979210	0.974319
2.9	0.995311	0.991476	0.987109	0.983193	0.979277
3.0	0.996204	0.993171	0.989572	0.986479	0.983280
3.1	0.996935	0.994457	0.991644	0.989048	0.986526
3.2	0.997545	0.995512	0.993305	0.991026	0.989149
3.3	0.998062	0.996360	0.994589	0.992818	0.991276
3.4	0.998456	0.997100	0.995624	0.994179	0.992889
3.5	0.998735	0.997699	0.996439	0.995265	0.994232
3.6	0.998968	0.998165	0.997119	0.996187	0.995334
3.7	0.999157	0.998505	0.997673	0.996912	0.996267
3.8	0.999335	0.998771	0.998138	0.997480	0.996978
3.9	0.999466	0.999032	0.998493	0.997975	0.997579
4.0	0.999559	0.999241	0.998766	0.998377	0.998021
4.1	0.999659	0.999382	0.999000	0.998678	0.998390
4.2	0.999727	0.999514	0.999200	0.998932	0.998715
4.3	0.999783	0.999620	0.999369	0.999157	0.998972
4.4	0.999830	0.999700	0.999499	0.999317	0.999174
4.5	0.999866	0.999767	0.999611	0.999453	0.999335
4.6	0.999894	0.999797	0.999694	0.999558	0.999481
4.7	0.999912	0.999822	0.999752	0.999640	0.999564
4.8	0.999927	0.999858	0.999798	0.999706	0.999635
4.9	0.999937	0.999879	0.999834	0.999763	0.999708
5.0	0.999953	0.999903	0.999857	0.999802	0.999764
5.1	0.999963	0.999920	0.999879	0.999852	0.999801
5.2	0.999968	0.999932	0.999902	0.999879	0.999850
5.3	0.999977	0.999951	0.999920	0.999902	0.999885
5.4	0.999981	0.999966	0.999939	0.999915	0.999909
5.5	0.999985	0.999973	0.999948	0.999930	0.999925
5.6	0.999985	0.999977	0.999960	0.999940	0.999935
5.7	0.999988	0.999982	0.999966	0.999952	0.999952

Table 3 (Continued)

LOD	Chromosome map length				
	50 cM	100 cM	150 cM	200 cM	250 cM
5.8	0.999991	0.999985	0.999974	0.999959	0.999958
5.9	0.999991	0.999990	0.999979	0.999965	0.999964
6.0	0.999992	0.999992	0.999982	0.999973	0.999975
6.1	0.999994	0.999994	0.999982	0.999976	0.999981
6.2	0.999995	0.999996	0.999987	0.999984	0.999984
6.3	0.999998	0.999997	0.999989	0.999986	0.999985
6.4	0.999998	0.999999	0.999990	0.999990	0.999987
6.5	0.999998	0.999999	0.999990	0.999991	0.999990

†The tail of the distribution with function values approximately from 0.9 to 0.99999.

‡Applicable to a recombinant inbred family of the fifth or higher generation.

For short distances the accumulated amount of recombination in late RI generations approaches twice the recombination frequency in a BC₁. Because also in an RI family two QTL classes are fitted, the c.d.f. of the maximum LOD on a chromosome of a certain length in an RI family is close to that of a chromosome of twice that length in a BC₁. Compare for instance the 50 cM column of Table 3 with the 100 cM column of Table 1.

Approximate multiple-QTL models

Multiple-QTL models are more powerful than single-QTL models when there are several segregating QTLs, but they require extreme computation times. As an alternative, Jansen (1992, 1993) and Zeng (1993, 1994), independently introduced approximate multiple-QTL models. Here, markers take over the role of the nearby QTLs and are fitted as cofactors while testing for a single QTL elsewhere in the genome. This way, the cofactors function as a genetic background control and absorb most of the genetic effects of their nearby QTLs from the residual variance. As a result, the power of the QTL analysis is enhanced, while reasonable computation times are retained.

In the mapping procedure with approximate multiple-QTL models (termed MQM mapping by Jansen, 1994), just as in interval mapping, tests for the presence of a single segregating QTL are performed at many positions in the genome. The difference between the two methods lies in the use of cofactor markers for background control of other segregating QTLs. The background control is part of both the null (no QTL) and the alternative (yes, a QTL) hypothesis. The tests in MQM mapping therefore have the same degrees of freedom as those in interval mapping. Simulation research of Jansen (1994) has shown that for MQM mapping the same LOD thresholds can be used as for interval mapping, under the condition 'that the residual degrees of freedom

for estimating the variance are adequate'. For testing under the presence of a linked QTL it is recommended that this linked QTL is flanked by two marker cofactors. Further, it is recommended that the number of parameters in the model is less than twice the square root of the number of individuals. In practice this means that this condition will be satisfied if there are not too many cofactors and the population is sufficiently large. The assignment of a marker cofactor essentially means that a QTL is concluded to be present. Experience so far has shown that the number of QTLs detected in a QTL analysis rarely exceeds 10, so that at least under current experimental practice the presented method of calculating LOD thresholds can also be applied to mapping with approximate multiple-QTL models.

Discussion

Using the presented method of calculating the LOD significance threshold will lead to a predictable rate of false-positive QTLs with reasonable accuracy. The values in the tables are accurate to about four decimal places (for more precise information about the accuracy, use can be made of the fact that each value in the tables is an estimate of a binomial probability). For the calculation of very high levels of significance use might be made of fitting some function through the tabulated data. As an alternative to a LOD threshold, the genome-wide significance level for the maximum LOD obtained in an analysis can be calculated with the method applied inversely. It must be realized that the calculated thresholds must be used as guidelines. Decisions with respect to further study, or utilization in breeding, of the particular genomic region should be based upon additional considerations, such as: What was the actual density of the markers used in the study? What generation was the RI family in? Does the trait behave according to normality? Does the estimated genetic effect of the QTL justify

Table 4 Cumulative distribution function† of the maximum LOD on a chromosome for QTL analysis based on four QTL genotypes‡

LOD	Chromosome map length				
	50 cM	100 cM	150 cM	200 cM	250 cM
2.4	0.893027	0.814419	0.743572	0.678096	0.619183
2.5	0.909560	0.841778	0.779605	0.721219	0.667719
2.6	0.923764	0.865540	0.810960	0.759701	0.712222
2.7	0.935659	0.885823	0.838892	0.794071	0.752156
2.8	0.946028	0.903357	0.862825	0.824132	0.787374
2.9	0.954848	0.918749	0.883863	0.850303	0.818580
3.0	0.962091	0.931499	0.901942	0.872949	0.845345
3.1	0.968357	0.942438	0.917332	0.892350	0.868863
3.2	0.973605	0.951750	0.930647	0.909149	0.889497
3.3	0.978022	0.959648	0.941813	0.923574	0.906987
3.4	0.981698	0.966257	0.951364	0.936037	0.921844
3.5	0.984802	0.971828	0.959340	0.946173	0.934490
3.6	0.987315	0.976535	0.966018	0.954972	0.945082
3.7	0.989433	0.980496	0.971783	0.962394	0.954153
3.8	0.991222	0.983759	0.976448	0.968699	0.961761
3.9	0.992726	0.986456	0.980513	0.974121	0.968106
4.0	0.993950	0.988738	0.983892	0.978466	0.973432
4.1	0.995015	0.990711	0.986652	0.982101	0.978033
4.2	0.995895	0.992391	0.988978	0.985089	0.981754
4.3	0.996690	0.993699	0.990894	0.987664	0.984927
4.4	0.997264	0.994794	0.992483	0.989880	0.987535
4.5	0.997774	0.995678	0.993753	0.991671	0.989704
4.6	0.998167	0.996432	0.994835	0.993113	0.991507
4.7	0.998468	0.997084	0.995784	0.994292	0.992961
4.8	0.998746	0.997579	0.996487	0.995314	0.994182
4.9	0.998959	0.998017	0.997101	0.996105	0.995195
5.0	0.999164	0.998380	0.997612	0.996821	0.996044
5.1	0.999306	0.998652	0.998020	0.997387	0.996782
5.2	0.999440	0.998920	0.998370	0.997867	0.997370
5.3	0.999538	0.999102	0.998659	0.998257	0.997811
5.4	0.999624	0.999279	0.998899	0.998577	0.998217
5.5	0.999699	0.999394	0.999093	0.998839	0.998540
5.6	0.999751	0.999502	0.999258	0.999046	0.998805
5.7	0.999798	0.999588	0.999396	0.999213	0.999002
5.8	0.999833	0.999665	0.999515	0.999365	0.999194
5.9	0.999870	0.999720	0.999606	0.999485	0.999328
6.0	0.999895	0.999768	0.999682	0.999571	0.999440
6.1	0.999912	0.999806	0.999741	0.999649	0.999545
6.2	0.999925	0.999843	0.999791	0.999718	0.999628
6.3	0.999935	0.999877	0.999820	0.999772	0.999695
6.4	0.999946	0.999899	0.999859	0.999826	0.999741
6.5	0.999956	0.999910	0.999879	0.999863	0.999790
6.6	0.999967	0.999936	0.999900	0.999891	0.999835
6.7	0.999972	0.999946	0.999918	0.999905	0.999857
6.8	0.999975	0.999958	0.999938	0.999926	0.999877
6.9	0.999983	0.999965	0.999948	0.999941	0.999896
7.0	0.999989	0.999967	0.999960	0.999951	0.999915
7.1	0.999990	0.999978	0.999968	0.999962	0.999933
7.2	0.999991	0.999984	0.999972	0.999968	0.999948
7.3	0.999992	0.999987	0.999979	0.999973	0.999953
7.4	0.999993	0.999991	0.999985	0.999980	0.999970
7.5	0.999993	0.999993	0.999988	0.999983	0.999973

Table 4 (Continued)

LOD	Chromosome map length				
	50 cM	100 cM	150 cM	200 cM	250 cM
7.6	0.999994	0.999997	0.999990	0.999986	0.999975
7.7	0.999995	0.999998	0.999992	0.999987	0.999979
7.8	0.999995	0.999999	0.999993	0.999987	0.999983
7.9	0.999997	0.999999	0.999994	0.999991	0.999986
8.0	0.999997	0.999999	0.999995	0.999994	0.999988
8.1	0.999998	0.999999	0.999996	0.999994	0.999989
8.2	0.999998	0.999999	0.999996	0.999994	0.999992

†The tail of the distribution with function values approximately from 0.9 to 0.99999.

‡Applicable to a full-sib family of a cross between noninbred genotypes of an outbreeding species.

further study? If several QTLs were detected, which ones explained most of the genetic variation?

The tables are based on simulations of a marker density of 1 cM. This is sufficiently representative for the dense map case in current experimental practice. Initial mapping populations usually consist of 100 up to 500 individuals. The effective deployment of higher marker densities requires much larger experimental populations. Such large populations must be used in the subsequent fine mapping. Lander & Kruglyak (1995) present LOD thresholds for QTL mapping in mouse and rat: for the backcross 3.3 LOD, for the F_2 intercross 4.3 LOD. According to the method presented here, these values are 3.1 and 4.0 LOD, respectively. The discrepancy is thought to be caused by differences in marker density: infinitely dense vs. 1 cM between markers; Lander & Kruglyak (1995) calculated a difference in LODs of about 7%, which corresponds to these findings.

Although the LOD score test appears to be reasonably robust against the data (after fitting the QTLs) having a skewed instead of a normal distribution (Doerge & Rebañ, 1996), other deviations from normality have not been investigated. Of course, the use of a permutation test avoids the problem of deviations from normality. An important drawback of the permutation test is that it will take several hours of computation time (on a 200-MHz PC) to obtain and analyse 1000 samples. This must be repeated for each trait. The question remains whether such sets of only 1000 samples would provide more accuracy than the use of the tables in this paper. In cases where the permutation test is going to be employed, the presented simulation results will be very useful for comparison.

Whether the calculated rate of false positives is acceptable depends on the general agreement on significance levels. In biology the usual rate is 5% for each experiment. In this respect, however, performing a QTL analysis is a peculiar kind of experiment. For a

segregating population, trait and marker data are determined. Each marker is tested for association with the trait. At a certain LOD threshold there exists a much larger opportunity for finding spurious linkage when many markers are tested because the investigated species has a large genome, than when few markers are tested. Now, what is considered an experiment in this QTL analysis? (a) The trait data plus one marker, (b) the trait data plus all markers on a single chromosome, or (c) the trait data plus all markers? This is important with respect to the experiment-wise 5% false-positives rate. There seems to be agreement that a whole genome scan (option (c)) should have a false-positives rate of 5%. However, this leads to the — at first sight strange — phenomenon that an Arabidopsis ($n=5$) geneticist may find a certain QTL effect that will be designated significant, whereas a wheat ($n=21$) geneticist detecting a similarly sized QTL effect cannot call it significant. Moreover, the wheat geneticist would have had to carry out a lot more work to obtain the results; that is, her/his experiment is much larger. On reflection it is clear that using a genome-wide 5% error rate should have such an effect: figuratively speaking, by allowing the collection of more wheat marker data the wheat geneticist simply gets many more shots at the bull's-eye.

Although Lander & Kruglyak (1995) use the same definition of genome-wide significance, their proposed classification of mapping results, suggestive and (highly) significant linkage, is based on the expected number of times that a LOD score above a certain threshold is obtained just by chance, in which multiple false positives per chromosome are allowed. From a statistical point of view it is an unusual approach; in statistics the definition of significance is based on a certain *probability* of obtaining false positives, rather than on a certain *expected number* of false positives. Although the result is not much different for a 5% probability of a false positive against an expected number of 0.05 false

positives, it is preferable to stick to the normal statistical approach of a significance level in the classification of mapping results, e.g. significant linkage should relate to a genome-wide 5% false-positives rate.

Lander & Kruglyak (1995) propose the term 'suggestive linkage' to allow for the publication of results that are not significant but point to a certain level of association between markers and trait. The use of a certain 'suggestive' level of significance is very appealing. The definition should be related to the fact that the analysis of the markers on a single chromosome can in a way be considered as a separate experiment. Because for each experiment a 5% error rate is an accepted rate, the term 'suggestive linkage' might be used for a chromosome-wide significance level of 5%. In recent years genetical research has discovered the potential power of comparative mapping (McKusick, 1997). For that purpose the results of mapping experiments must be comparable across species boundaries in an objective fashion. Because chromosome map length varies considerably across species, using a standard chromosome length of 100 cM in the definition of 'suggestive linkage' will allow an objective comparison of mapping results across species boundaries. Therefore, the proposal is to define the term 'suggestive linkage' for a chromosome-wide significance level of 5% for a standard chromosome length of 100 cM. For the various experimental population types that correspond to the four tables in this paper, the LOD thresholds for suggestive linkage are the fixed LOD values 1.9 (BC), 2.7 (F₂), 2.1 (RI) and 3.2 (FS), respectively.

The presented method provides reasonably accurate approximations to LOD significance thresholds. Mathematical formulae would have presented a more elegant solution, though these would probably be rather complex and are presently not available for some of the usual experimental situations. Genetical research in many species is expanding and is certainly in need of convenient ways to calculate significance thresholds applicable to the species under study with its own specific genome size. Therefore, the current results provide an equivalent, easy and pragmatic solution.

Acknowledgements

I thank Hans Jansen, Michiel Jansen, Chris Maliepaard and Piet Stam for critical comments and useful discussions.

References

BURR, B. AND BURR, F. A. 1991. Recombinant inbreds for molecular mapping in maize: theoretical and practical considerations. *Trends Genet.*, **7**, 55–60.

- CHURCHILL, G. A. AND DOERGE, R. W. 1994. Empirical threshold values for quantitative trait mapping. *Genetics*, **138**, 963–971.
- DOERGE, R. W. AND CHURCHILL, G. A. 1996. Permutation tests for multiple loci affecting a quantitative character. *Genetics*, **142**, 285–294.
- DOERGE, R. W. AND REBAI, A. 1996. Significance thresholds for QTL interval mapping tests. *Heredity*, **77**, 459–464.
- DUPUIS, J. AND SIEGMUND, D. 1999. Statistical methods for mapping quantitative trait loci from a dense set of markers. *Genetics*, **151**, 373–386.
- DYER, A. F. 1979. *Investigating Chromosomes*. Edward Arnold, London.
- FEINGOLD, E., BROWN, P. O. AND SIEGMUND, D. 1993. Gaussian models for genetic linkage analysis using complete high-resolution maps of identity by descent. *Am. J. Hum. Genet.*, **53**, 234–251.
- GRATTAPAGLIA, D. AND SEDEROFF, R. 1994. Genetic linkage maps of *Eucalyptus grandis* and *Eucalyptus urophylla* using a pseudo-testcross: mapping strategy and RAPD markers. *Genetics*, **137**, 1121–1137.
- GRATTAPAGLIA, D., BERTOLUCCI, F. L. G., PENCHEL, R. AND SEDEROFF, R. R. 1996. Genetic mapping of quantitative trait loci controlling growth and wood quality traits in *Eucalyptus grandis* using a maternal half-sib family and RAPD markers. *Genetics*, **144**, 1205–1214.
- HALEY, C. S. AND KNOTT, S. A. 1992. A simple regression method for mapping quantitative trait loci in line crosses using flanking markers. *Heredity*, **69**, 315–324.
- HEMMAT, M., WEEDEN, N. F., MANGANARIS, A. G. AND LAWSON, D. M. 1994. A molecular marker linkage map for apple. *J. Hered.*, **85**, 4–11.
- JANSEN, R. C. 1992. A general mixture model for mapping quantitative trait loci by using molecular markers. *Theor. Appl. Genet.*, **85**, 252–260.
- JANSEN, R. C. 1993. Interval mapping of multiple quantitative trait loci. *Genetics*, **135**, 205–211.
- JANSEN, R. C. 1994. Controlling the type I and type II errors in mapping quantitative trait loci. *Genetics*, **138**, 871–881.
- JOHN, B. AND LEWIS, K. R. 1975. *Chromosome Hierarchy*. Clarendon Press, Oxford.
- KRUGLYAK, L. AND LANDER, E. S. 1995. A nonparametric approach for mapping quantitative trait loci. *Genetics*, **139**, 1421–1428.
- LANDER, E. S. AND BOTSTEIN, D. 1989. Mapping Mendelian factors underlying quantitative traits using RFLP linkage maps. *Genetics*, **121**, 185–199.
- LANDER, E. AND KRUGLYAK, L. 1995. Genetic dissection of complex traits: guidelines for interpreting and reporting linkage results. *Nature Genet.*, **11**, 241–247.
- MALIEPAARD, C., JANSEN, J. AND VAN OOLJEN, J. W. 1997. Linkage analysis in a full-sib family of an outbreeding plant species: overview and consequences for applications. *Genet. Res.*, **70**, 237–250.
- MALIEPAARD, C., ALSTON, F. H., VAN ARKEL, G., BROWN, L. M., CHEVREAU, E., DUNEMANN, F. ET AL. 1998. Aligning male and female linkage maps of apple (*Malus pumila* Mill.) using multi-allelic markers. *Theor. Appl. Genet.*, **97**, 60–73.

- MARTINEZ, O. AND CURNOW, R. N. 1992. Estimating the locations and the sizes of the effects of quantitative trait loci using flanking markers. *Theor. Appl. Genet.*, **85**, 480–488.
- McKUSICK, V. A. 1997. Genomics: structural and functional studies of genomes. *Genomics*, **45**, 244–249.
- REBAÏ, A., GOFFINET, B. AND MANGIN, B. 1994. Approximate thresholds of interval mapping tests for QTL detection. *Genetics*, **138**, 235–240.
- VAN OOIJEN, J. W. 1992. Accuracy of mapping quantitative trait loci in autogamous species. *Theor. Appl. Genet.*, **84**, 803–811.
- ZENG, Z.-B. 1993. Theoretical basis for separation of multiple linked gene effects in mapping quantitative trait loci. *Proc. Natl. Acad. Sci. U.S.A.*, **90**, 10972–10976.
- ZENG, Z.-B. 1994. Precision mapping of quantitative trait loci. *Genetics*, **136**, 1457–1468.