

Sequential sampling in determining linkage between marker loci and quantitative trait loci

U. Motro¹ and M. Soller²

¹ Department of Evolution, Systematics and Ecology, and Department of Statistics

² Department of Genetics, The Hebrew University of Jerusalem, Israel

Received August 28, 1991; Accepted June 26, 1992
Communicated by A. L. Kahler

Summary. As compared to classical, fixed sample size techniques, simulation studies showed that a proposed sequential sampling procedure can provide a substantial decrease (up to 50%, in some cases) in the mean sample size required for the detection of linkage between marker loci and quantitative trait loci. Sequential sampling with truncation set at the required sample size for the non-sequential test, produced a modest further decrease in mean sample size, accompanied by a modest increase in error probabilities. Sequential sampling with observations taken in groups produced a noticeable increase in mean sample size, with a considerable decrease in error probabilities, as compared to straightforward sequential sampling. It is concluded that sequential sampling has a particularly useful application to experiments aimed at investigating the genetics of differences between lines or strains that differ in some single outstanding trait.

Key words: Sequential sampling – Quantitative trait loci – Genetic markers – QTL mapping – Linkage

Introduction

The uncovering of ever more prolific sources of highly polymorphic genetic markers, in both plant and animal populations (Beckmann and Soller 1983, 1990; Bernatzky and Tanksley 1986; Botstein et al. 1980; Chang et al. 1988; Fries et al. 1989; Georges et al. 1990; Helentjaris et al. 1986; Landry et al. 1987) makes it feasible to carry out experiments intended to map the

polygenic loci affecting quantitative traits (quantitative trait loci, QTL) in agricultural populations. In principle, mapping programs of this sort will involve two stages. The first stage will consist of a general testing of the entire genome in order to identify chromosomal regions affecting the trait (or traits) of interest. First-stage studies will often allow rough estimation of QTL position within these chromosomal regions (Weller 1987; Paterson et al. 1988; Lander and Botstein 1989). However, this is not the primary purpose of such experiments, and more precise estimates of QTL locations within the identified chromosomal regions, will generally require a second stage, involving additional experimental material specifically designed or collected for this purpose. In this study we consider the application of sequential sampling methods to the first, or testing, stage of QTL mapping programs. This will be particularly important in allowing an early decision to discontinue scoring markers in QTL-negative chromosomal regions. It was found that utilization of sequential sampling can indeed effect a major reduction in sample sizes required for detection of marker-QTL linkage, as compared to non-sequential methods. The effects of grouping of observations and of sample truncation were also examined.

Theory

In experiments designed to test for the presence of a QTL (or group of QTL) affecting a trait of interest in some chromosomal region, a population segregating for one or more genetic markers in that region is produced or obtained. A significant difference in quantitative trait value between marker genotypes is taken as an indication of the presence of one or more

QTL affecting trait value in the near vicinity of the marker (Sax 1923; Thoday 1961; Spickett and Thoday 1966; Kahler and Wehrhahn 1986; Edwards et al. 1987; Weller 1987; Paterson et al. 1988; Weller et al. 1988). A number of theoretical studies have been conducted to explore the power of such experiments in crosses between inbred lines (Soller et al. 1976) and segregating populations (Beckmann and Soller 1988), or in analyses carried out within segregating populations (Soller and Genizi 1978; Amos and Elston 1989; Weller et al. 1990). The effect of utilizing selected tails of populations (Lebowitz et al. 1987; Lander and Botstein 1989; Darvasi and Soller 1992), marker intervals (Lander and Botstein 1989), replicated progenies (Knapp and Bridges 1990; Soller and Beckmann 1990), and likelihood ratio tests as compared to ANOVA (Jensen 1989; Lander and Botstein 1989; Simpson 1989; Knapp et al. 1990) have also been explored. Numbers required for adequate power range from a few hundreds (crosses between inbred lines and segregating populations) to a few or many thousands (studies within segregating animal or human populations), depending on gene effect at the QTL relative to error variance and other factors.

In experiments involving large numbers of samples, marker information will be obtained in a sequential manner, as DNA samples from the experimental populations are scored for marker genotype. Thus, it would seem plausible to utilize statistical methods of sequential analysis to more rapidly obtain a decision to discontinue scoring the marker in QTL-negative chromosomal regions.

We consider here the simplest design for the detection of marker-QTL linkage, namely, the F_2 or backcross of a cross between inbred lines. In this design, if we denote alternative alleles at the marker locus, M and m , there are two classes of informative progeny: MM and mm in an experiment based on F_2 progeny, and Mm and mm in an experiment based on backcross progeny. In a typical simple random sampling experiment, a given total number of progeny, n , are produced. All progeny are evaluated with respect to quantitative traits of interest, and then scored one at a time or in small batches for the markers. When all progeny have been scored for both quantitative traits and markers, a marker-QTL linkage analysis is carried out.

Progeny showing one or other of the informative marker genotypes will be considered as belonging to one or other of two distinct populations, with a possible difference in the mean value of the quantitative trait of interest. More specifically, we consider two normally distributed random variables, x_1 and x_2 , which represent the value of the trait in each population, with unknown means μ_1 and μ_2 , but a known variance σ^2 . When two parental lines differ markedly with respect to QTL affecting a particular trait (e.g., in

crosses between resistant and sensitive populations, or between selection lines), the expected direction of effect associated with particular parental marker genotypes is known. In these cases, we test $H_0: \mu_1 = \mu_2$ against the one-sided alternative $H_1: \mu_1 > \mu_2$. In other cases, the direction of QTL effect in the two parental lines is not known, so we test the two-sided alternative $H_1: \mu_1 \neq \mu_2$. Practically, the acceptance of H_0 will be considered as an error if $(\mu_1 - \mu_2)/\sqrt{2}\sigma \geq \delta$ for a one-tailed test or if $|\mu_1 - \mu_2|/\sqrt{2}\sigma \geq \delta$ for a two-tailed test, where δ is a preassigned positive number. By appropriately choosing the units, we can set σ^2 equal to 1.

Let α and β denote the Type I and Type II error probabilities, respectively. Thus, if $\mu_1 - \mu_2 = 0$, the probability of rejecting H_0 does not exceed α , and whenever $\mu_1 - \mu_2 \geq \sqrt{2}\delta$ (for the one-tailed test) or $|\mu_1 - \mu_2| \geq \sqrt{2}\delta$ (for the two-tailed test), the probability of accepting H_0 does not exceed β .

Sequential sampling

Following Wald (1947), let $A = (1 - \alpha)/\beta$ and $B = \beta/(1 - \alpha)$. At the k^{th} ($k \geq 1$) stage, that is, after $2k$ observations (k from each population) have already been obtained, we compute a test statistic, such that if the test statistic is $\leq \log B$, we accept H_0 ; if it is $\geq \log A$, we reject H_0 ; while if the test statistic is between $\log B$ and $\log A$, we continue and take two more observations, one from each population.

Test statistics for the one-tailed and for the two-tailed tests can be derived from the likelihood ratio, in the following manner. Suppose we draw two samples (x_{11}, \dots, x_{1k} and x_{21}, \dots, x_{2k}), one from each of the two populations, each sample having the same size, k . The combined likelihood is

$$(1/2\pi)^k \exp \left[-\frac{1}{2} \sum_{i=1}^k (x_{1i} - \mu_1)^2 - \frac{1}{2} \sum_{j=1}^k (x_{2j} - \mu_2)^2 \right].$$

Denote $\bar{\mu} = (\mu_1 + \mu_2)/2$, then under H_0 the likelihood is

$$C \equiv (1/2\pi)^k \exp \left[-\frac{1}{2} \sum_{i=1}^k (x_{1i} - \bar{\mu})^2 - \frac{1}{2} \sum_{j=1}^k (x_{2j} - \bar{\mu})^2 \right].$$

If $\mu_1 - \mu_2 = \sqrt{2}\delta$, then $\mu_1 = \bar{\mu} + \delta/\sqrt{2}$ and $\mu_2 = \bar{\mu} - \delta/\sqrt{2}$, and the likelihood is

$$D \equiv (1/2\pi)^k \exp \left[-\frac{1}{2} \sum_{i=1}^k (x_{1i} - \bar{\mu} - \delta/\sqrt{2})^2 - \frac{1}{2} \sum_{j=1}^k (x_{2j} - \bar{\mu} + \delta/\sqrt{2})^2 \right],$$

whereas if $\mu_1 - \mu_2 = -\sqrt{2}\delta$, then $\mu_1 = \bar{\mu} - \delta/\sqrt{2}$ and

$\mu_2 = \bar{\mu} + \delta/\sqrt{2}$, and the likelihood is

$$E \equiv (1/2\pi)^k \exp \left[-\frac{1}{2} \sum_{i=1}^k (x_{1i} - \bar{\mu} + \delta/\sqrt{2})^2 - \frac{1}{2} \sum_{j=1}^k (x_{2j} - \bar{\mu} - \delta/\sqrt{2})^2 \right].$$

For the one-tailed test, the likelihood ratio is D/C , which, after some algebra, simplifies into

$$\exp [(S_{1,k} - S_{2,k})\delta/\sqrt{2} - \frac{1}{2}k\delta^2],$$

where $S_{1,k} = \sum_{i=1}^k x_{1i}$ and $S_{2,k} = \sum_{j=1}^k x_{2j}$ are the sum of observations of each sample. Hence the test statistic for the one-tailed test is

$$(S_{1,k} - S_{2,k})\delta/\sqrt{2} - \frac{1}{2}k\delta^2. \tag{1}$$

For the two-tailed test, two different approaches will be considered. The first is based on Wald's (1947) procedure for a single-sample, two-tailed test. Adapting Wald's single-sample procedure to our two-sample case, we thus consider the ratio $\frac{1}{2}(D + E)/C$, which simplifies into

$$\frac{1}{2} \exp(-\frac{1}{2}k\delta^2) \{ \exp[(S_{1,k} - S_{2,k})\delta/\sqrt{2}] + \exp[-(S_{1,k} - S_{2,k})\delta/\sqrt{2}] \}.$$

Hence the test statistic for the two-tailed test is

$$\log \{ \cosh[(S_{1,k} - S_{2,k})\delta/\sqrt{2}] \} - \frac{1}{2}k\delta^2, \tag{2}$$

where $\cosh x = \frac{1}{2}(e^x + e^{-x})$.

The second approach, or Armitage's (1975) procedure, actually considers

$$|S_{1,k} - S_{2,k}| \delta/\sqrt{2} - \frac{1}{2}k\delta^2 \tag{3}$$

as the test statistic for the two-tailed test.

Expected number of observations under sequential sampling

(1) One-tailed test. An approximation for the expected number of observations (taken from each population) required by the sequential procedure (when the difference of means $\mu_1 - \mu_2$ is θ) is (Wald 1947)

$$E_{\theta}(n) \approx \frac{P(\theta) \log A + [1 - P(\theta)] \log B}{(\theta - \delta/\sqrt{2})\delta/\sqrt{2}}, \tag{4}$$

where $P(\theta)$ is the probability that, at the final step, the test statistic takes a value $\geq \log A$, and $1 - P(\theta)$ is the probability that it takes a value $\leq \log B$. Hence, if $\mu_1 - \mu_2 = 0$,

$$E_0(n) \approx -[\alpha \log A + (1 - \alpha) \log B] / \frac{1}{2}\delta^2,$$

and if $\mu_1 - \mu_2 = \sqrt{2\delta}$,

$$E_{\sqrt{2\delta}}(n) \approx [(1 - \beta) \log A + \beta \log B] / \frac{1}{2}\delta^2.$$

(2) Two-tailed test. For the two-tailed sequential test there is no easy way to obtain an approximation for the expected sample size, hence for both Wald's and Armitage's procedures this parameter was estimated using computer simulations.

Effect of random presentation of populations

Testing equality of means by the procedure proposed in this study requires that the two samples, one from each population, will always be of equal size. In the present application it is not possible to determine to which population an individual belongs until after it is drawn. In such cases, if the test calls for, say, n individuals from each population, one should be prepared to end with examining more than $2n$ individuals. If each individual has the same probability of belonging to either population, the expected excess over $2n$ is (see Appendix)

$$2n \binom{2n}{n} \frac{1}{2^{2n}},$$

and for a large n it is, approximately, $2\sqrt{n}/\sqrt{\pi}$.

Effect of grouping

Suppose observations are taken in groups, each of size m , say. The sequential procedure is then as follows: At the k^{th} stage, that is, after sampling $2k$ groups of size m (k groups from each population), we compute the appropriate test statistic for the one-tailed [Eq. (1)] or the two-tailed [Eq. (2) or Eq. (3)] test, where k is replaced by km , and accept H_0 if the test statistic is $\leq \log B$; reject H_0 if it is $\geq \log A$; or continue and take two more groups of size m , one from each population, if the test statistic is between $\log B$ and $\log A$.

When such a procedure is followed, the expected sample size will be larger than when observations are taken singly. Indeed, the expected number (from each population) can be increased by an amount which is even larger than $m - 1$. This reflects the fact that after the first passage either into the acceptance region or into the rejection region, there is still a positive probability of re-entering the non-decision region (that is, for our test statistic to be back between $\log B$ and $\log A$). There is no easy way to assess the effect of grouping on the realized error probabilities α and β . Consequently, the effect of grouping was determined by simulation.

Effect of truncation

In some situations, sample progeny are taken from very large populations that are produced for a variety of commercial or experimental reasons, and there is no effective limit to the number of individuals that can

be sampled. This would be the case, for example, in analyses of dairy cattle data, where very large numbers of individuals are produced and evaluated with respect to quantitative traits, for routine reasons of herd management. In such cases, an unlimited sequential sampling procedure can be followed. In many cases, however, an experimental population will be produced particularly for purposes of marker-QTL analysis. In such cases, it is important to set some upper limit, say N , for the size of the population to be produced, and to terminate sequential sampling when this limit is reached. This is termed sequential sampling with 'truncation'. In this procedure we follow the usual sequential test procedure as long as $k \leq N - 1$. At the N^{th} stage, if and whenever it is reached, we either accept H_0 (if the test statistic is ≤ 0) or reject H_0 (if the test statistic is > 0).

Theoretically, truncation is expected to increase both Type I and Type II error probabilities. If N , however, is large enough, the effect of truncation on both types of error can be negligible. Again, there is no easy way to assess the effect of truncation on the realized α and β , and this was determined by simulation.

Numerical results

Tables 1 and 2 show estimates (based on computer simulations) of the mean sample size and of the Type I and Type II error probabilities for the proposed one-tailed and two-tailed sequential tests.

Table 1 considers the features of the sequential tests as a function of δ , the parameter which reflects

Table 1. Estimates of the mean size $[\hat{E}(n) \pm \text{SE}]$ of the sample drawn from each population, and of the realized error probabilities ($\hat{\alpha}$ and $\hat{\beta}$), for the sequential tests (without grouping or truncation), as a function of δ . The tests were constructed with the assumed probabilities 0.05 for each type of error, and the estimates were obtained by computer simulations, 2000 samples drawn for each case. (Since in the one-tailed test, the symmetric assumption $\alpha = \beta$ brings about mean sample sizes and realized error probabilities which are the same under $\theta = 0$ as under $\theta = \sqrt{2}\delta$, estimations for the one-tailed test were performed only under $\theta = 0$.) The required sample sizes for the classical, non-sequential test are given for comparison, together with the expected savings due to the sequential procedure

(a) One-tailed test

Resolution	$\delta = 0.2$	$\delta = 0.4$	$\delta = 0.6$	$\delta = 0.8$	$\delta = 1.0$
Sequential sampling					
$\hat{E}(n)$	136.1 \pm 2.1	37.3 \pm 0.5	17.1 \pm 0.2	10.3 \pm 0.2	7.4 \pm 0.1
Saving	50%	45%	43%	39%	33%
$\hat{\alpha}, \hat{\beta}$	0.0495	0.0355	0.0310	0.0265	0.0305
Fixed-size sampling					
n	271	68	30	17	11

(b) Two-tailed test

Resolution	$\delta = 0.2$	$\delta = 0.4$	$\delta = 0.6$	$\delta = 0.8$	$\delta = 1.0$
Sequential sampling, Wald's procedure					
$\theta = 0$					
$\hat{E}(n)$	209.8 \pm 2.0	55.5 \pm 0.5	26.6 \pm 0.3	14.2 \pm 0.1	9.5 \pm 0.1
Saving	35%	32%	26%	29%	27%
$\hat{\alpha}$	0.0455	0.0585	0.0500	0.0350	0.0235
$\theta = \sqrt{2}\delta$					
$\hat{E}(n)$	175.0 \pm 2.4	45.5 \pm 0.6	21.1 \pm 0.3	12.3 \pm 0.2	8.2 \pm 0.1
Saving	46%	44%	41%	39%	37%
$\hat{\beta}$	0.0420	0.0415	0.0300	0.0300	0.0345
Sequential sampling, Armitage's procedure					
$\theta = 0$					
$\hat{E}(n)$	247.8 \pm 2.7	61.6 \pm 0.7	30.9 \pm 0.4	16.8 \pm 0.2	10.9 \pm 0.1
Saving	24%	24%	14%	16%	16%
$\hat{\alpha}$	0.0555	0.0485	0.0635	0.0335	0.0245
$\theta = \sqrt{2}\delta$					
$\hat{E}(n)$	167.4 \pm 2.5	43.6 \pm 0.7	20.1 \pm 0.4	12.7 \pm 0.2	8.4 \pm 0.1
Saving	48%	46%	44%	37%	35%
$\hat{\beta}$	0.0475	0.0275	0.0345	0.0220	0.0140
Fixed-size sampling					
n	325	81	36	20	13

Table 2. The effect of grouping and of truncation. Estimates of the mean sample size [$\widehat{E}(n) \pm SE$] and of the realized error probabilities ($\hat{\alpha}$ and $\hat{\beta}$) are given for sequential tests with observations taken in groups (of size m), and for sequential tests subject to truncation (T). The tests were constructed for $\delta = 0.2$, with the assumed probabilities 0.05 for each type of error, and the estimates were obtained by computer simulations, 2000 samples drawn for each case. Truncation was at the required sample size for the classical, non-sequential test, i.e., at 271 for the one-tailed and at 325 for the two-tailed tests. Mean sample size is also shown as a proportion (R) of the mean size for $m = 1$ without truncation. (Note that values for $m = 1$ without truncation are the same values as for $\delta = 0.2$ in Table 1.) The proportion of truncated samples was 0.0895 for the one-tailed test, and 0.1070 and 0.0970 for the two-tailed test under H_0 and H_1 , respectively

(a) One-tailed test

Group Size (m)	1	1, T	10	50	100
$\widehat{E}(n)$	136.1 \pm 2.1	128.4 \pm 1.6	156.8 \pm 2.4	197.2 \pm 2.8	227.8 \pm 3.0
R	1.00	0.94	1.15	1.45	1.67
$\hat{\alpha}, \hat{\beta}$	0.0495	0.0575	0.0375	0.0210	0.0165

(b) Two-tailed test (Wald's procedure)

Group Size (m)	1	1, T	10	50	100
$\theta = 0$					
$\widehat{E}(n)$	209.8 \pm 2.0	203.6 \pm 1.4	223.2 \pm 2.1	268.4 \pm 2.6	296.1 \pm 2.9
R	1.00	0.97	1.06	1.28	1.41
$\hat{\alpha}$	0.0455	0.0515	0.0360	0.0225	0.0175
$\theta = \sqrt{2}\delta$					
$\widehat{E}(n)$	175.0 \pm 2.4	165.7 \pm 1.9	189.3 \pm 2.5	229.7 \pm 3.2	257.5 \pm 3.1
R	1.00	0.95	1.08	1.31	1.47
$\hat{\beta}$	0.0420	0.0615	0.0425	0.0280	0.0230

the required degree of resolution of the test. As expected, mean sample size increases with increase in the required degree of resolution (i.e., with decrease in δ), and sample sizes are greater for two-tailed than for one-tailed tests. Sequential sample sizes, nevertheless, are considerably smaller than those required by the classical, fixed sample size procedure. Yet, although the approximation for the mean sample size of the one-tailed test [Eq. (4)] predicts (under H_0 , for example) a saving of about 51% in sample size, the simulations show a somewhat smaller reduction, from 33% (if $\delta = 1.0$) to 50% (if $\delta = 0.2$). This, however, is accompanied by reduced error probabilities, which are generally less than those set in the simulation. A similar tendency is also displayed by the two-tailed tests.

It is not easy to decide which of the two-tailed tests is preferable. However, the simulations suggest that Wald's procedure has an advantage over Armitage's procedure under H_0 , whereas Armitage's procedure seems to be slightly more advantageous than Wald's under H_1 .

Table 2 shows the effect of grouping and the effect of truncation on sample sizes and error probabilities for the one-tailed and for the two-tailed (Wald's procedure) sequential tests. Truncation, which was set at the required sample size for the classical, non-sequential test, produced a modest decrease in mean sample size, accompanied by a modest increase in error probabilities. Thus, truncation at that level does not appear to produce a major decrease in overall accuracy of the sequential procedure.

As expected, grouping produced a noticeable enlargement of mean sample size, an enlargement which increases with group size, m . In each case, moreover, the proportional enlargement was almost twice as great for the one-tailed as for the two-tailed test. Notwithstanding, grouping produced a considerable decrease in error probabilities, to values well below those set in the simulation.

Discussion

The results of this study show that sequential sampling procedures can provide a substantial decrease in the mean sample sizes required for marker-QTL linkage determination.

Morton (1955) examined the mean sample size required by sequential tests for detection of linkage in human sibship data, as compared to various fixed size tests. He found that when linkage is present, the sequential test required only half as many observations as a fixed size test; while when linkage was not present, only one-third as many observations were required. These savings are greater than those obtained in the present study, but this is probably due to the weakness of the fixed size tests that he investigated.

Morton (1955) was also the first to develop user-friendly sequential probability ratio tests for the determination of linkage between marker loci and disease loci in human pedigrees. At present a considerable effort is being invested in developing user-friendly

likelihood ratio tests for determination of linkage between marker loci and QTL (Jensen 1989; Lander and Botstein 1989; Knapp et al. 1990). When these methods are available, they will be readily amenable to application of sequential sampling schemes, since all that will be required will be a comparison to the *A* and *B* factors of Wald (as defined above), of the likelihood ratio obtained at each stage. Although calculation of expected sample sizes for such test statistics will require extensive simulations, it is clear from the results of the present, as well as other, studies, that considerable savings can be confidently expected. Thus, it would appear that routine implementation of sequential sampling schemes with truncation and mild grouping can be recommended, even without further analysis. It should not be necessary to know the exact extent of savings, in order to benefit from whatever savings are to be obtained.

In practice, marker-QTL linkage analyses can be expected to involve many dozens or more of markers. Hence, the actual number of progeny that are raised and evaluated with respect to the quantitative traits of interest will not necessarily be reduced, since at least some of the markers can be expected to require considerably larger than average sample sizes for reaching a decision. The reduction in sample size, therefore, will be obtained as a reduction in the number of individuals that are scored, on the average, for any particular marker. That is, the reduction will be in the total number of individual marker determinations. When sample sizes are small and costs of marker scoring are only a small fraction of the overall costs of a study, as in the mapping of genes responsible for human diseases, it may be most convenient to score all progeny in parallel for most marker genotyping systems (e.g., RFLPs or automated PCR-based typing). However, agricultural or human studies involving polygenic loci with relatively small effects may require marker genotyping of many hundreds or thousands of individuals. In these cases, and in marker-assisted selection programs (Kashi et al. 1990), genotyping will be a major component of the overall cost of the study, and sequential sampling methods can be expected to yield substantial savings.

When many traits are considered, and all are evaluated simultaneously for marker-QTL linkage, it can also be expected that even for any particular marker, a decision for some traits will require a larger sample size than for others. In this case, the required number of marker identifications is determined by the trait which takes the longest to reach a decision, and this number can be distinctly larger than the number expected in sequential sampling for a single trait. Thus, in this case, very little can be gained by sequential procedures. Major savings, therefore, will be obtained in those instances where the experiment is aimed at

identifying QTL affecting a single trait. However, such experiments can be expected to be a major component in the total context of marker-QTL linkage analyses, including as they do all of those instances of marker-QTL analyses that are aimed at investigating the genetics of differences between lines or strains that differ in some single outstanding trait – as in the case of resistances to specific diseases [e.g., trypanotolerant N'Dama vs the sensitive Zebu (Soller and Beckmann 1987)], breeds characterized by an extreme phenotype for a particular trait (e.g., the high fertility Chinese swine, or broiler breeds of poultry as compared to layer breeds), or experimental selection lines that differ with respect to a single trait. In such cases, it has been shown that the use of selective genotyping (basing the analysis on selected extreme phenotypic tails of the population) can lead to a three- to four-fold reduction in the number of individuals assayed for the markers (Lander and Botstein 1989; Darvasi and Soller 1992). Thus, a combination of sequential sampling and selective genotyping can be expected to reduce the overall number of marker evaluations by a factor of almost eight-fold, as compared to the classical, fixed-size sampling from an unselected offspring population.

Acknowledgements. This research was supported by a grant from the US-Israel Binational Agricultural Research and Development Fund (BARD).

References

- Amos CI, Elston RC (1989) Robust methods for the detection of genetic linkage for quantitative data from pedigrees. *Genet Epidemiol* 6:349–360
- Armitage P (1975) *Sequential medical trials*. 2nd edn. Wiley, New-York
- Beckmann JS, Soller M (1983) Restriction fragment length polymorphisms in genetic improvement: methodologies, mapping and costs. *Theor Appl Genet* 67:35–43
- Beckmann JS, Soller M (1988) Detection of linkage between marker loci and loci affecting quantitative traits in crosses between segregating populations. *Theor Appl Genet* 76: 228–236
- Beckmann JS, Soller M (1990) Toward a unified approach to genetic mapping of eukaryotes based on sequence-tagged microsatellite sites. *Bio/Technology* 8:930–932
- Bernatzky T, Tanksley SD (1986) Towards a saturated linkage map of tomato based on isozymes and random cDNA sequences. *Genetics* 112:887–898
- Botstein D, White RL, Skolnick M, Davis RW (1980) Construction of a genetic linkage map in man using restriction fragment length polymorphisms. *Am J Hum Genet* 32: 314–331
- Chang C, Bowman JL, DeJohn AW, Lander ES, Meyerowitz EM (1988) Restriction fragment length polymorphisms map for *Arabidopsis thaliana*. *Proc Natl Acad Sci USA* 85: 6856–6860
- Darvasi A, Soller M (1992) Selective genotyping for determination of linkage between a marker locus and a quantitative trait locus. *Theor Appl Genet* 85:353–359

- Edwards MD, Stuber CW, Wendel JF (1987) Molecular-marker facilitated investigations of quantitative trait loci in maize. I. Numbers, genomic distribution and types of gene action. *Genetics* 116:113–125
- Fries R, Beckmann JS, Georges M, Soller M, Womack J (1989) The bovine gene map. *Anim Genet* 20:3–29
- Georges M, Mishra A, Sargeant L, Steele M, Zhao X (1990) Progress toward a primary DNA marker map in cattle. Proc 4th World Cong Genetics applied to Livestock Production, Edinburgh, 23–27 July, vol 13:107–112
- Helentjaris T, Slocum M, Wright S, Schaefer A, Nienhuis J (1986) Construction of genetic linkage maps in maize and tomato using restriction fragment length polymorphisms. *Theor Appl Genet* 72:761–769
- Jensen J (1989) Estimation of recombination parameters between a quantitative trait locus (QTL) and two marker gene loci. *Theor Appl Genet* 78:613–618
- Kahler AL, Wehrhahn CF (1986) Associations between quantitative traits and enzyme loci in the F_2 population of a maize hybrid. *Theor Appl Genet* 72:15–26
- Kashi Y, Hallerman E, Soller M (1990) Marker-assisted selection of candidate bulls for progeny testing programmes. *Anim Prod* 51:63–74
- Knapp SJ, Bridges WC (1990) Using molecular markers to estimate quantitative trait locus parameters: Power and genetic variances for unreplicated and replicated progeny. *Genetics* 126:769–777
- Knapp SJ, Bridges WC, Brikes D (1990) Quasi-Mendelian analyses of 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
- Landry BS, Kesseli RV, Farrara B, Michelmore RW (1987) A genetic map of lettuce (*Lactuca sativa*) with restriction fragment length polymorphisms, isozyme, disease resistance and morphological markers. *Genetics* 116:331–337
- Lebowitz RJ, Soller M, Beckmann JS (1987) Trait-based analyses for the detection of linkage between marker loci and quantitative trait loci in crosses between inbred lines. *Theor Appl Genet* 73:556–562
- Morton NE (1955) Sequential tests for the detection of linkage. *Am J Hum Genet* 7:277–318
- Paterson AH, Lander ES, Hewitt JD, Peterson S, Lincoln SE, Tanksley SD (1988) Resolution of quantitative traits into Mendelian factors by using a complete linkage map of restriction fragment length polymorphisms. *Nature* 335:721–726
- Sax K (1923) The association of size differences with seed-coat pattern and pigmentation in *Phaseolus vulgaris*. *Genetics* 8:552–560
- Simpson SP (1989) Detection of linkage between quantitative trait loci and restriction fragment length polymorphisms using inbred lines. *Theor Appl Genet* 77:815–819
- Soller M, Beckmann JS (1987) Toward an understanding of the genetic basis of trypanotolerance in the N'Dama cattle of West Africa. Consultation report submitted to FAO, Rome, March 1987
- Soller M, Beckmann JS (1990) Marker-based mapping of quantitative trait loci using replicated progenies. *Theor Appl Genet* 80:205–208
- Soller M, Genizi A (1978) The efficiency of experimental designs for the detection of linkage between a marker locus and a locus affecting a quantitative trait in segregating populations. *Biometrics* 34:47–55
- Soller M, Brody T, Genizi A (1976) On the power of experimental design for the detection of linkage between marker loci and quantitative loci in crosses between inbred lines. *Theor Appl Genet* 47:35–39
- Spickett SG, Thoday JM (1966) Regular responses to selection. 3. Interactions between located polygenes. *Genet Res* 7:96–121
- Thoday JM (1961) Location of polygenes. *Nature* 191:368–370
- Wald A (1947) *Sequential analysis*. Wiley, New-York
- 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, Soller M, Brody T (1988) Linkage analysis of quantitative traits in an interspecific cross of tomato (*Lycopersicon esculentum* × *Lycopersicon pimpnellifolium*) by means of genetic markers. *Genetics* 118:329–339
- Weller JI, Kashi Y, Soller M (1990) Power of 'daughter' and 'granddaughter' designs for mapping of quantitative traits in dairy cattle using genetic markers. *J Dairy Sci* 73:2525–2532

Appendix

Effect of random drawing of sample populations on total sample size

Suppose that we already have a certain number of pairs. We can either be in a balanced situation, where we have the same number of individuals from each population, or we can be in an unbalanced situation. If we wish to obtain an additional pair, and if we are in an unbalanced situation, we need to get exactly one more individual from whichever population is in deficiency. The number of additional draws then has a geometric distribution with the parameter $1/2$, and its expectation is 2. If we are in a balanced situation, after a single draw we will shift into an unbalanced situation. Thus, the expected number of draws needed in this case to obtain an additional pair is $1 + 2 = 3$.

When obtaining the k^{th} pair, we are in a balanced situation with probability

$$p_k = 2 \binom{2k-1}{k-1} \frac{1}{2^{2k}} = \binom{2k}{k} \frac{1}{2^{2k}},$$

and in an unbalanced situation with probability $1 - p_k$. The expected number of additional draws needed to obtain the $(k+1)^{\text{th}}$ pair is $3p_k + 2(1 - p_k) = 2 + p_k$. Hence, the expected number of draws needed to obtain n pairs is

$$\sum_{k=0}^{n-1} (2 + p_k) = 2n + \sum_{k=0}^{n-1} p_k,$$

which is larger than $2n$ by an excess of $\sum_{k=0}^{n-1} p_k$. Note that for $k \geq 1$, $p_k = p_{k-1}(1 - 1/2k)$, that is,

$$2kp_k = 2kp_{k-1} - p_{k-1}.$$

Summing both sides of this equation over k from 1 to n yields

$$\sum_{k=0}^{n-1} p_k = 2np_n.$$

Hence, for obtaining n pairs, the expected excess in the number of draws is $2n \binom{2n}{n} \frac{1}{2^{2n}}$.

Using Stirling's formula, we see that for a large n , the expected excess is approximately $2\sqrt{n}/\sqrt{\pi}$.